



An Introduction to Reinforcement Learning

1

Yuan YAO

HKUST

Supervised Learning

- **Data:** (x, y)
x is input, y is output/response (label)
- **Goal:** Learn a *function* to map $x \rightarrow y$
- **Examples:**
 - Classification,
 - regression,
 - object detection,
 - semantic segmentation,
 - image captioning, etc.



➤ Cat

Unsupervised Learning

- **Data:** x
Just input data, no output labels!
- **Goal:** Learn some underlying hidden *structure* of the data
- **Examples:**
 - Clustering,
 - dimensionality reduction (manifold learning),
 - feature learning,
 - density estimation,
 - Generative models and GANs, etc.

Generative Models

Given training data, generate new samples from same distribution



Training data $\sim p_{\text{data}}(x)$

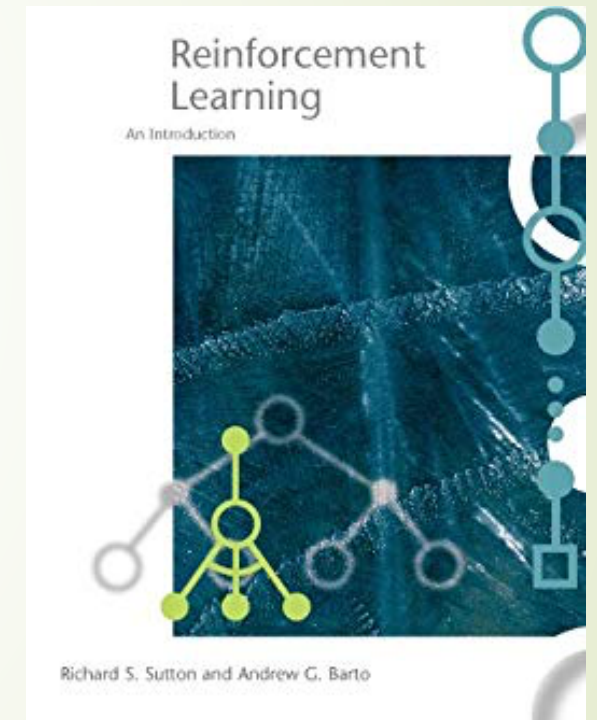
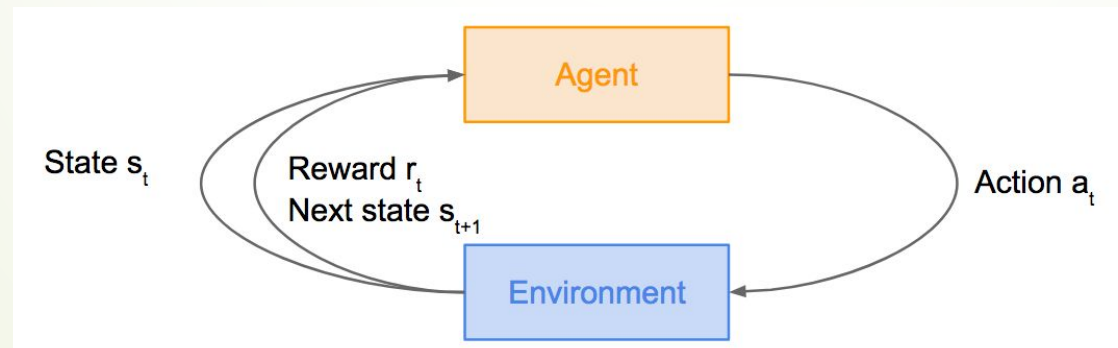


Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

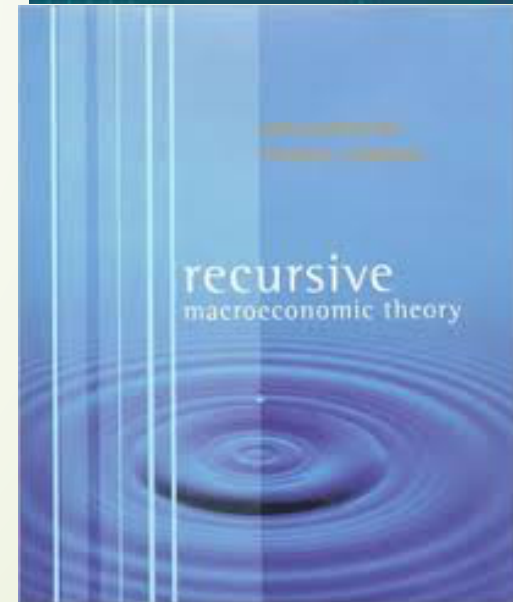
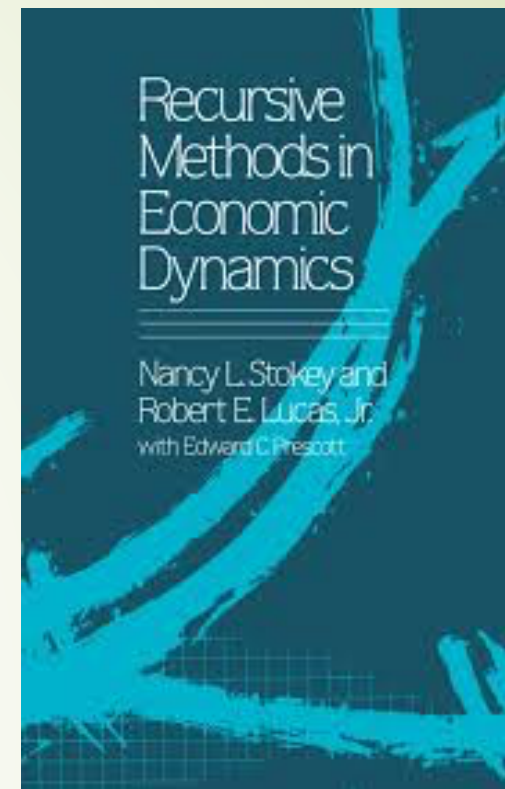
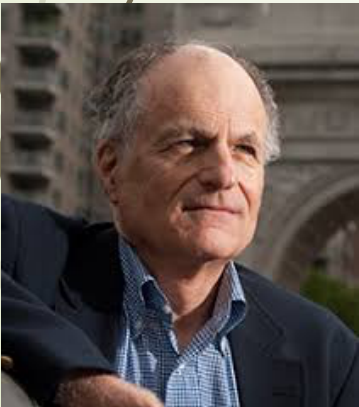
Today: Reinforcement Learning

- Problems involving an **agent**
- interacting with an **environment**,
- which provides numeric **reward** signals
- **Goal:**
 - Learn how to take actions in order to maximize reward in dynamic scenarios



