# Reinforcement Learning

—

# Exercise 4: Model-free Prediction

Nico Meyer

Fraunhofer
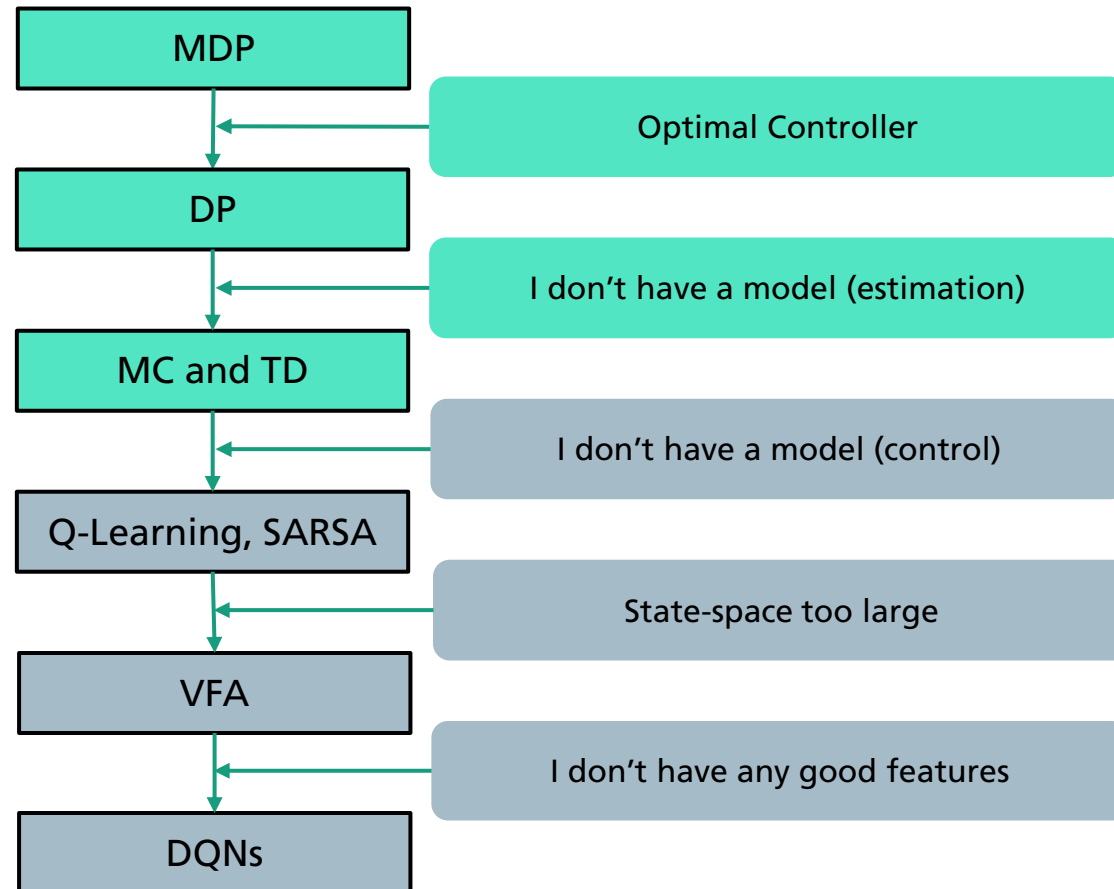
IIS

Fraunhofer-Institut für Integrierte
Schaltungen IIS

# Overview
## Exercise Content

| Week | Date | Topic | Material | Who? |
|------|------|-------|----------|------|
| 0 | | | *no exercises* | |
| 1 | 23.04. | MDPs | | Nico |
| 2 | 30.04. | Dynamic Programming | | Alex |
| 3 | 07.05. | OpenAI Gym, PyTorch-Intro | | Alex |
| 4 | 14.05. | TD-Learning | | Nico |
| 5 | 22.05. | Practical Session (zoom@home) | **Attention: Lecture Slot!** | Nico + Alex |
| 6 | 28.05. | TD-Control | | Nico |
| 7 | 04.06. | DQN | | Nico |
| 8 | 11.06. | VPG | | Alex |
| 9 | 18.06. | A2C | | Nico |
| 10 | 25.06. | Multi-armed Bandits | | Alex |
| 11 | 02.07. | RND/ICM | | Alex |
| 12 | 09.07. | MCTS | | Alex |
| 13 | 16.07. | BCQ | | Nico |

Fraunhofer

IIS

# Overview
## Overall Picture

# Recap
## Model-free Prediction

Fraunhofer
IIS

# Recap
## Monte Carlo and TD Methods

- So far: We know our MDP model $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$.
  - Planning by using dynamic programming
  - Solve a known MDP

- What if we don't know the model, i.e., $\mathcal{P}$ or $\mathcal{R}$ or both?

