

Reinforcement learning: day 2

Variability in modular interchangeable RL methods



Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

- Fly stunt manoeuvres in a helicopter
- Defeat the world champion at Backgammon
- Manage an investment portfolio
- Control a power station
- Make a humanoid robot walk
- Play many different Atari games better than humans
- Industrial control
- Production control
- Automotive control
- Autonomous vehicles control
- Logistics
- Telecommunication networks
- Sensor networks
- Finance
- Games

Intro

Applications

Interdisciplinary context

CS context

Review

Policy iteration

Q-learning

λ methods

SARSA(λ)

RL context

Feature-based

Example features

Value approximation

Which functions?

Architectures?

Linear / Non-linear?

Gradient descent

Linear regression

Update rules

Approximate control

SARSA linear reg

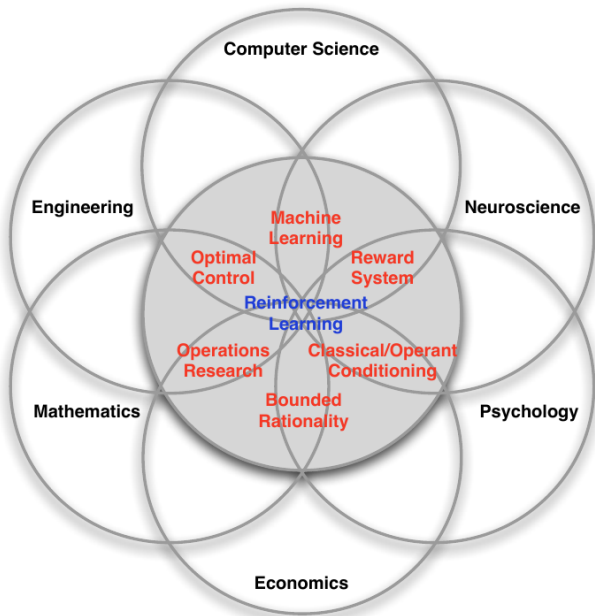
Policy-gradient methods

Examples

Games

Control

Interdisciplinary context



Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

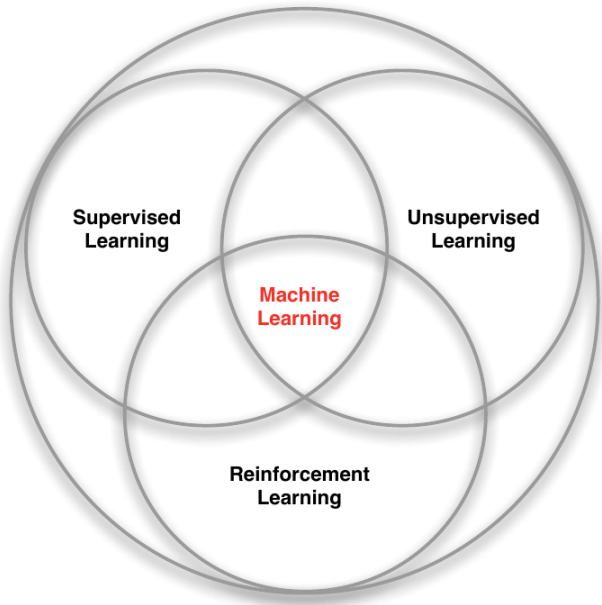
Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control



Intro

Applications
Interdisciplinary
context

CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

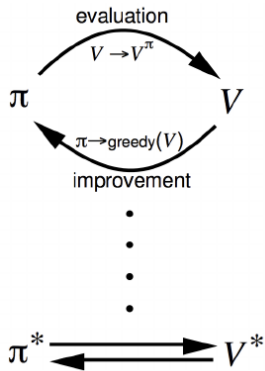
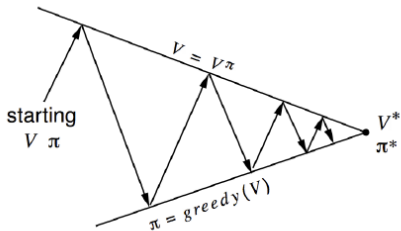
Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Review: Policy iteration