Markov Decision Process /Dynamic Programming in Economics

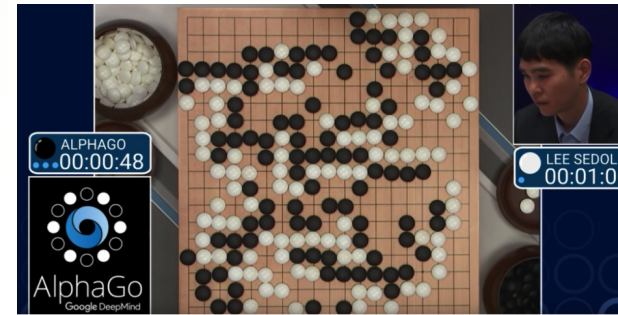
- ▶ The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 1995 was awarded to **Robert E. Lucas Jr. "for having developed and applied the hypothesis of rational expectations, and thereby having transformed macroeconomic analysis and deepened our understanding of economic policy"**.
- ▶ **Thomas John Sargent** was awarded the Nobel Memorial Prize in Economics in 2011 together with Christopher A. Sims for their **"empirical research on cause and effect in the macroeconomy"**



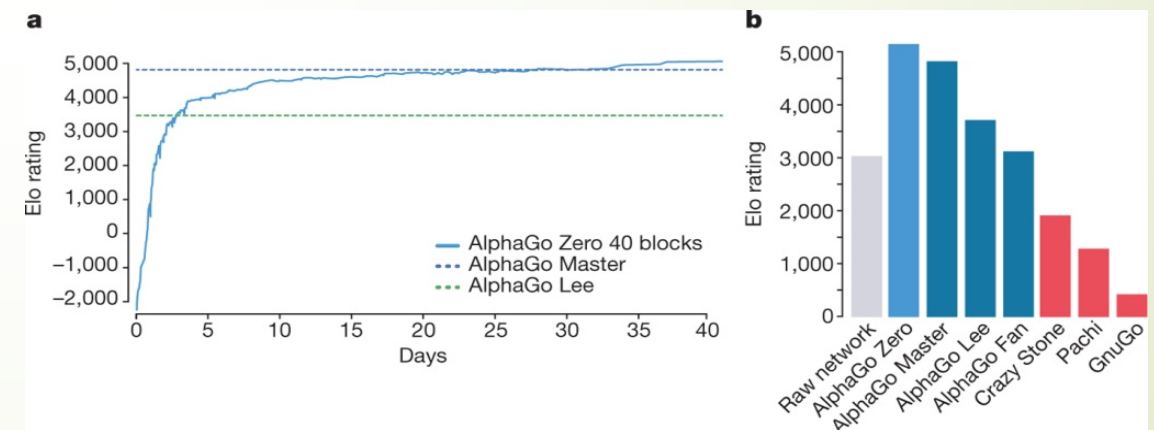
Playing games against human champions



Deep Blue in 1997



AlphaGo "LEE" 2016





Outline

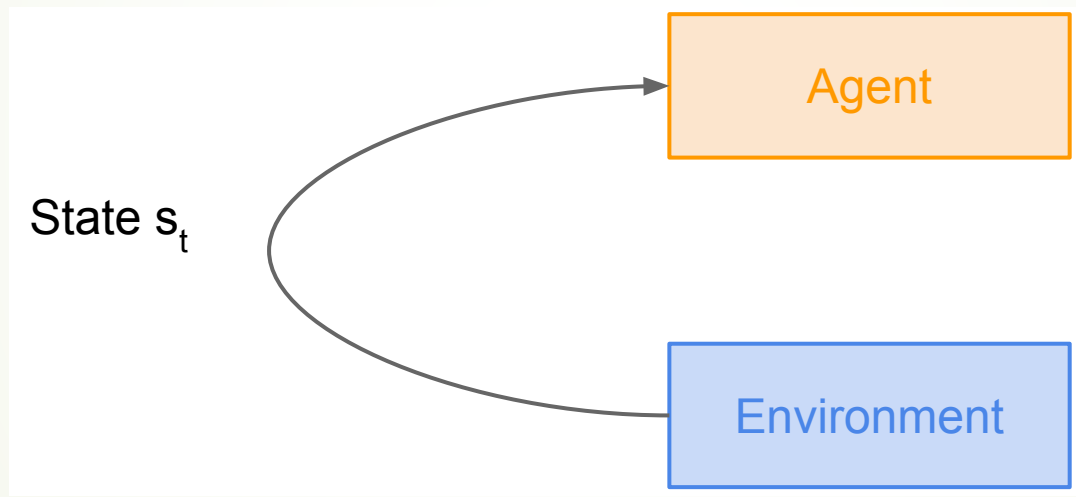


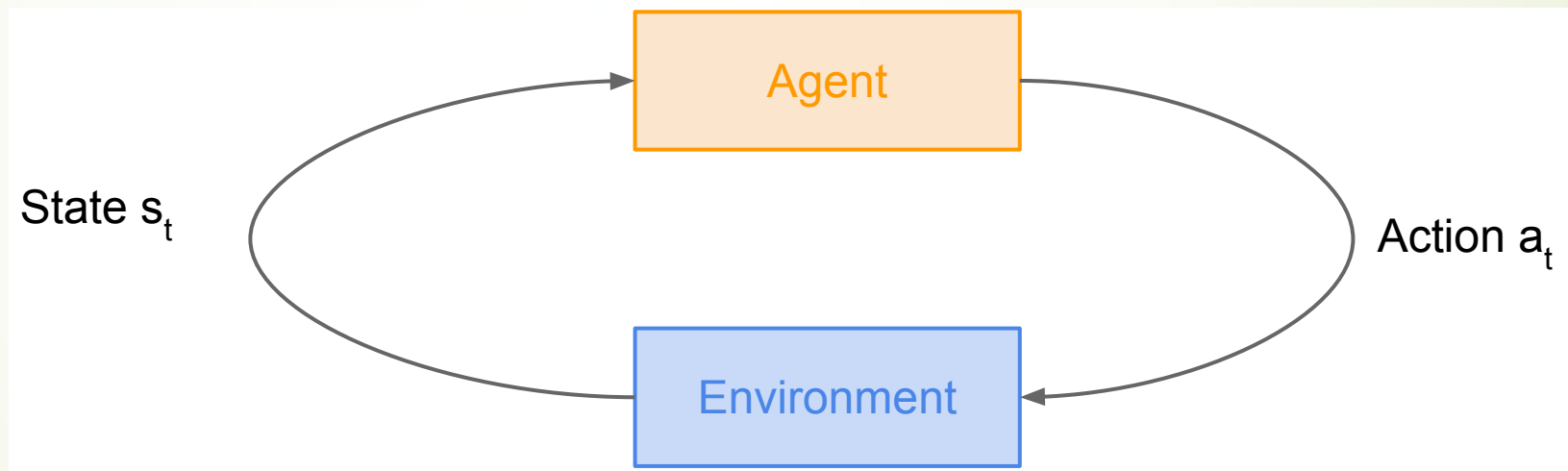
- What is Reinforcement Learning?
- Markov Decision Processes
- Bellman Equation as Linear Programming
- Q-Learning
- Policy Gradients
- Actor-Critics (Q-learning+Policy gradient)
- An Example of Order Book Optimization via Discrete Q-Learning by Prof. Michael Kearns

