

Markov Decision Process

Markov State \longrightarrow **State Transition**

$$P_{\{s_{t+1}|s_t\}} = P_{\{s_{t+1}|s_1, \dots, s_t\}} \quad P_{ss'} = P_{\{s_{t+1} = s' | s_t = s\}}$$

Markov Process

+

Reward Function + **Return**

$$R_s = E[R_{t+1} | S_t = s] \quad G_t = R_{t+1} + \gamma R_{t+2} + R_{t+3} + \dots$$

↓

Value Function $\xrightarrow{\text{Bellman}}$ $\left\{ \begin{matrix} R_{t+1} \\ \gamma V(S_{t+1}) \end{matrix} \right\} \longrightarrow V(s) = \mathcal{R}_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$

$$V(s) = E[G_t | S_t = s]$$

Markov Reward Process

+

Action and Policy
 $\pi(a|s) = P\{A_t = a | S_t = s\}$

↓

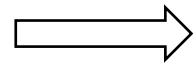
Value Function

$$V_\pi(s) = E_\pi[G_t | S_t = s] \quad q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a] \xrightarrow{\text{Bellman}} \begin{cases} V_\pi(s) = \sum_{a \in A} \pi(a|s) q_\pi(s, a) \\ q_\pi(s, a) = \mathcal{R}_s^a + \gamma P_{ss'}^a v_\pi(s') \end{cases}$$

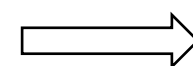
The optimal value functions are recursively related by the Bellman optimality equations.

$$v_*(s) = \max_a \mathcal{R}_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')$$

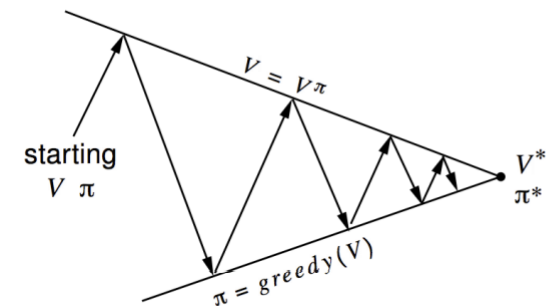
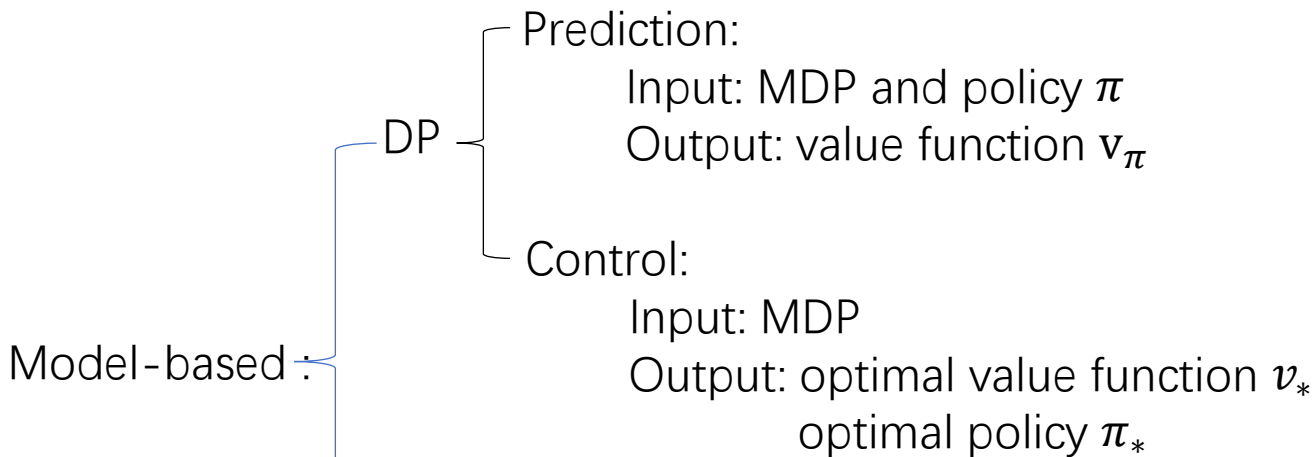
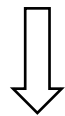
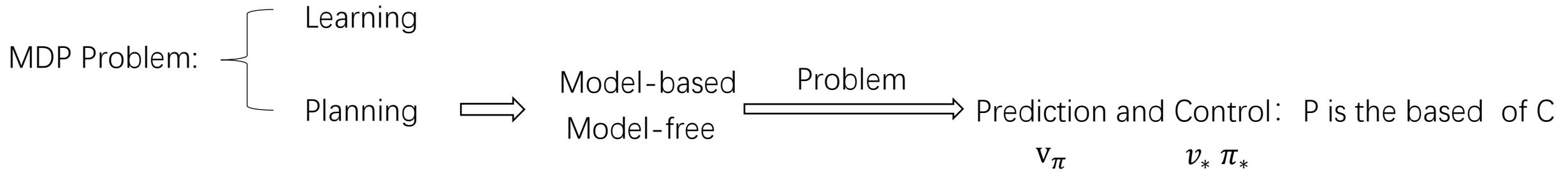
$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} q_*(s', a')$$