- Policy evaluation** Estimate v_π
Any policy evaluation algorithm
- Policy improvement** Generate $\pi' \geq \pi$
Any policy improvement algorithm

Intro

- Applications
- Interdisciplinary context
- CS context

Review

- Policy iteration**
- Q-learning

λ methods

- SARSA(λ)
- RL context

Feature-based

- Example features
- Value approximation
- Which functions?
- Architectures?
- Linear / Non-linear?
- Gradient descent
- Linear regression
- Update rules
- Approximate control
- SARSA linear reg

Policy-gradient methods

Examples

- Games
- Control

Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$;

 until S is terminal

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Recall: To maximally use all your data, it makes sense to update as many of your states as possible for each new data point.

Intro

- Applications
- Interdisciplinary context
- CS context

Review

- Policy iteration
- Q-learning

λ methods

- SARSA(λ)
- RL context

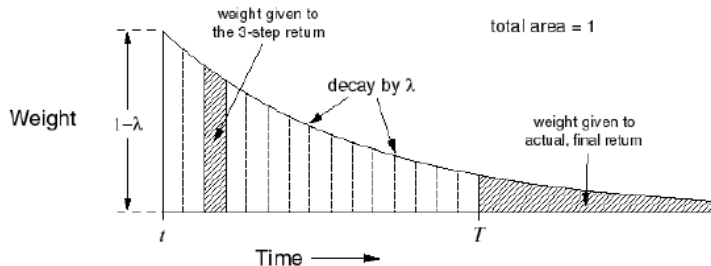
Feature-based

- Example features
- Value approximation
- Which functions?
- Architectures?
- Linear / Non-linear?
- Gradient descent
- Linear regression
- Update rules
- Approximate control
- SARSA linear reg

Policy-gradient methods

Examples

- Games
- Control



Change Q or V values for states visited farther back less.

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

SARSA(λ) keeps an eligibility trace

initialize $Q[S,A]$ arbitrarily

initialize $e[s,a]=0$ for all s,a

observe current state s

select action a using a policy based on Q

repeat forever:

 carry out an action a

 observe reward r and state s'

 select action a' using a policy based on Q

$$\delta \leftarrow r + \gamma Q[s',a'] - Q[s,a]$$

$$e[s,a] \leftarrow e[s,a] + \delta$$

 for all s'',a'' :

$$Q[s'',a''] \leftarrow Q[s'',a''] + \alpha \delta e[s'',a'']$$

$$e[s'',a''] \leftarrow \gamma \lambda e[s'',a'']$$

$s \leftarrow s'$

$a \leftarrow a'$

end-repeat

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

SARSA(λ) keeps an eligibility trace

Initialize $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat (for each episode):

$E(s, a) = 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Initialize S, A

Repeat (for each step of episode):

Take action A , observe R, S'

Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)

$\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$

$E(S, A) \leftarrow E(S, A) + 1$

For all $s \in \mathcal{S}, a \in \mathcal{A}(s)$:

$Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$

$E(s, a) \leftarrow \gamma \lambda E(s, a)$

$S \leftarrow S'; A \leftarrow A'$

until S is terminal

This will be important for your assignment!

