## Agent-Environment Interface
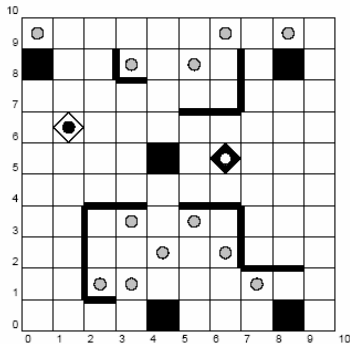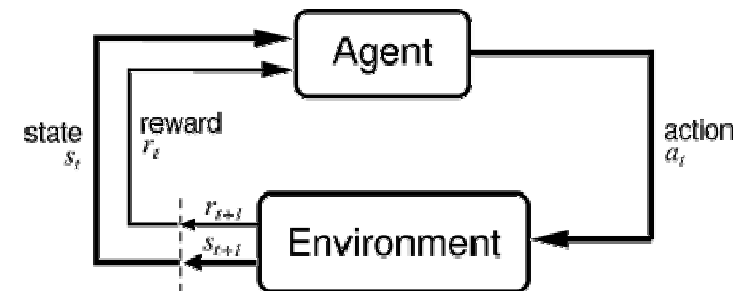
Markov Decision Processes, Dynamic Programming, and Reinforcement Learning in R

Jeffrey Todd Lins
Thomas Jakobsen
Saxo Bank A/S

jtl@saxobank.com, tj@saxobank.com

Source: Sutton & Barto, 2001

useR! 2006
Vienna, June 15-17, 2006

SAXO BANK

useR! 2006
Vienna, June 15-17, 2006

SAXO BANK

## Markov Decision Process

We define a Markov Decision Process as a tuple $(\mathcal{S}, \mathcal{A}, T, R)$ where

- $\mathcal{S}$ is a finite set of states
- $\mathcal{A}$ is a finite set of actions
- $T : \mathcal{S} \times \mathcal{A} \to \Pi(S)$ is the transition model giving a probability distribution over all states for ending in a future state, $s'$, given that an agent takes action, $a$ in state $s$.
- $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a real-valued reward function yielding the immediately expected reward for taking each action in each state.

## Dynamic Programming

- Deterministic Policy
$$\pi : \mathcal{S} \to \mathcal{A}$$
- Stochastic Policy
$$\pi : \mathcal{S} \to \Pi(\mathcal{A})$$
- State Value Function
$$V_\pi(s) = E_\pi\Big[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_t = s\Big],$$

where $0 < \gamma < 1$ is a discount factor that controls how much influence future rewards have, and $r_t$ is the reward received at time $t$.

- State-Action Value Function
$$Q_\pi(s,a) = E_\pi\Big[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_t = s, a_t = a\Big].$$

useR! 2006
Vienna, June 15-17, 2006

SAXO BANK

useR! 2006
Vienna, June 15-17, 2006

SAXO BANK

## Bellman Equation

- Bellman Equation

$$Q_\pi(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') Q_\pi(s',a).$$

or in matrix notation

$$Q_\pi = R + \gamma \mathbf{T} \mathbf{\Pi}_\pi Q_\pi.$$

- Now we consider that an optimal policy, $\pi^*$, will satisfy

$$\pi^* = \arg\max_\pi E\Big[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi\Big].$$

## Bellman Optimality Equation

- It can now be shown that there is a policy,

$$\pi^* = \arg\max_a E\Big[\sum_{s'} T(s,a,s') V(s') \mid \pi\Big]$$

that is optimal for the value in a *subsequent* state.

- Bellman Optimality Equation

$$Q_{\pi^*}(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_a Q_{\pi^*}(s',a).$$

or in matrix notation

$$Q_{\pi^*} = R + \gamma \mathbf{T} \mathbf{\Pi}_{\max} Q_{\pi^*}.$$

## Value Iteration

```
// S : States
// A : Actions
// T : Transition Model
// R : Reward Function
// γ : Discount Factor
// Q₀ : Initial State-Action Value Function
// ε : Stopping criterion
Q' ← Q₀
repeat
   Q ← Q'
   Qπ ← R + γTᵐᵃˣQ
until |Q − Q'| < ε
∀s ∈ S, π'(s) ← arg maxₐ Q'(s,a)
return  π
```

## Policy Iteration

```
// S : States
// A : Actions
// T : Transition Model
// R : Reward Function
// γ : Discount Factor
// π₀ : Initial Policy

π' ← π₀
repeat
   π ← π'
   Qπ ← (I − γTπ)⁻¹R
   ∀s ∈ S, π'(s) ← arg maxₐ Qπ(s,a)
until π ← π'
return  π
```

# Reinforcement Learning

- **Temporal Difference (TD) Learning**(Sutton, 1988) yields the state value function, $V_\pi$, for a fixed policy, given a sample set $(s, a, r, s')$

$$\widehat{V}_{t+1}(s) = \widehat{V} + \alpha\left[r + \gamma\widehat{V}_t(s') - \widehat{V}_t(s)\right]$$

.