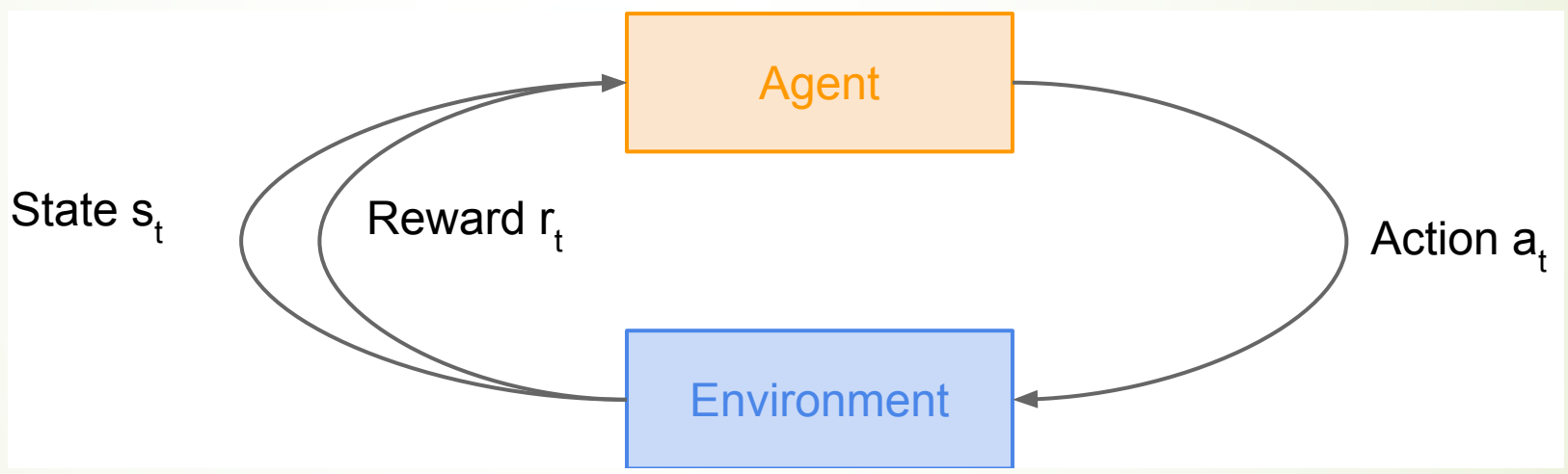


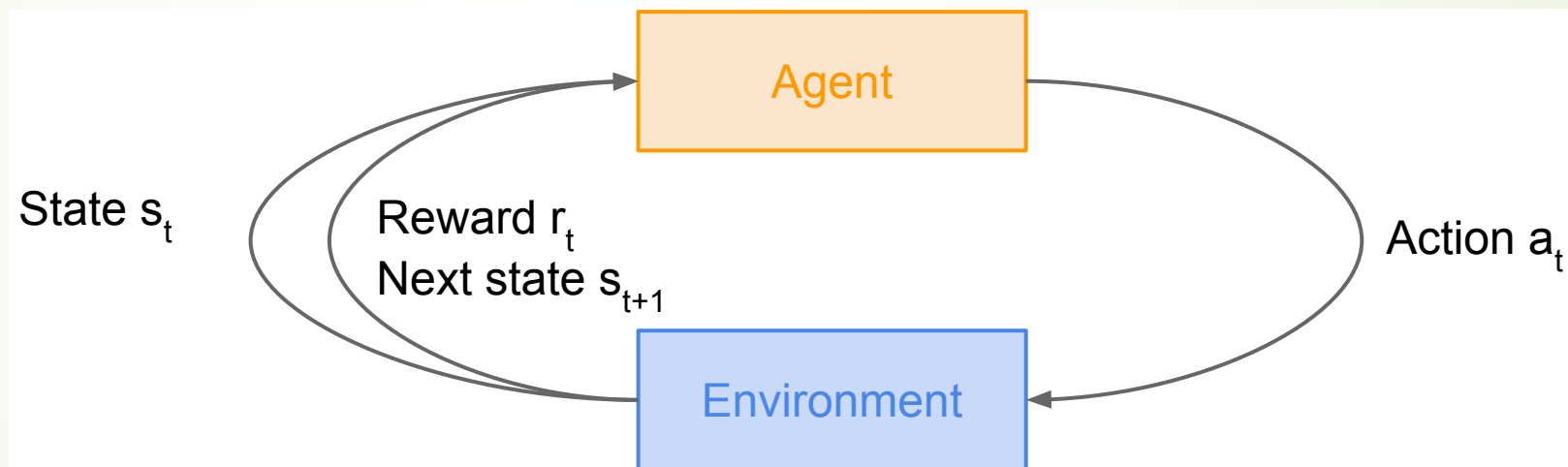
Agent

Environment

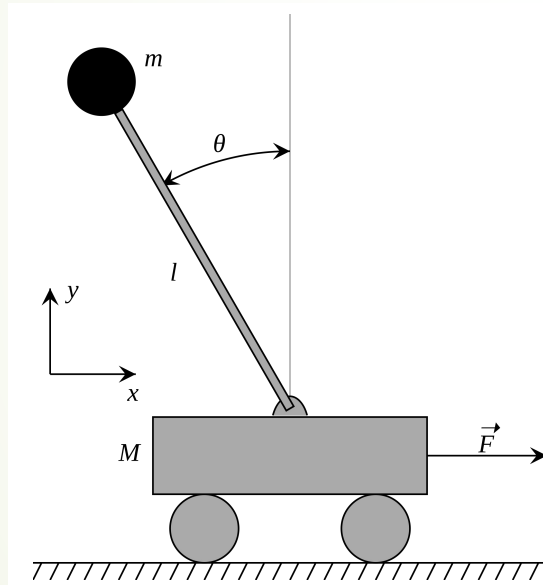








Car-Pole Control Problem

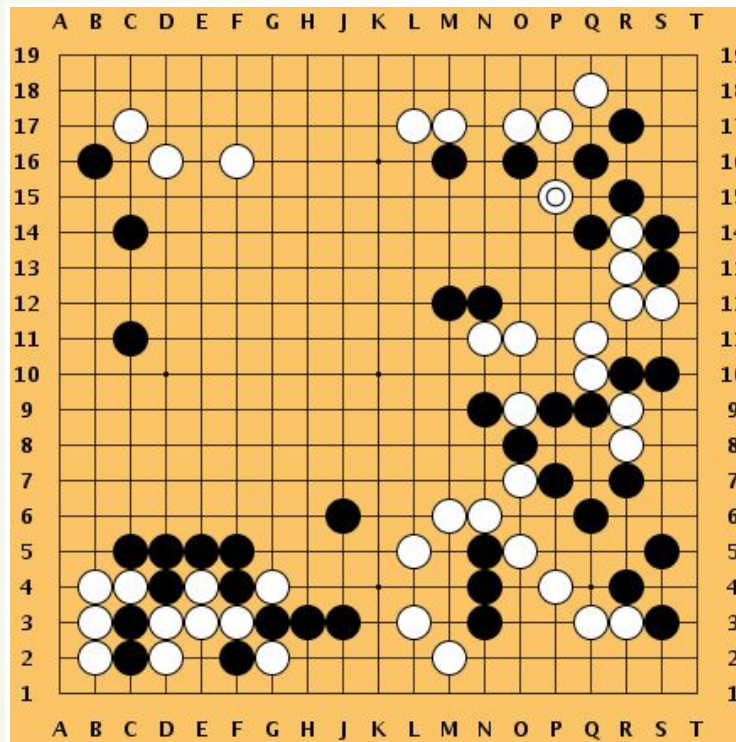


Objective: Balance a pole on top of a movable cart

State: angle, angular speed, position, horizontal velocity

Action: horizontal force applied on the cart

Reward: 1 at each time step if the pole is upright



Objective: Win the game!