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

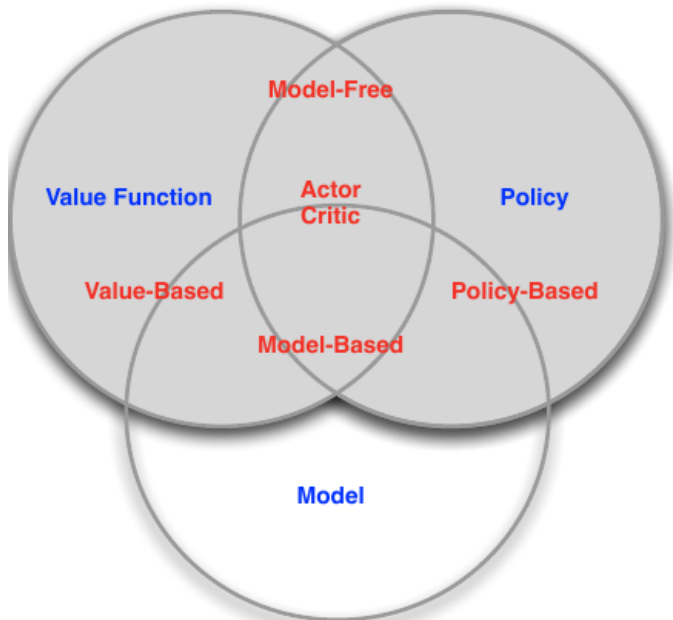
Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control



Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)

RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

What about very large state spaces and continuous problems?

Intro

- Applications
- Interdisciplinary context
- CS context

Review

- Policy iteration
- Q-learning

λ methods

- SARSA(λ)
- RL context

Feature-based

- Example features
- Value approximation
- Which functions?
- Architectures?
- Linear / Non-linear?
- Gradient descent
- Linear regression
- Update rules
- Approximate control
- SARSA linear reg

Policy-gradient methods

Examples

- Games
- Control

Reinforcement learning can be used to solve large problems, e.g.,

Backgammon: 10^{20} states

Computer Go: 10^{170} states

Helicopter: continuous state space

How can we scale up the model-free methods for prediction and control?

Intro

- Applications
- Interdisciplinary context
- CS context

Review

- Policy iteration
- Q-learning

λ methods

- SARSA(λ)
- RL context

Feature-based

- Example features
- Value approximation
- Which functions?
- Architectures?
- Linear / Non-linear?
- Gradient descent
- Linear regression
- Update rules
- Approximate control
- SARSA linear reg

Policy-gradient methods

Examples

- Games
- Control

- **flat** or modular or hierarchical
- explicit states or **features** or individuals and relations
- static or finite stage or **indefinite stage or infinite stage**
- **fully observable** or partially observable
- deterministic or **stochastic** dynamics
- goals or **complex preferences**
- **single agent** or multiple agents
- knowledge is given or **knowledge is learned**
- **perfect rationality** or bounded rationality

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

- Usually we don't want to reason in terms of states, but in terms of features.
- In state-based methods, information about one state cannot be used by similar states.
- If there are too many parameters to learn, it takes too long.
- **Idea:** Express the value function as a function of the features. Most typical is a linear function of the features.
- $Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$
- $V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Example Features

- $F_1(s, a) = 1$ if a goes from state s into a monster location and is 0 otherwise.
- $F_2(s, a) = 1$ if a goes into a wall, is 0 otherwise.
- $F_3(s, a) = 1$ if a goes toward a prize.
- $F_4(s, a) = 1$ if the agent is damaged in state s and action a takes it toward the repair station.
- $F_5(s, a) = 1$ if the agent is damaged and action a goes into a monster location.
- $F_6(s, a) = 1$ if the agent is damaged.
- $F_7(s, a) = 1$ if the agent is not damaged.
- $F_8(s, a) = 1$ if the agent is damaged and there is a prize in direction a .
- $F_9(s, a) = 1$ if the agent is not damaged and there is a prize in direction a .

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features

Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Example Features

- $F_{10}(s, a)$ is the distance from the left wall if there is a prize at location P_0 , and is 0 otherwise.
- $F_{11}(s, a)$ has the value $4 - x$, where x is the horizontal position of state s if there is a prize at location P_0 ; otherwise is 0.
- $F_{12}(s, a)$ to $F_{29}(s, a)$ are like F_{10} and F_{11} for different combinations of the prize location and the distance from each of the four walls.
For the case where the prize is at location P_0 , the y -distance could take into account the wall.

