

# APPROXIMATE Q-LEARNING

Scott O'Hara

Metrowest Developers Machine Learning Group

06/10/2020

## REFERENCES

**CS188 - Introduction to Artificial Intelligence** course at University of California, Berkeley:

- ▶ <http://ai.berkeley.edu/home.html>
- ▶ <http://gamescrafters.berkeley.edu/~cs188/sp20/>

**CS181 - Machine Learning** course at Harvard University:

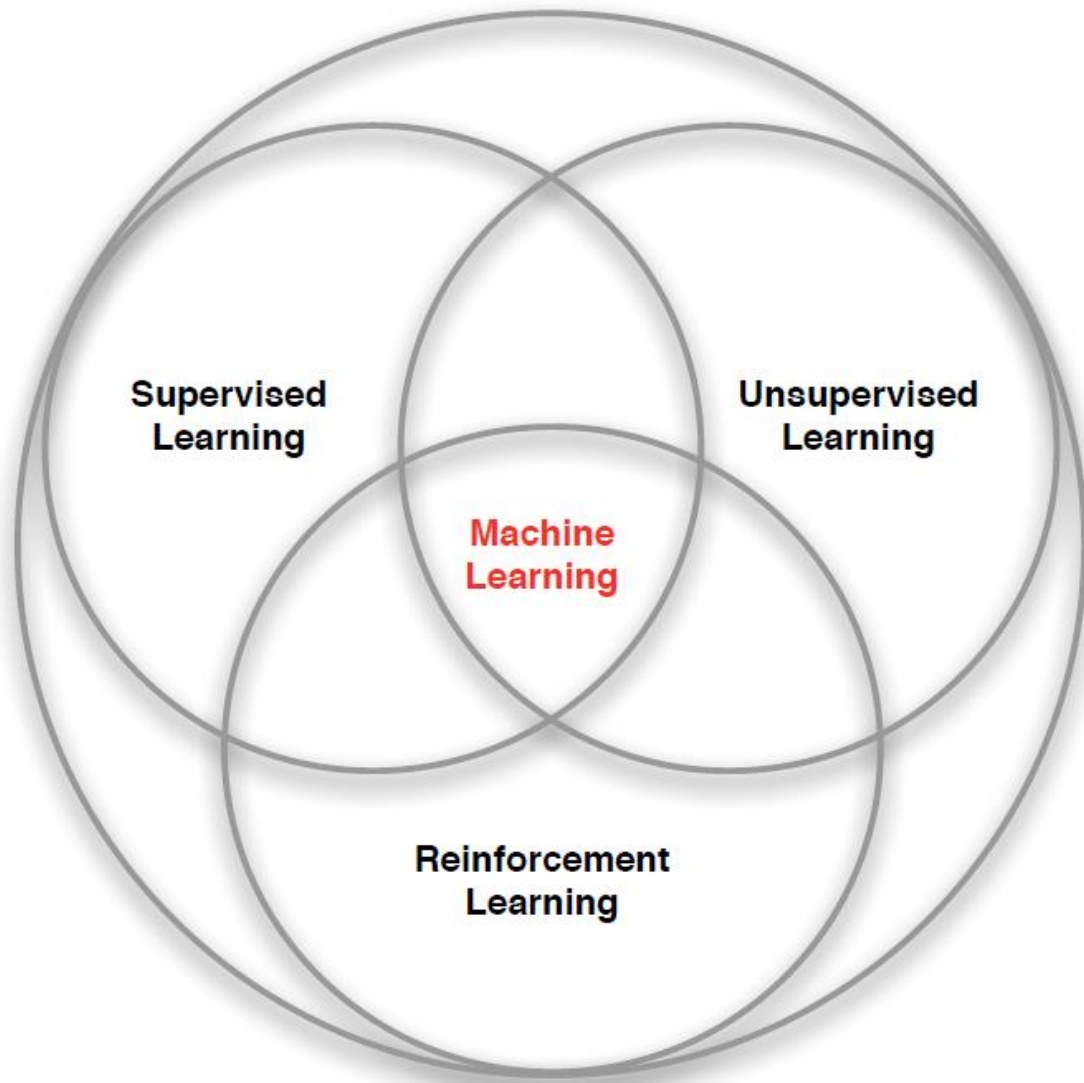
- ▶ *Lectures and notes from multiple offerings of CS181: Spring 2009, 2011, 2014 and 2017.*

***The Hundred-Page Machine Learning Book***, A. Burkov,  
Publisher: Andriy Burkov, 2019. <http://themlbook.com/>.

***Reinforcement learning: an Introduction***, R. S. Sutton and A. G. Barto, Second edition. Cambridge, Massachusetts: The MIT Press, 2018.