How to solve q_* and v_*

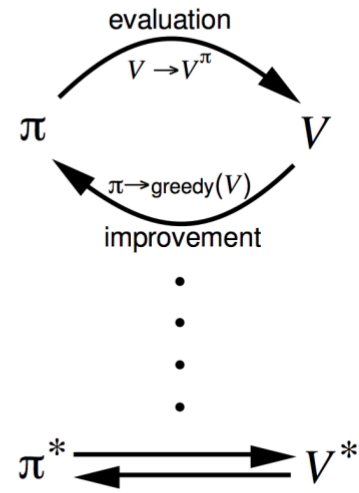


- Value Iteration
- Policy Iteration
- Q-learning
- Sarsa



Policy evaluation Estimate v_{π}
Iterative policy evaluation

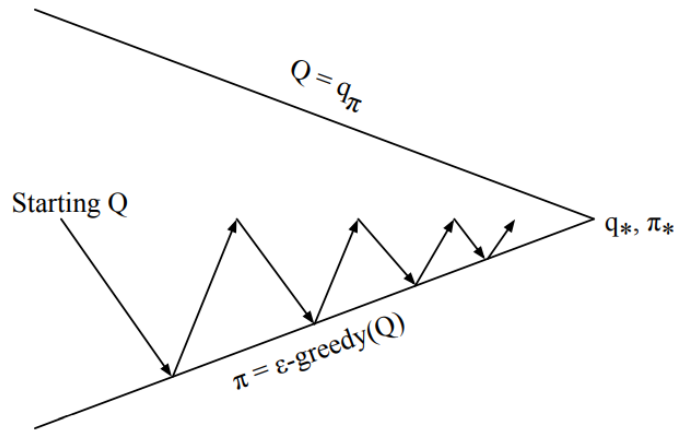
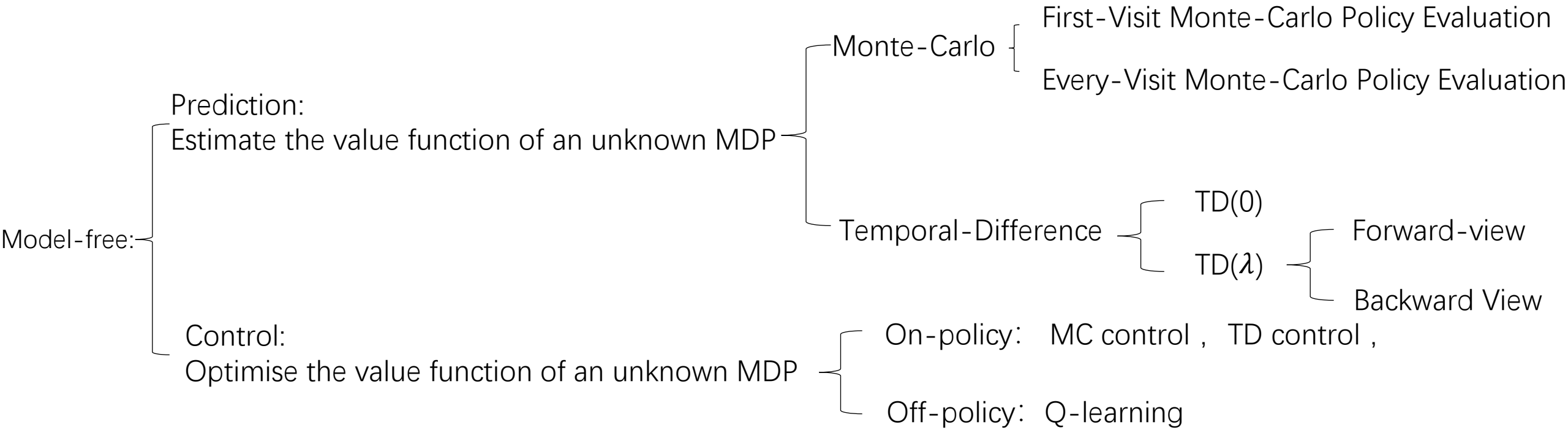
Policy improvement Generate $\pi' \geq \pi$
Greedy policy improvement