- We distinguish between 2 problems for unknown MDPs:

  - **Model-free Prediction:** Evaluate the future, given the policy $\pi$.
    *(estimate the value function)*

  - **Model-free Control:** Optimize the future by finding the best policy $\pi$.
    *(optimize the value function)*

Fraunhofer
IIS

# Recap
## Monte Carlo Policy Evaluation

- MC Policy Evaluation
  - MC methods learn from episodes of experience under policy $\pi$:

  $$s_t, a_t, r_t, s_{t+1}, \ldots, s_{T-1}, a_{T-1}, r_{T-1}, s_T \sim \pi$$

  - To evaluate a state $s \in \mathcal{S}$ we keep track of the rewards received from that state onwards.

- First-Visit Monte-Carlo Policy Evaluation:
  - First time-step $t$ that state $s$ is visited in an episode
    - Increment counter $N(s) \leftarrow N(s) + 1$,
    - Increment total return $S(s) \leftarrow S(s) + G_t$,
    - Value is estimated by mean return: $V(s) = S(s)/N(s)$
  - Our estimation $V(s)$ will come close to $V^{\pi}(s)$ as $N(s) \rightarrow \infty$.
    (considering the law of large numbers)

# Recap
## Temporal Difference Policy Evaluation

- Temporal-Difference Learning
  - Breaks up episodes and makes use of the intermediate returns
  - Learns from incomplete episodes (bootstrapping)
  - **We update a guess towards a guess**

$$V^\pi(s) = \underbrace{r(s,\pi(s))}_{} + \gamma \sum_{s' \in S} \underbrace{\boxed{\mathcal{P}(s'|s,\pi(s))} V^\pi(s')}_{}$$

We don't know the transition model

$$\boxed{(s, a, r, s')}$$

But we have real transitions available

$$\boxed{V^\pi(s) = r + \gamma V^\pi(s')}$$

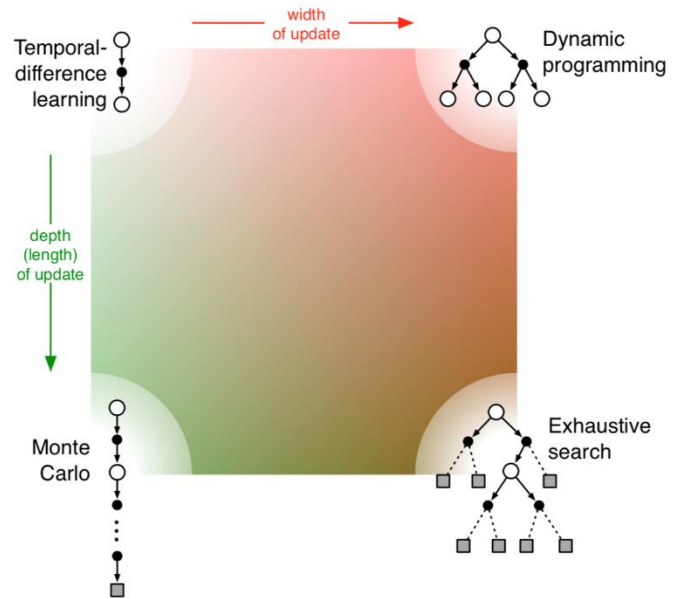Let's assume that the reality is the transition we observed

$$\boxed{V(s) \leftarrow V(s) + \alpha(r + \gamma V(s') - V(s))}$$

→ and update our old estimate "a bit" in this direction
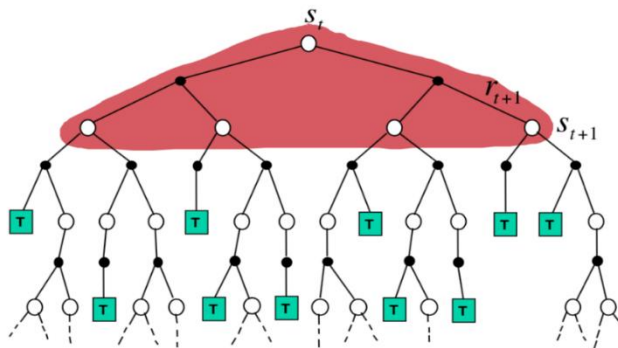
Fraunhofer
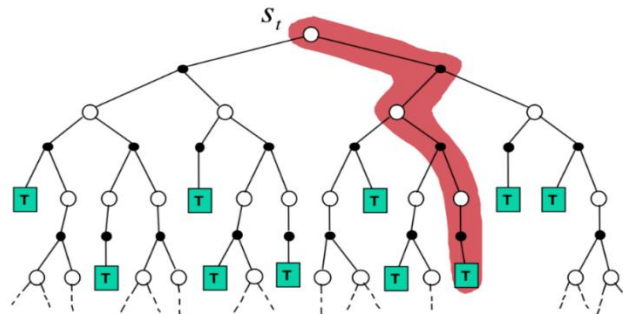IIS