# REINFORCEMENT LEARNING





# 3 BRANCHES OF MACHINE LEARNING

**Credit:** adapted from lecture slides David Silver, DeepMind, "Introduction to Reinforcement Learning"

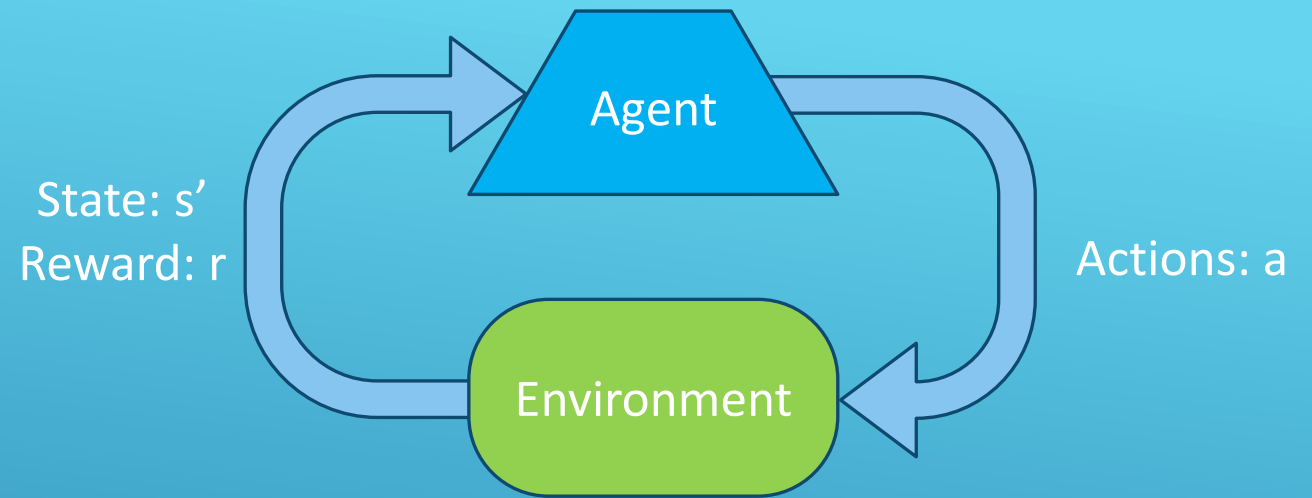
# 1. SUPERVISED LEARNING

- **Labeled Data** =  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  with **feature vector**  $x_i$  and **label**  $y_i$ .
- **Train** a **model**  $f$  on the labeled data.
- **Predict** label  $y = f(x)$  for unlabeled data.
- Examples: linear regression, decision trees, SVMs, k-nearest neighbors.

## 2. UNSUPERVISED LEARNING

- *Unlabeled Data* =  $\{x_1, \dots, x_n\}$  of **feature vectors**  $x_i$ .
- Create a **model**  $f$  on the unlabeled data.
- **Transform** feature vector  $x$  into:
  - **ID**  $y = f(x)$  (**Clustering**)
  - **Vector**  $x' = f(x)$  (**Dimensionality Reduction**)
- Clustering examples: K-means, HAC
- Dimensionality reduction examples: PCA, UMAP

### 3. REINFORCEMENT LEARNING



- Agent “lives” in an environment and can perceive the **state** of that environment as a **vector of features**.
- Agent learns to act to **maximize the expected average reward**
- Agent knows the current state  $s$ , takes an action  $a$ , receives a reward  $r$  and observes the next state  $s'$ .