State: Position of all pieces

Action: Where to put the next piece down

Reward: 1 if win at the end of the game, 0 otherwise



Mathematical Formulation of Reinforcement Learning

- **Markov property:** Current state completely characterizes the state of the world

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$


\mathcal{S} : set of possible states

\mathcal{A} : set of possible actions

\mathcal{R} : distribution of reward given (state, action) pair

\mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair

γ : discount factor

- 
- At time step $t=0$, environment samples initial state $s_0 \sim p(s_0)$
 - Then, for $t=0$ until done:
 - Agent selects action a_t
 - Environment samples reward $r_t \sim R(\cdot | s_t, a_t)$
 - Environment samples next state $s_{t+1} \sim P(\cdot | s_t; a_t)$
 - Agent receives reward r_t and next state s_{t+1}
 - A policy π is a function from S to A that specifies what action to take in each state
 - **Objective:** find policy that maximizes the cumulated discounted reward

A simple MDP: Grid World

actions = {

1. right 

2. left 

3. up 

4. down 

}

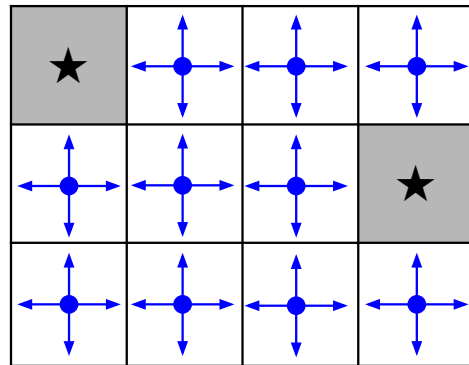
states

★			
			★

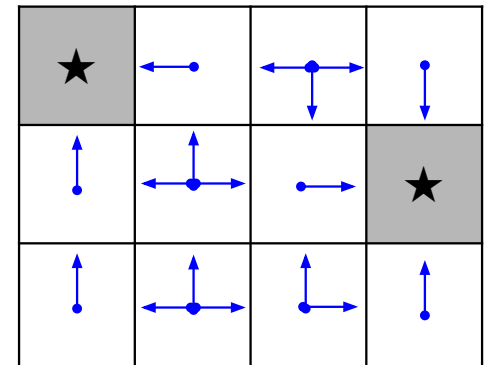
Set a negative “reward”
for each transition
(e.g. $r = -1$)

Objective: reach one of terminal states (greyed out) in
least number of actions

A simple MDP: Grid World



Random Policy



Optimal Policy

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?
Maximize the **expected sum of rewards!**

Formally: $\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi \right]$ with $s_0 \sim p(s_0)$, $a_t \sim \pi(\cdot | s_t)$, $s_{t+1} \sim p(\cdot | s_t, a_t)$

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

How good is a state?

The **value function** at state s , is the expected cumulative reward from following the policy from state s :

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$$

How good is a state-action pair?

The **Q-value function** at state s and action a , is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

Bellman Equation of Optimal Value

Optimal Value Function $V^* : \mathcal{S} \rightarrow \mathcal{R} = x^*$ satisfied the following nonlinear fixed point equation

$$x^*(i) = \max_{a \in \mathcal{A}} \left\{ r_a(i) + \gamma \sum_{j \in \mathcal{S}} P_a(i, j) x^*(j) \right\}$$

where a policy π^* is an optimal policy if and only if it attains the optimality of the Bellman equation.

Remarks

- In the continuous-time analog of MDP, i.e., stochastic optimal control, the Bellman equation is the HJB
- Exact solution methods: value iteration, policy iteration, variational analysis
- What makes things hard:

Curse of dimensionality + Modeling Uncertainty

Bellman Equation as LP (Farias and Van Roy, 2003)

The Bellman equation is equivalent to

$$\begin{aligned} & \text{minimize} && e^T x \\ & \text{subject to} && (I - \gamma P_a)x - r_a \geq 0, \quad a \in \mathcal{A}, \quad \sum_{i \in \mathcal{S}} e(i) = 1, e > 0. \end{aligned}$$

- Exact policy iteration is a form of simplex method and exhibits strongly polynomial performance (Ye 2011)
- Again, curse of dimensionality:
- Variable dimension = $|\mathcal{S}|$.
- Number of constraints = $|\mathcal{S}| \times |\mathcal{A}|$.

Duality between Value Function and Policy

Let $\lambda_{i,a} \geq 0$ be the multiplier associated with the i -th row of the primal constraint $\gamma P_a x + r_a \leq x$. The dual problem is

$$\begin{aligned} & \text{maximize} && \lambda_a^T r_a, \quad a \in \mathcal{A} \\ & \text{subject to} && \sum_{a \in \mathcal{A}} (I - \gamma P_a^T) \lambda_a = e, \quad \lambda_a \geq 0, \quad a \in \mathcal{A} \end{aligned}$$

where the dual variable is high-dimensional $\lambda = (\lambda_a)_{a \in \mathcal{A}} \in \mathbb{R}^{|\mathcal{A}||\mathcal{S}|}$.

Theorem

The optimal dual solution $\lambda^* = (\lambda_{i,a}^*)_{i \in \mathcal{S}, a \in \mathcal{A}}$ is **sparse** and has exact $|\mathcal{S}|$ nonzeros. It satisfies

$$(\lambda_{i, \mu^*(i)}^*)_{i \in \mathcal{S}} = (I - \alpha P_{\mu^*}^T)^{-1} e,$$

and $\lambda_{i,a}^* = 0$ if $a \neq \mu^*(i)$.

Finding the optimal policy $\mu^* =$ Finding the basis of the dual solution λ^*

Online Value-Policy Iteration (Mengdi Wang 2017, arXiv:1704.01869)

Stochastic primal-dual (value-policy) algorithm

- **Input:** Simulation Oracle \mathcal{M} , $n = |\mathcal{S}|$, $m = |\mathcal{A}|$, $\alpha \in (0, 1)$.
- Initialize $x^{(0)}$ and $\lambda = (\lambda_u^{(0)} : u \in \mathcal{A})$ arbitrarily.
- For $k = 1, 2, \dots, T$
 - Sample i_k uniformly from \mathcal{S} and sample u_k uniformly from \mathcal{A} .
 - **Sample next state j_k and immediate reward $g_{i_k j_k u_k}$ conditioned on (i_k, u_k) from \mathcal{M} .**
 - Update the iterates by

$$x^{(k-\frac{1}{2})} = x^{(k-1)} - \gamma_k \left(-e + m\lambda_{u_k}^{(k-1)} - \alpha mn \left(\lambda_{u_k}^{(k-1)} \cdot e_{i_k} \right) e_{j_k} \right),$$

$$\lambda_{u_k}^{(k-\frac{1}{2})} = \lambda_{u_k}^{(k-1)} + m\gamma_k \left(x^{(k-1)} - \alpha n \left(x^{(k-1)} \cdot e_{j_k} \right) e_{i_k} - n g_{i_k j_k u_k} e_{i_k} \right),$$

$$\lambda_u^{(k-\frac{1}{2})} = \lambda_u^{(k-1)}, \quad \forall u \neq u_k,$$

- Project the iterates orthogonally to some regularization constraints

$$x^{(k)} = \Pi_X x^{(k-\frac{1}{2})}, \quad \lambda^{(k)} = \Pi_\Lambda \lambda^{(k-\frac{1}{2})}.$$