Example function:

$$Q(s, a) = 2.0 - 1.0 * F_1(s, a) - 0.4 * F_2(s, a) - 1.3 * F_3(s, a) - 0.5 * F_4(s, a) - 1.2 * F_5(s, a) - 1.6 * F_6(s, a) + 3.5 * F_7(s, a) \dots$$

Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features

Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Value Function Approximation

So far we have represented value function by a lookup table

- Every state s has an entry $V(s)$
- Or every state-action pair s, a has an entry $Q(s, a)$

Problem with large MDPs:

- There are too many states and/or actions to store in memory
- It is too slow to learn the value of each state individually

Solution for large MDPs:

- Estimate value function with function approximation
 $\hat{v}(s, \mathbf{w}) \approx v_{\pi}(s)$
 or $\hat{q}(s, a, \mathbf{w}) \approx q_{\pi}(s, a)$
- Generalise from seen states to unseen states
- Update parameter w using MC or TD learning

Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Which Function Approximator?

There are many function approximators, e.g.

Linear combinations of features

Neural network

Decision tree

Nearest neighbour

Fourier / wavelet bases

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation

Which functions?

Architectures?
Linear / Non-linear?

Gradient descent

Linear regression

Update rules

Approximate control

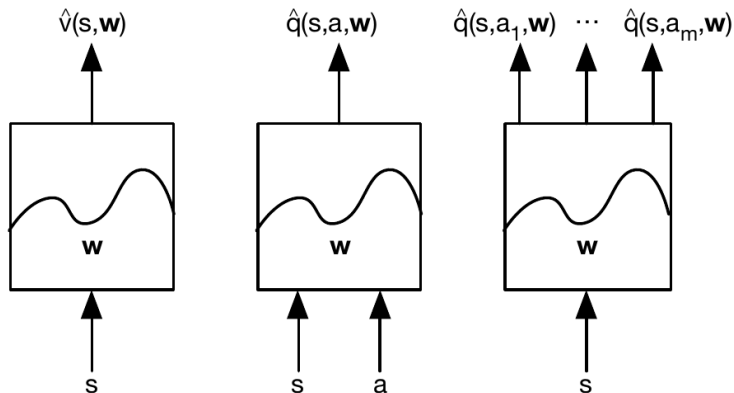
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Function approximation architectures



- 1) v 2) q , action-in 3) q , action-out (deepmind)

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

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

Not all function approximation methods will converge with RL-algorithms. Non-linear methods like multi-layer networks are particularly problematic.

Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?

Linear / Non-linear?

Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

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

Some new TD methods are more robust

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

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

(✓) = chatters around near-optimal value function

Still, the non-linear methods are hard.

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?

Linear / Non-linear?

Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Gradient descent

To find a (local) minimum of a real-valued function $f(x)$:

- assign an arbitrary value to x
- repeat

$$x \leftarrow x - \eta \frac{df}{dx}$$

where η is the step size

To find a local minimum of real-valued function $f(x_1, \dots, x_n)$:

- assign arbitrary values to x_1, \dots, x_n
- repeat:

for each x_i

$$x_i \leftarrow x_i - \eta \frac{\partial f}{\partial x_i}$$

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?

Gradient descent

Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Linear Regression

- A linear function of variables x_1, \dots, x_n is of the form

$$f^{\bar{w}}(x_1, \dots, x_n) = w_0 + w_1x_1 + \dots + w_nx_n$$

$\bar{w} = \langle w_0, w_1, \dots, w_n \rangle$ are weights. (Let $x_0 = 1$).