$$S_0, A_0, R_0, S_1, A_1, R_1, S_2, A_2, R_2, \dots, S_n, A_n, R_n, S_T$$

# MARKOV DECISION PROCESSES

- **States:**  $s_1, \dots, s_n$
- **Actions:**  $a_1, \dots, a_m$

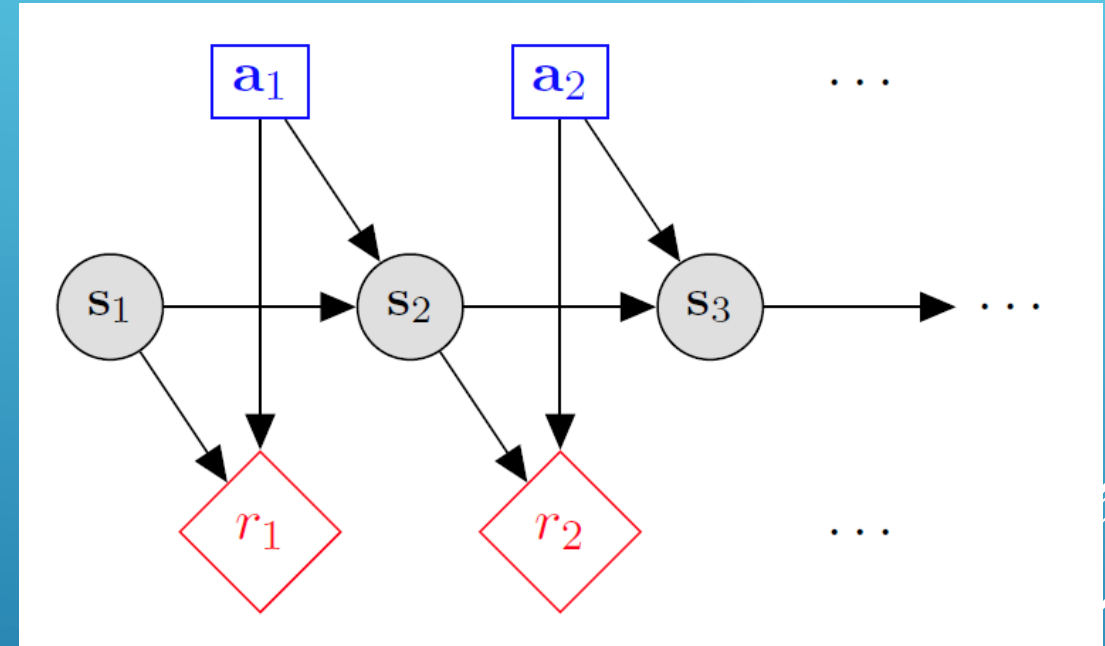
- **Reward model:**

$$R(s, a, s') \in R$$

- **Transition model:**

$$T(s, a, s') = P(s'|s, a)$$

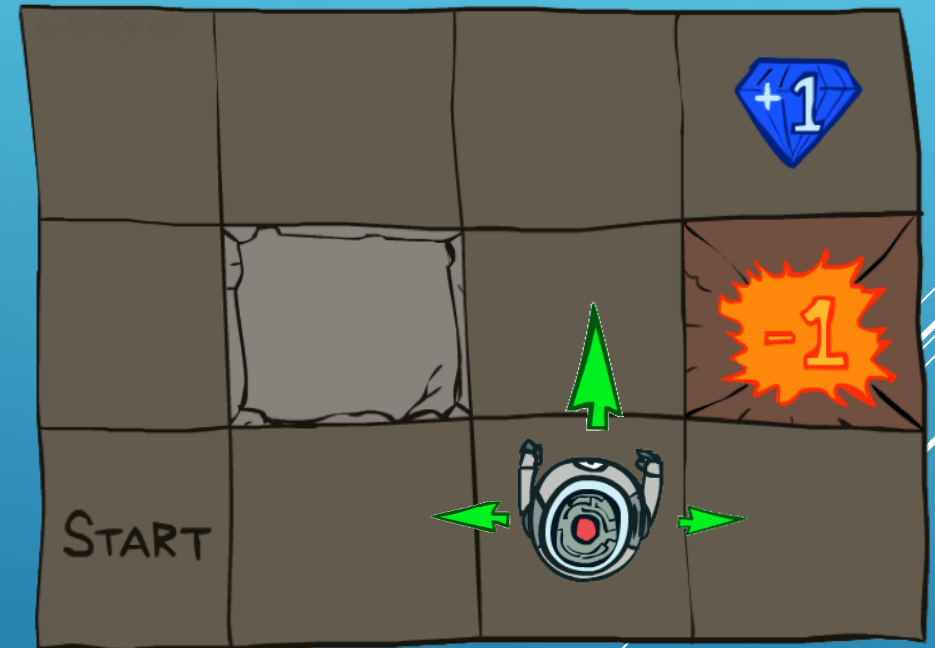
- **Discount factor:**  $\gamma \in [0, 1]$