- **Output:** Averaged dual iterate $\hat{\lambda} = \frac{1}{T} \sum_{k=1}^T \lambda^{(k)}$

Near Optimal Primal-Dual Algorithms

Method	Setting	Sample Complexity	Run-Time Complexity	Space Complexity	Reference
Phased Q-Learning	γ discount factor, ϵ -optimal value	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3\epsilon^2} \ln \frac{1}{\delta}$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3\epsilon^2} \ln \frac{1}{\delta}$	$ \mathcal{S} \mathcal{A} $	[17]
Model-Based Q-Learning	γ discount factor, ϵ -optimal value	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3\epsilon^2} \ln \frac{ \mathcal{S} \mathcal{A} }{\delta}$	NA	$ \mathcal{S} ^2 \mathcal{A} $	[1]
Randomized P-D	γ discount factor, ϵ -optimal policy	$\frac{ \mathcal{S} ^3 \mathcal{A} }{(1-\gamma)^6\epsilon^2}$	$\frac{ \mathcal{S} ^3 \mathcal{A} }{(1-\gamma)^6\epsilon^2}$	$ \mathcal{S} \mathcal{A} $	[25]
Randomized P-D	γ discount factor, τ -stationary, ϵ -optimal policy	$\tau^4 \frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^4\epsilon^2}$	$\tau^4 \frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^4\epsilon^2}$	$ \mathcal{S} \mathcal{A} $	[25]
Randomized VI	γ discount factor, ϵ -optimal policy	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^4\epsilon^2}$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^4\epsilon^2}$	$ \mathcal{S} \mathcal{A} $	[23]
Primal-Dual π Learning	τ -stationary, t_{mix}^* -mixing, ϵ -optimal policy	$\frac{(\tau \cdot t_{mix}^*)^2 \mathcal{S} \mathcal{A} }{\epsilon^2}$	$\frac{(\tau \cdot t_{mix}^*)^2 \mathcal{S} \mathcal{A} }{\epsilon^2}$	$ \mathcal{S} \mathcal{A} $	This Paper

Table 1: Complexity Results for Sampling-Based Methods for MDP. The sample complexity is measured by the number of queries to the \mathcal{SO} . The run-time complexity is measured by the total run-time complexity under the assumption that each query takes $\tilde{\mathcal{O}}(1)$ time. The space complexity is the additional space needed by the algorithm in addition to the input.

Q-Learning

Bellman equation

The optimal Q-value function Q^* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

Q^* satisfies the following **Bellman equation**:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

Intuition: if the optimal state-action values for the next time-step $Q^*(s', a')$ are known, then the optimal strategy is to take the action that maximizes the expected value of $r + \gamma Q^*(s', a')$

The optimal policy π^* corresponds to taking the best action in any state as specified by Q^*



Solving for the optimal policy

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} \left[r + \gamma \max_{a'} Q_i(s', a') | s, a \right]$$

Q_i will converge to Q^* as $i \rightarrow \text{infinity}$

Solving for the optimal policy

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} \left[r + \gamma \max_{a'} Q_i(s', a') | s, a \right]$$

Q_i will converge to Q^* as $i \rightarrow \infty$

What's the problem with this?

Not scalable. Must compute $Q(s,a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Solution: use a function approximator to estimate $Q(s,a)$. E.g. a neural network!



Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

If the function approximator is a deep neural network => **deep q-learning!**

Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

Forward Pass

Loss function: $L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$

Backward Pass

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right] \nabla_{\theta_i} Q(s, a; \theta_i)$$

[Mnih et al. NIPS Workshop 2013; Nature 2015]

Case Study: Playing Atari Games



Objective: Complete the game with the highest score

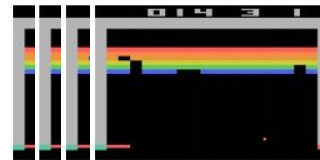
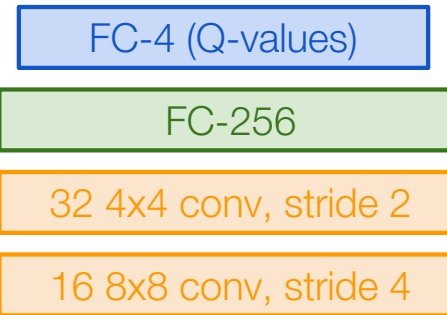
State: Raw pixel inputs of the game state

Action: Game controls e.g. Left, Right, Up, Down

Reward: Score increase/decrease at each time step

Q-network Architecture

$Q(s, a; \theta)$:
neural network
with weights θ



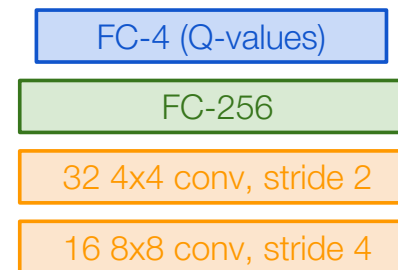
Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

[Mnih et al. NIPS Workshop 2013; Nature 2015]

Q-network Architecture

$Q(s, a; \theta)$:
neural network
with weights θ

A single feedforward pass
to compute Q-values for all
actions from the current
state => efficient!



Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

← Last FC layer has 4-d
output (if 4 actions),
corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2)$, $Q(s_t, a_3)$,
 $Q(s_t, a_4)$

Number of actions between 4-18
depending on Atari game

Training the Q-network: Loss function (from before)

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

Forward Pass

Loss function: $L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} \left[(y_i - Q(s, a; \theta_i))^2 \right]$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$

Iteratively try to make the Q-value close to the target value (y_i) it should have, if Q-function corresponds to optimal Q^* (and optimal policy π^*)

Backward Pass

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right] \nabla_{\theta_i} Q(s, a; \theta_i)$$

Training the Q-network: Experience Replay

- ▶ Learning from batches of consecutive samples is problematic:
 - ▶ Samples are correlated => inefficient learning
 - ▶ Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops
- ▶ Address these problems using **experience replay**
 - ▶ Continually update a **replay memory** table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
 - ▶ Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

Each transition can also contribute to multiple weight updates
=> greater data efficiency

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for



Example

- ▶ Google DeepMind's Deep Q-learning playing Atari Breakout:
 - ▶ <https://www.youtube.com/watch?v=V1eYniJ0Rnk>
 - ▶ Google DeepMind created an artificial intelligence program using deep reinforcement learning that plays Atari games and improves itself to a superhuman level. It is capable of playing many Atari games and uses a combination of deep artificial neural networks and reinforcement learning. After presenting their initial results with the algorithm, Google almost immediately acquired the company for several hundred million dollars, hence the name Google DeepMind. Please enjoy the footage and let me know if you have any questions regarding deep learning!