- Given a set E of examples.
Example e has input $x_i = e_i$ for each i and observed value, o_e :

$$Error_E(\bar{w}) = \sum_{e \in E} (o_e - f^{\bar{w}}(e_1, \dots, e_n))^2$$

- Minimizing the error using gradient descent, each example should update w_i using:

$$w_i \leftarrow w_i - \eta \frac{\partial Error_E(\bar{w})}{\partial w_i}$$

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Gradient Descent for Linear Regression

Given E : set of examples over n features
 each example e has inputs (e_1, \dots, e_n) and output o_e :
 Assign weights $\bar{w} = \langle w_0, \dots, w_n \rangle$ arbitrarily
repeat:

For each example e in E :
 let $\delta = o_e - f^{\bar{w}}(e_1, \dots, e_n)$
For each weight w_i :
 $w_i \leftarrow w_i + \eta \delta e_i$

Intro

Applications
 Interdisciplinary
 context
 CS context

Review

Policy iteration
 Q-learning

λ methods

SARSA(λ)
 RL context

Feature-based

Example features
 Value approximation
 Which functions?
 Architectures?
 Linear / Non-linear?
 Gradient descent
Linear regression
 Update rules
 Approximate control
 SARSA linear reg

Policy-gradient methods

Examples

Games
 Control

Update rules for function approximation

With V:

$$\nabla_w \hat{v}(S, w) = x(S)$$

Update = step-size (reward + prediction error x feature values)

$$\begin{aligned} \Delta_w &= \alpha(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla_w \hat{v}(S_t, \mathbf{w}) \\ &= \alpha \delta x(S) \end{aligned}$$

or with Q:

$$\Delta_w = \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{v}(S_t, A_t, \mathbf{w})) \nabla_w \hat{v}(S, A_t, \mathbf{w})$$

Note: \mathbf{w} is a vector, which means this happens on a matrix or array all at once!

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

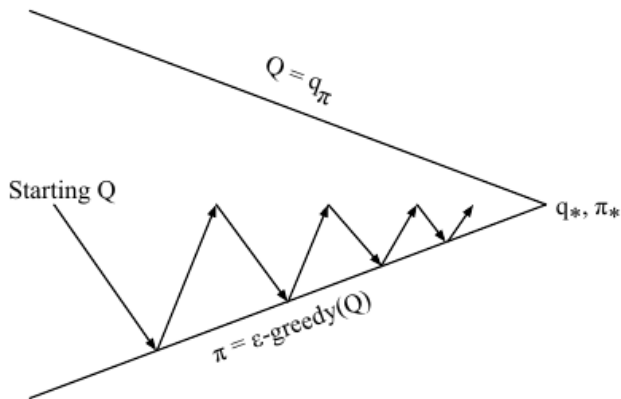
Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control



Every **time-step**:

Policy evaluation **Sarsa**, $Q \approx q_\pi$

Policy improvement ϵ -greedy policy improvement

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules

Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

- One step backup provides the examples that can be used in a linear regression.
- Suppose F_1, \dots, F_n are the features of the state and the action.
- So $Q_{\bar{w}}(s, a) = w_0 + w_1 F_1(s, a) + \dots + w_n F_n(s, a)$
- An experience $\langle s, a, r, s', a' \rangle$ provides the “example”:
 - ▶ **old predicted value:** $Q_{\bar{w}}(s, a)$
 - ▶ **new “observed” value:** $r + \gamma Q_{\bar{w}}(s', a')$

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

SARSA with linear function approximation

Given γ :discount factor; η :step size

Assign weights $\bar{w} = \langle w_0, \dots, w_n \rangle$ arbitrarily

observe current state s

select action a

repeat forever:

 carry out action a

 observe reward r and state s'

 select action a' (using a policy based on $Q_{\bar{w}}$)

 let $\delta = r + \gamma Q_{\bar{w}}(s', a') - Q_{\bar{w}}(s, a)$

 For $i = 0$ to n

$$w_i \leftarrow w_i + \eta \delta F_i(s, a)$$

$s \leftarrow s'$

$a \leftarrow a'$

Bonus points: If you can get this working in the assignment

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

controller SARSA-FA(F, γ, η)