- **Q-Learning**(Watkins, 1989) yields an optimal policy, $\pi^*$, by an approximation of $Q_{\pi^*}$, for a fixed policy, given a sample set $(s, a, r, s')$

$$\widehat{Q}_{t+1}(s, a) = \widehat{Q} + \alpha\left[r + \gamma\widehat{Q}_t(s', a') - \widehat{Q}_t(s, a)\right]$$

# Temporal Difference Learning

```
// D : Samples (s, a, r, s')
// A : Actions
// α₀ : Initial Learning Rate
// γ : Discount Factor
// V₀ : Initial Value Function
// π :Policy
```

$\tilde{V} \leftarrow V_0, \alpha \leftarrow \alpha_0, t \leftarrow 0$
**for** $(s, a, r, s') \in D(\pi)$ **do**
    $\tilde{V}(s) \leftarrow \tilde{V}(s) + \alpha(r + \gamma\tilde{V}(s') - \tilde{V}(s,))$
    $\alpha \leftarrow \sigma(\alpha, \alpha_0, t)$
    $t \leftarrow t + 1$
**end for**

**return** $\tilde{V}$

# Q-Learning

```
// D : Samples (s, a, r, s')
// A : Actions
// α₀ : Initial Learning Rate
// γ : Discount Factor
// Q₀ : Initial State-Action Value Function
// π : Exploration Policy
```

$\tilde{Q} \leftarrow Q_0, \alpha \leftarrow \alpha_0, t \leftarrow 0$
**for** $(s, a, r, s') \in D(\pi, Q)$ **do**
    $\tilde{Q}(s, a) \leftarrow \tilde{Q}(s, a) + \alpha(r + \gamma\max\tilde{Q}(s', a') - \tilde{Q}(s, a))$
    $\alpha \leftarrow \sigma(\alpha, \alpha_0, t)$
    $t \leftarrow t + 1$
**end for**

**return** $\tilde{Q}$

# Linear Architectures

**Linear Approximation Architectures**
- Basis Functions: $\phi(s, a)$
- Weights: $w_i$
- $k$: column vector of size $|\mathcal{S}||\mathcal{A}|$

$$Q_{\pi,w}(s, a) = \sum_{i=1}^{k}\phi_i(s, a)w_{i,\pi} = \phi(s, a)^\top w_\pi.$$

## Least Squares TD Learning

```
// D : Samples (s, a, r, s')
// A : Actions
// k : Number of basis functions
// γ : Discount Factor
// V₀ : Initial Value Function
// π : Policy
```

$$\mathbf{A} \leftarrow 0, b \leftarrow 0$$
$$\mathbf{for}\ (s, a, r, s') \in D(\pi)\ \mathbf{do}$$
$$\quad \tilde{A} \leftarrow \tilde{A} + \phi(s)(\phi(s) - \phi(s'))$$
$$\quad \tilde{b} \leftarrow \tilde{b} + \phi(s)r$$
$$\quad w_\pi \leftarrow \tilde{A}^{-1}\tilde{b}$$
$$\mathbf{end\ for}$$

$$\mathbf{return}\ \ \tilde{w}_\pi$$

## Examples of RL in Finance

*Performance Functions and Reinforcement Learning for Trading Systems and Portfolios.*

John Moody, Lizhong Wu, Yuansong Liao & Matthew Saffell. Journal of Forecasting, Volume 17, Pages 441-470, 1998.

*Intraday FX trading: Reinforcement learning vs evolutionary learning.*

M. A. H. Dempster, T. W. Payne, & V. S. Romahi. Working Paper No. 23/01, Judge Institute of Management, University of Cambridge, 2001.

## Advantages of RL in R

- Vectorized Programming

- Flexible, Interactive Simulation Environment

- Wide Range of Possibilities for Linear Basis Functions

- Interface to Existing Packages: HMMs, SVMs, GAs, Neural Networks

## References

Richard Sutton and Andrew Barto. Reinforcement Learning: An Introduction. The MIT Press, Cambridge, Massachusetts, 1998.

Michail G. Lagoudakis and Ronald Parr. "Least-Squares Policy Iteration," *Journal of Machine Learning Research*, 4, 2003, pp. 1107-1149.