# Recap
## DP vs. MC vs. TD



### DP Backup

$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$
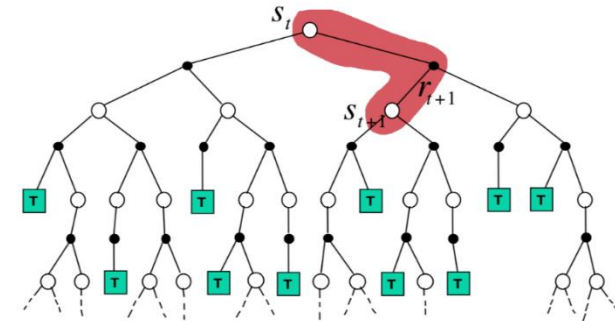


### MC Backup

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



### TD Backup

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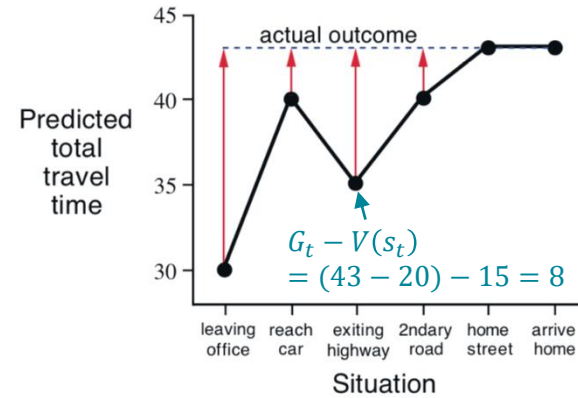


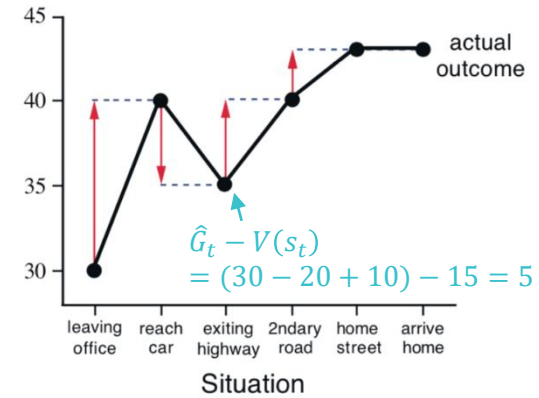*Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.*

# Recap
## TD and MC Algorithms

MC ($\alpha = 1$)



$$G_t - V(s_t)$$
$$= (43 - 20) - 15 = 8$$

TD ($\alpha = 1$)

$$\hat{G}_t - V(s_t)$$
$$= (30 - 20 + 10) - 15 = 5$$

*Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.*

---

**Tabular TD(0) for estimating $v_\pi$**

Input: the policy $\pi$ to be evaluated
Algorithm parameter: step size $\alpha \in (0, 1]$
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        $A \leftarrow$ action given by $\pi$ for $S$
        Take action $A$, observe $R$, $S'$
        $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$
        $S \leftarrow S'$
    until $S$ is terminal

**First-visit MC prediction, for estimating $V \approx v_\pi$**

Input: a policy $\pi$ to be evaluated
Initialize:
    $V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$
    $Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):
    Generate an episode following $\pi$: $S_0, A_0, R_1, S_1, A_1, R_2, \ldots, S_{T-1}, A_{T-1}, R_T$
    $G \leftarrow 0$
    Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
        $G \leftarrow \gamma G + R_{t+1}$
        Unless $S_t$ appears in $S_0, S_1, \ldots, S_{t-1}$:
            Append $G$ to $Returns(S_t)$
            $V(S_t) \leftarrow$ average($Returns(S_t)$)

*Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.*

# Recap
## Advantages and Disadvantages of MC and TD

- Which one should I use? Does it make any difference?

  - Bias/Variance Trade-Off

  - MC has high variance, but zero bias
    - good convergence (even with FA)
    - insensitive to initialization (no bootstrapping), simple to understand
    - only works for episodic problems (must wait until end of episode for update)
    - more efficient in non-Markov environments

  - TD has low variance, but some bias
    - TD(0) converges to $\pi_v(s)$ (be careful with FA: bias is a risk)
    - sensitive to initialization (because of the bootstrapping)
    - update after each step
    - exploits Markov property and is more efficient in Markov environment
    - **usually more efficient in practice**

Fraunhofer
IIS

# Exercise Sheet 4
## Model-free Prediction

Thank you for your attention!