2: **Inputs**

3: $F = \langle F_1, \dots, F_n \rangle$: a set of features

4: $\gamma \in [0, 1]$: discount factor

5: $\eta > 0$: step size for gradient descent

6: **Local**

7: weights $w = \langle w_0, \dots, w_n \rangle$, initialized arbitrarily

8: observe current state s

9: select action a

10: **repeat**

11: carry out action a

12: observe reward r and state s'

13: select action a' (using a policy based on Q_w)

14: let $\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$

15: **for** $i=0$ to n **do**

16: $w_i \leftarrow w_i + \eta \delta F_i(s, a)$

17:

18: $s \leftarrow s'$

19: $a \leftarrow a'$

20: **until** termination

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Value Based

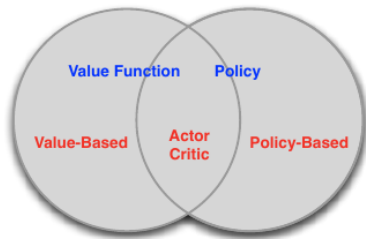
- Learn Value Function
- Implicit policy
- (e.g. ϵ -greedy)

Policy Based

- No Value Function
- Learn Policy

Actor-Critic

- Learn Value Function
- Learn Policy



Improve on evolutionary methods mentioned at beginning of last class

Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Games

Fly stunt manoeuvres in a helicopter

Defeat the world champion at Backgammon

Manage an investment portfolio

Control a power station

Make a humanoid robot walk

Play many different Atari games better than humans

Industrial control

Production control

Automotive control

Autonomous vehicles control

Logistics

Telecommunication networks

Sensor networks

Finance

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Current game software (not all RL)

Program	Level of Play	Program to Achieve Level
Checkers	Perfect	<i>Chinook</i>
Chess	Superhuman	<i>Deep Blue</i>
Othello	Superhuman	<i>Logistello</i>
Backgammon	Superhuman	<i>TD-Gammon</i>
Scrabble	Superhuman	<i>Maven</i>
Go	Grandmaster	<i>MoGo</i> ¹ , <i>Crazy Stone</i> ² , <i>Zen</i> ³
Poker ⁴	Superhuman	<i>Polaris</i>

¹9 × 9

²9 × 9 and 19 × 19

³19 × 19

⁴Heads-up Limit Texas Hold'em

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Current RL game software

Program	Level of Play	RL Program to Achieve Level
Checkers	Perfect	<i>Chinook</i>
Chess	International Master	<i>KnightCap / Meep</i>
Othello	Superhuman	<i>Logistello</i>
Backgammon	Superhuman	<i>TD-Gammon</i>
Scrabble	Superhuman	<i>Maven</i>
Go	Grandmaster	<i>MoGo</i> ¹ , <i>Crazy Stone</i> ² , <i>Zen</i> ³
Poker ⁴	Superhuman	<i>SmooCT</i>

¹9 × 9

²9 × 9 and 19 × 19

³19 × 19

⁴Heads-up Limit Texas Hold'em

Though this is a recent table, DeepMind and AlphaGo were more recently

Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Program	Input features	Value Fn	RL	Training	Search
Chess <i>Meep</i>	Binary <i>Pieces, pawns, ...</i>	Linear	TreeStrap	Self-Play / Expert	$\alpha\beta$
Checkers <i>Chinook</i>	Binary <i>Pieces, ...</i>	Linear	TD leaf	Self-Play	$\alpha\beta$
Othello <i>Logistello</i>	Binary <i>Disc configs</i>	Linear	MC	Self-Play	$\alpha\beta$
Backgammon <i>TD Gammon</i>	Binary <i>Num checkers</i>	Neural network	TD(λ)	Self-Play	$\alpha\beta$ / MC
Go <i>MoGo</i>	Binary <i>Stone patterns</i>	Linear	TD	Self-Play	MCTS
Scrabble <i>Maven</i>	Binary <i>Letters on rack</i>	Linear	MC	Self-Play	MC search
Limit Hold'em <i>SmooCT</i>	Binary <i>Card abstraction</i>	Linear	MCTS	Self-Play	-

Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control

Fashion models and financial models are similar.
They bear a similar relationship to everyday world.