A control process

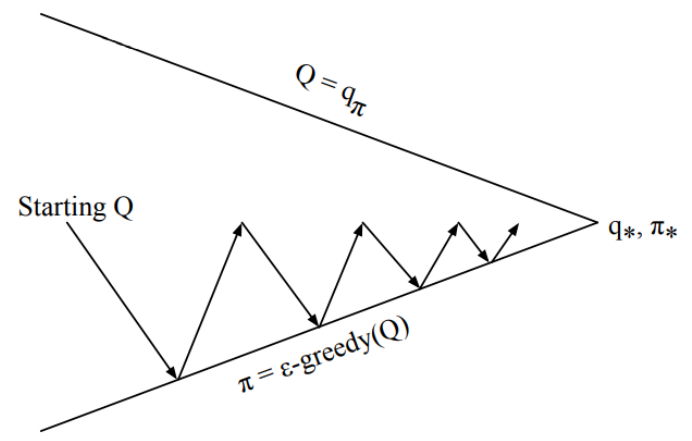
Algorithms

Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	Bellman Expectation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration



Every episode:
 Policy evaluation Monte-Carlo policy evaluation, $Q \approx q_\pi$
 Policy improvement ϵ -greedy policy improvement

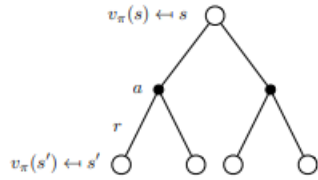

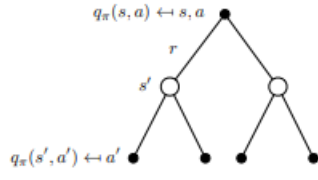
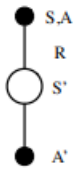
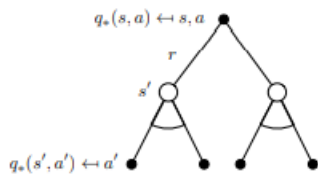
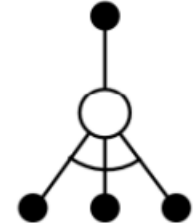
MC control



Every time-step:
 Policy evaluation Sarsa, $Q \approx q_\pi$
 Policy improvement ϵ -greedy policy improvement

Sarsa

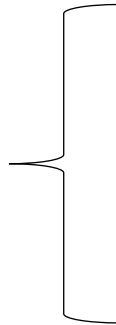
Relationship Between DP and TD

	<i>Full Backup (DP)</i>	<i>Sample Backup (TD)</i>
Bellman Expectation Equation for $v_{\pi}(s)$	 <p>Iterative Policy Evaluation</p>	 <p>TD Learning</p>
Bellman Expectation Equation for $q_{\pi}(s, a)$	 <p>Q-Policy Iteration</p>	 <p>Sarsa</p>
Bellman Optimality Equation for $q_*(s, a)$	 <p>Q-Value Iteration</p>	 <p>Q-Learning</p>

So what we've been talking about is the same thing: "discrete case" or **called lookup table**



Value Function Approximation
(Calculate w)



- a. Linear combinations of features: GD/LS
- b. Neural network
- c. Decision tree
- d. Nearest neighbour
- e. Fourier / wavelet bases
- ...

+ training method

- a. Incremental method
- b. Batch method

	"Compound"	Algorithm
Prediction →	GD + a	MC(on/off), TD(0)(on),
	LS + b	MC(on/off), LSMC(on/off), TD(on), LSTD(on/off)
	b + b	DQN
Control →	GD + a	MC, Sarsa, Gradient Q-learning
	LS + b	MC, Sarsa, LSPI

What we do is using **Approximate** in place of "previous V_{π}, q_{π} " and using training method to get w

Other form of table

Prediction

GD

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	TD(0)	✓	✓	✗
	TD(λ)	✓	✓	✗
Off-Policy	MC	✓	✓	✓
	TD(0)	✓	✗	✗
	TD(λ)	✓	✗	✗

LS

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✓	✗
Off-Policy	LSTD	✓	✓	-
	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✗	✗
	LSTD	✓	✓	-

Control

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-learning	✓	✗	✗
Gradient Q-learning	✓	✓	✗

(✓) = chatters around near-optimal value function

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-learning	✓	✗	✗
LSPI	✓	(✓)	-

(✓) = chatters around near-optimal value function