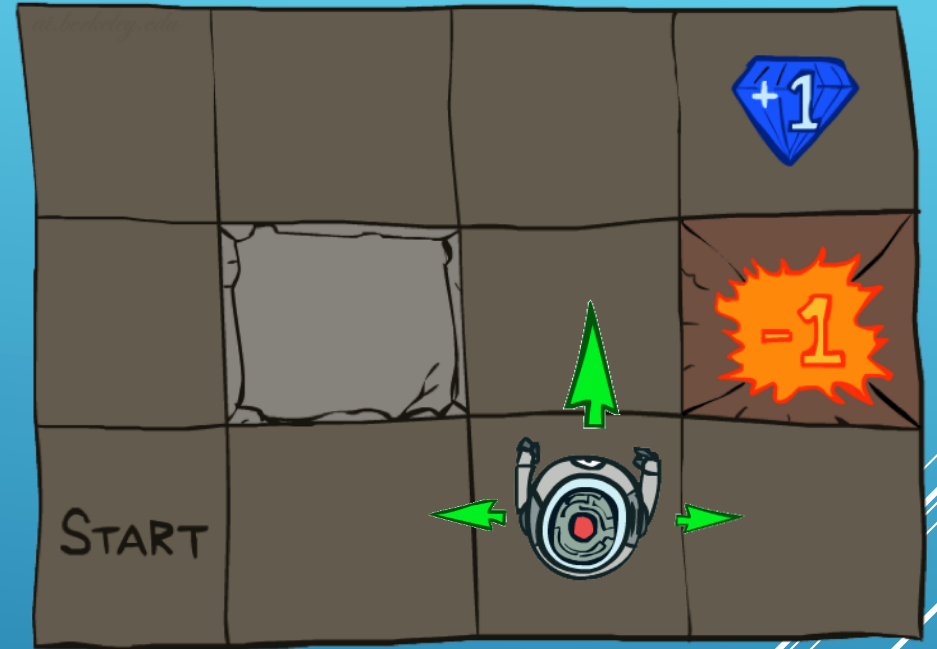
# EXAMPLE: GRIDWORLD

- A maze-like problem
  - The agent lives in a grid
  - Walls block the agent's path
- Noisy movement: actions do not always go as planned
  - 80% of the time, the action North takes the agent North (if there is no wall there)
  - 10% of the time, North takes the agent West; 10% East
  - If there is a wall in the direction the agent would have been taken, the agent stays put
- The agent receives rewards each time step
  - Small "living" reward each step (usually negative)
  - Big rewards come at the end (good or bad)
- Goal: maximize sum of rewards



# EXAMPLE: GRIDWORLD

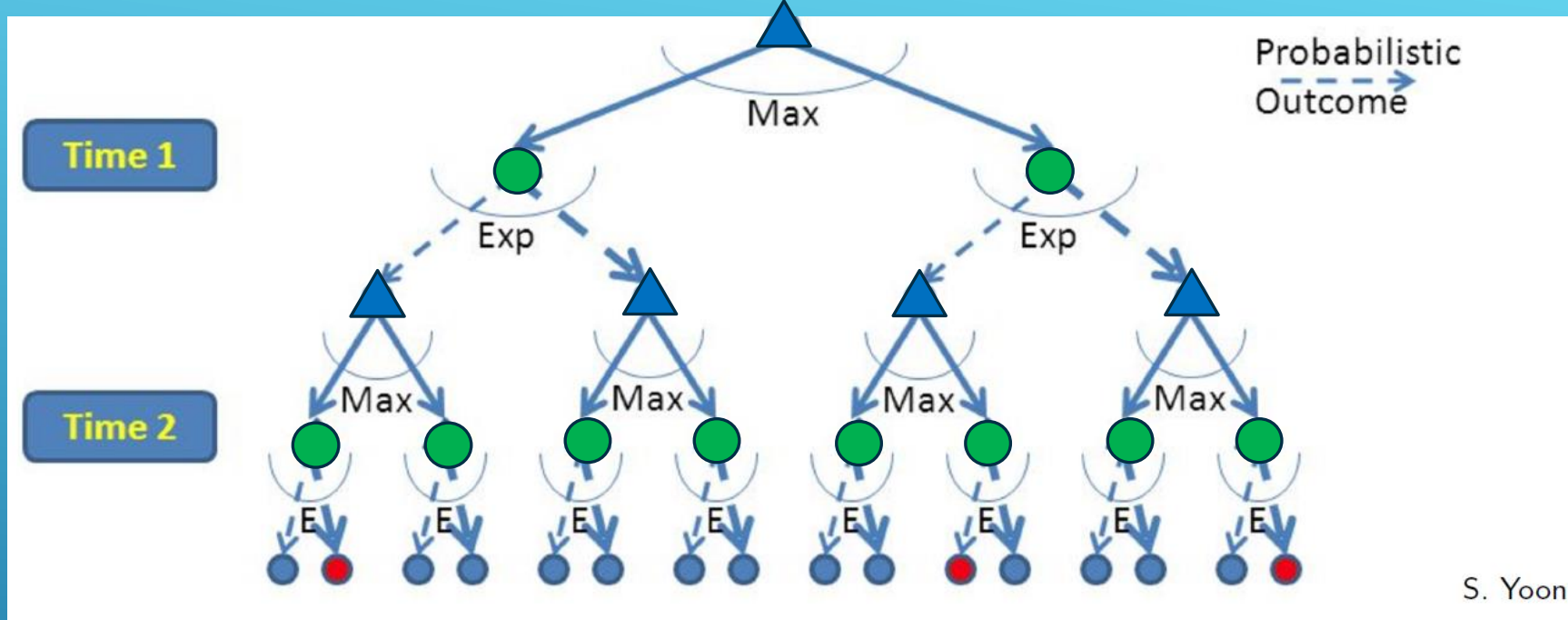
- **States:**  $\langle x, y \rangle$  locations
- **Actions:** move **north**, **south**, **east** or **west**.
- **Reward model:**
  - +1 if robot moves to  $\langle 3, 2 \rangle$
  - 1 if robot moves  $\langle 3, 1 \rangle$
  - otherwise, get “living reward”.
- **Transition model:**
  - $n$  – probability of unintended (noisy) action.
  - $1-n$  probability of intended action.
  - stay put if you move into a wall.