Like supermodels, financial models are idealized
representations of the real world,
they are not real, they don't quite work the way that the
real world works.

There is celebrity in both worlds.
In the end, there is the same inevitable disappointment”
- Satyajit Das

Some popular deep methods are **Q-learning with an action-out convolutional network as the feature approximator**

Intro

- Applications
- Interdisciplinary context
- CS context

Review

- Policy iteration
- Q-learning

λ methods

- SARSA(λ)
- RL context

Feature-based

- Example features
- Value approximation
- Which functions?
- Architectures?
 - Linear / Non-linear?
- Gradient descent
- Linear regression
- Update rules
- Approximate control
- SARSA linear reg

Policy-gradient methods

Examples

- Games
- Control

Games

Fly stunt manoeuvres in a helicopter

Defeat the world champion at Backgammon

Manage an investment portfolio

Control a power station

Make a humanoid robot **walk**

Play many different Atari games better than humans

Industrial **control**

Production **control**

Automotive **control**

Autonomous vehicles **control**

Logistics

Telecommunication networks

Sensor networks

Finance

Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

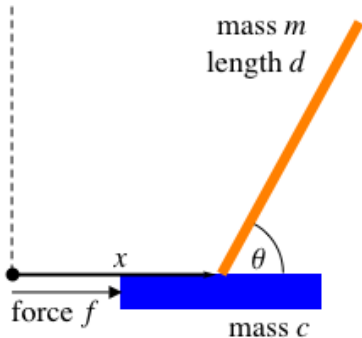
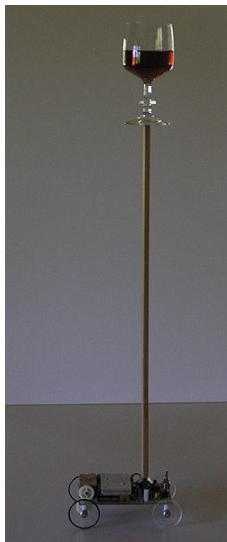
Policy-gradient methods

Examples

Games
Control

Pole-cart / inverted pendulum

p. 39



Intro

Applications
Interdisciplinary
context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

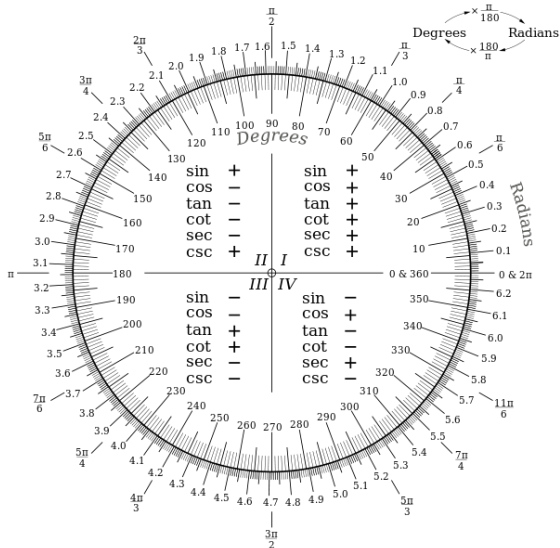
Games
Control

https://www.youtube.com/watch?v=Ep21NMic_fk



Pole-cart / inverted pendulum

Your assignment: Keep the pole vertical at $\pi/2$



Intro

Applications
Interdisciplinary context
CS context

Review

Policy iteration
Q-learning

λ methods

SARSA(λ)
RL context

Feature-based

Example features
Value approximation
Which functions?
Architectures?
Linear / Non-linear?
Gradient descent
Linear regression
Update rules
Approximate control
SARSA linear reg

Policy-gradient methods

Examples

Games
Control