Policy Gradients

- ▶ What is a problem with Q-learning?
The Q-function can be very complicated!
- ▶ Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair
- ▶ But the policy can be much simpler: just close your hand
Can we learn a policy directly, e.g. finding the best policy from a collection of policies?



Policy Gradients

Formally, let's define a class of parametrized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta \right]$$

We want to find the optimal policy $\theta^* = \arg \max_{\theta} J(\theta)$

How can we do this?

Gradient ascent on policy parameters!



REINFORCE algorithm

Mathematically, we can write:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] \\ &= \int_{\tau} r(\tau) p(\tau; \theta) d\tau \end{aligned}$$

Where $r(\tau)$ is the reward of a trajectory $\tau = (s_0, a_0, r_0, s_1, \dots)$

Expected reward:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$$
$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

Now let's differentiate this: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$ **Intractable! Gradient of an expectation is problematic when p depends on θ**

However, we can use a nice trick: $\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$
If we inject this back:

$$\nabla_{\theta} J(\theta) = \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau$$
$$= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)]$$

Can estimate with Monte Carlo sampling

REINFORCE algorithm

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau \\ &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)]\end{aligned}$$

Can we compute those quantities without knowing the transition probabilities?

We have: $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$

Thus: $\log p(\tau; \theta) = \sum_{t \geq 0} \log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t)$

And when differentiating: $\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

Doesn't depend on transition probabilities!

Therefore when sampling a trajectory τ , we can estimate $J(\theta)$ with

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Intuition

Gradient estimator:
$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. **But in expectation, it averages out!**

However, this also suffers from high variance because credit assignment is really hard. Can we help the estimator?

Variance reduction

Gradient estimator: $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

First idea: Push up probabilities of an action seen, only by the cumulative future reward from that state

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Second idea: Use discount factor γ to ignore delayed effects

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Variance reduction: Baseline

Problem: The raw value of a trajectory isn't necessarily meaningful. For example, if rewards are all positive, you keep pushing up probabilities of actions.

What is important then? Whether a reward is better or worse than what you expect to get

Idea: Introduce a baseline function dependent on the state. Concretely, estimator is now:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

How to choose the baseline?

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

A simple baseline: constant moving average of rewards experienced so far from all trajectories

Variance reduction techniques seen so far are typically used in “Vanilla REINFORCE”

How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action a_t in a state s_t if $Q^\pi(s_t, a_t) - V^\pi(s_t)$ is large. On the contrary, we are unhappy with an action if it's small.

Using this, we get the estimator:
$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)) \nabla_\theta \log \pi_\theta(a_t | s_t)$$

Actor-Critic Algorithm

Problem: we don't know Q and V. Can we learn them?

Yes, using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q-function).

- The actor decides which action to take, and the critic tells the actor how good its action was and how it should adjust
- Also alleviates the task of the critic as it only has to learn the values of (state, action) pairs generated by the policy
- Can also incorporate Q-learning tricks e.g. experience replay
- **Remark:** we can define by the **advantage function** how much an action was better than expected

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

Actor-Critic Algorithm

Initialize policy parameters θ , critic parameters ϕ

For iteration=1, 2 ... **do**

Sample m trajectories under the current policy

$\Delta\theta \leftarrow 0$

For $i=1, \dots, m$ **do**

For $t=1, \dots, T$ **do**

$$A_t = \sum_{t' \geq t} \gamma^{t'-t} r_{t'}^i - V_\phi(s_t^i)$$

$$\Delta\theta \leftarrow \Delta\theta + A_t \nabla_\theta \log(a_t^i | s_t^i)$$

$$\Delta\phi \leftarrow \sum_i \sum_t \nabla_\phi \|A_t^i\|^2$$

$$\theta \leftarrow \alpha \Delta\theta$$

$$\phi \leftarrow \beta \Delta\phi$$

End for

REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

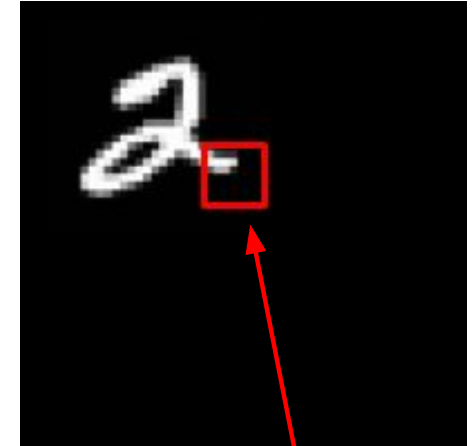
Take a sequence of “glimpses” selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

Reward: 1 at the final timestep if image correctly classified, 0 otherwise

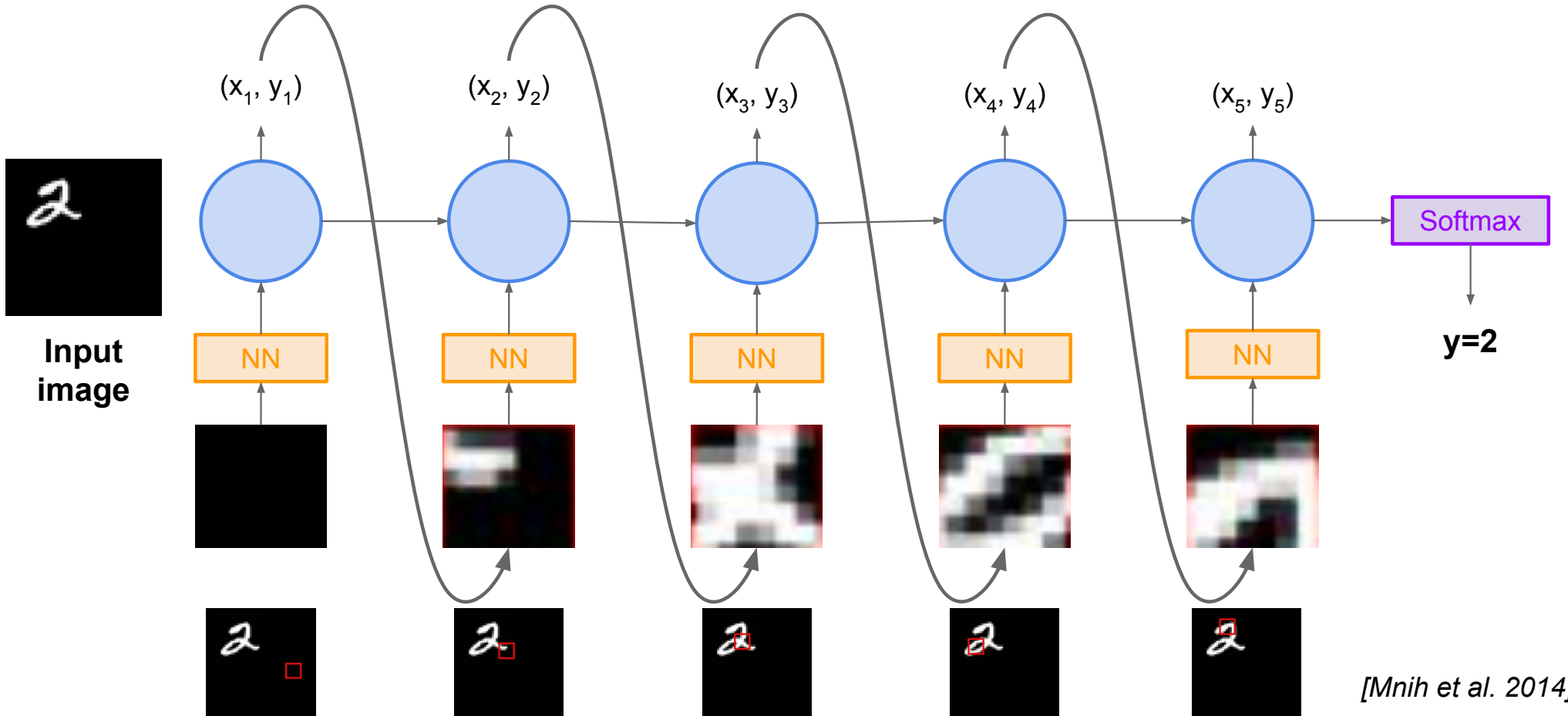


glimpse

Glimpsing is a non-differentiable operation => learn policy for how to take glimpse actions using REINFORCE
Given state of glimpses seen so far, use RNN to model the state and output next action

[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



Pytorch Implementation

- <https://github.com/kevinzakka/recurrent-visual-attention>
- A Pytorch implementation for the paper, [Recurrent Models of Visual Attention](#) by Volodymyr Mnih, Nicolas Heess, Alex Graves and Koray Kavukcuoglu, NIPS 2014.



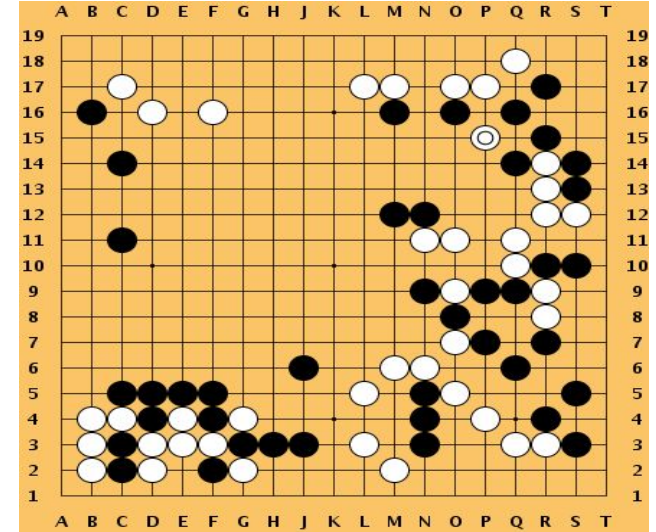
More policy gradients: AlphaGo

Overview:

- Mix of supervised learning and reinforcement learning
- Mix of old methods (Monte Carlo Tree Search) and recent ones (deep RL)

How to beat the Go world champion:

- Featurize the board (stone color, move legality, bias, ...)
- Initialize policy network with supervised training from professional go games, then continue training using policy gradient (play against itself from random previous iterations, +1 / -1 reward for winning / losing)
- Also learn value network (critic)
- Finally, combine combine policy and value networks in a Monte Carlo Tree Search algorithm to select actions by lookahead search



[Silver et al.,
Nature 2016]

This image is CC0 public domain



Summary



- ▶ **Policy gradients:** very general but suffer from high variance so requires a lot of samples. **Challenge:** sample-efficiency
- ▶ **Q-learning:** does not always work but when it works, usually more sample-efficient. **Challenge:** exploration
- ▶ Guarantees:
 - ▶ **Policy Gradients:** Converges to a local minima, often good enough!
 - ▶ **Q-learning:** Zero guarantees since you are approximating Bellman equation with a complicated function approximator

Optimized Execution, Market Microstructure and Reinforcement Learning

[Y. Nevmyvaka, Y. Feng, MK; ICML 2006]

[MK, Y. Nevmyvaka; In "High Frequency Trading", O'Hara et al. eds, Risk Books 2013]



Michael Kearns, University of Pennsylvania, ICML 2014, Beijing

A Brief Field Guide to Wall Street

- “Buy Side”: Attempt to outperform market via proprietary research
 - Includes hedge funds, mutual funds, statistical arbitrage, HFT, prop trading groups
 - May or may not be quantitative and automated
 - Have investors but not clients
 - Take and hold positions → risk
 - Generation of “alpha” still more art than science
- “Sell Side”: Provide brokerage and execution services
 - Includes bank and independent brokerages, exchanges
 - Almost entirely quantitative and automated
 - Clients are the buy side
 - Do not hold risk; paid via fees/commissions/etc.
- In reality, alpha and execution are blurred
 - Especially at shorter holding periods (e.g. HFT)

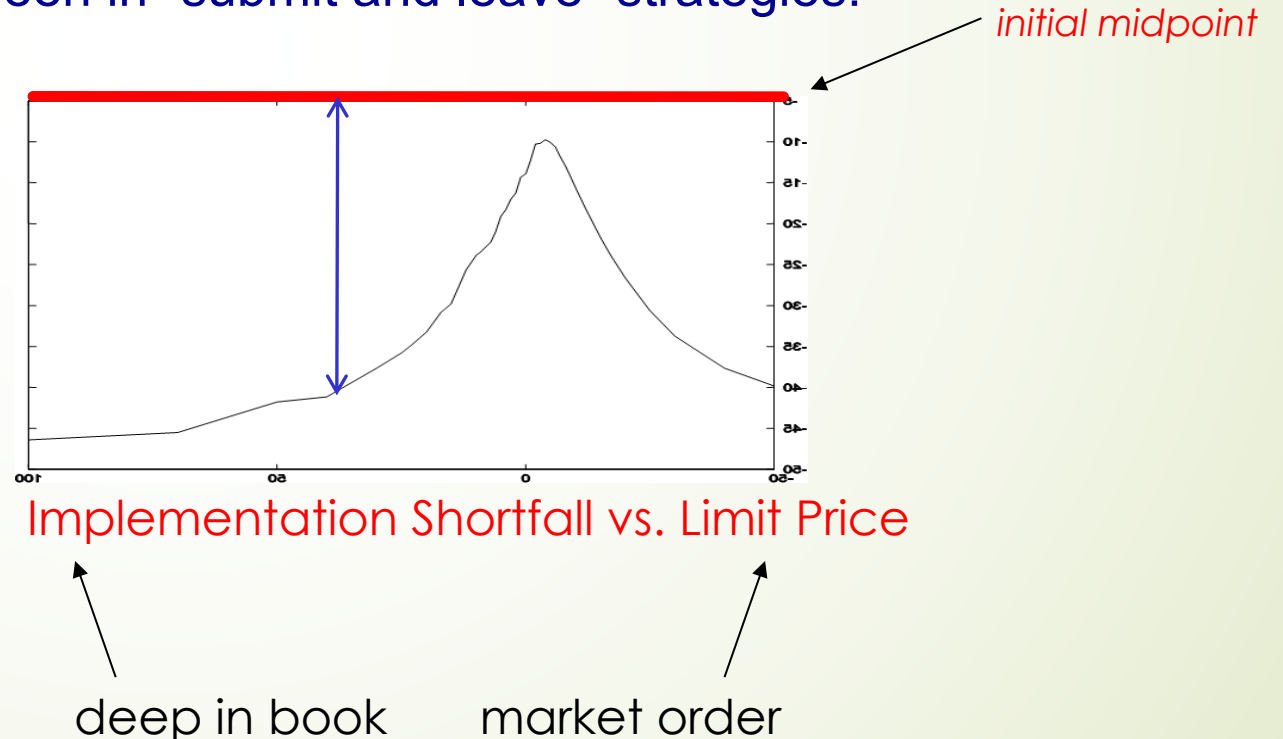
A Canonical Trading Problem

- Goal (buy side to sell side): Sell V shares in T time steps; maximize revenue
- Strategy Evaluation Metric Benchmarks:
 - Volume Weighted Average Price (VWAP)
 - Time Weighted Average Price (TWAP)
 - Implementation Shortfall (midpoint of bid-ask spread at beginning)
- Natural to view as a problem of *state-based control (RL)*
 - State variables: inventory V and time remaining T (discretized)
 - Features capturing market activity?

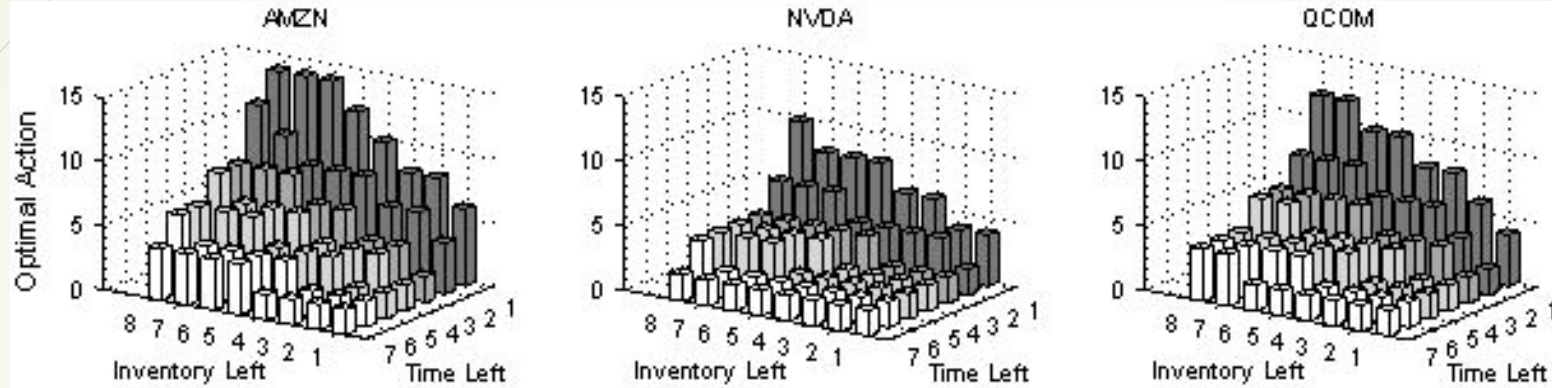
Market Microstructure

refresh island home disclaimer help			
GET STOCK			
MSFT		MSFT	go
Symbol Search			
LAST MATCH		TODAY'S ACTIVITY	
Price	23.7790	Orders	1,630
Time	9:01:55.614	Volume	44,839
BUY ORDERS		SELL ORDERS	
SHARES	PRICE	SHARES	PRICE
1,000	23.7600	100	23.7800
3,087	23.7500	800	23.7990
200	23.7500	500	23.8000
100	23.7400	1,720	23.8070
1,720	23.7280	900	23.8190
2,000	23.7200	200	23.8500
1,000	23.7000	1,000	23.8500
100	23.7000	1,000	23.8500
100	23.7000	1,000	23.8600
800	23.6970	200	24.0000
500	23.6500	500	24.0000
3,000	23.6500	1,000	24.0300
4,300	23.6500	200	24.0300
2,000	23.6500	1,100	24.0400
200	23.6200	500	24.0500
(195 more)		(219 more)	

- Continuous double auction with limit orders: buy orders decreasing; sell orders increasing
- Volatile and dynamic; sub-millisecond time scale
- Cancellations, revisions, partial executions
- How do individual orders (micro) influence aggregate market behavior (macro)?
- Tradeoff between *immediacy* and *price*
- Seen in “submit and leave” strategies:



Policies Learned: Time and Volume Remaining



- Experimental framework
 - Full historical order book reconstruction and simulation
 - Learn optimal policy on 1 year training; test on following 6 months
 - Pitfalls: directional drift, “counterfactual” market impact
- Overall shape is consistent and sensible
 - Become more aggressive (spread crossing) as time runs out or inventory is too large
 - Learning optimizes this qualitative schedule

Additional Improvement From Order Book Features

Bid Volume	-0.06%	Ask Volume	-0.28%
Bid-Ask Volume Misbalance	0.13%	Bid-Ask Spread	7.97%
Price Level	0.26%	Immediate Market Order Cost	4.26%
Signed Transaction Volume	2.81%	Price Volatility	-0.55%
Spread Volatility	1.89%	Signed Incoming Volume	0.59%
Spread + Immediate Cost	8.69%	Spread+ImmCost+Signed Vol	12.85%

Some Idealized Trading Scenarios and Risks

- Assume all the transactions cross the bid/ask spread at approximate midpoint (median) price
 - Example: $V=\{1,0,-1\}$ (long/nothing/short), $T=1$ min
- *Return* maximization with **no-regret** sequential (online) strategies:
 - Compete with best single strategy in hindsight
 - Unfortunately methods work poorly in practice
- Could ask for no-regret to best strategy in *risk-adjusted metrics*:
 - Sharpe Ratio: $\mu(\text{returns})/\sigma(\text{returns})$
 - Mean-Variance: $\mu(\text{returns}) - \sigma(\text{returns})$
- Yet strong negative results in risk-adjusted metrics:
 - No-regret provably impossible
 - $1 + \epsilon$ lower bound on competitive ratio
- Intuition: Volatility terms σ introduce additional costs that one has to pay
- Loss design should incorporate risk measurements, or internalize risks in strategies



Online Tutorials

- ▶ A GitHub repo for *deep reinforcement learning strategies and environments for quantitative trading*
 - ▶ <https://github.com/Ceruleanacg/Personae/blob/master/README.md>
 - ▶ This is a good start for the application of deep reinforcement learning in algorithmic trading
 - ▶ Can you **reproduce** the results there?
- ▶ next week, we shall have a tutorial about Reinforcement Learning for Quantitative Finance.

Thank you!