# TABULAR Q-LEARNING



# RL IS LIKE A GAME AGAINST NATURE



- Reinforcement learning is like a game-playing algorithm.
- Nodes where you move are called **states**:  $S$  ( $\triangle$ )
- Nodes where nature “moves” are called **Q-states**:  $\langle S, A \rangle$  ( $\bullet$ )

# QUANTITIES TO OPTIMIZE

- The **state-value function**:

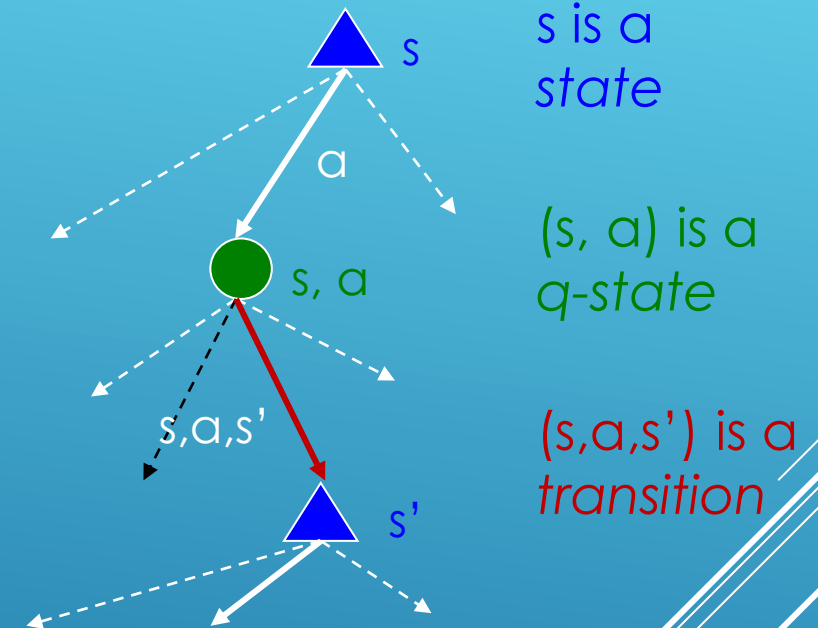
$V(s)$  = the expected utility of starting in state  $s$  and acting optimally afterwards.

- The **action-value function**:

$Q(s,a)$  = the expected utility of starting in state  $s$ , taking action  $a$ , and acting optimally afterwards.

- The **policy**:

$\pi(s)$  = defines the action to take in every state  $s$ .



# THE BELLMAN OPTIMALITY EQUATIONS

- ▶ The Bellman Optimality Equations define a relationship that, when satisfied, guarantees that the *state-value function*  $V^*(s)$  and the *state-action function*  $Q^*(s, a)$  are optimal for every state and action.
- ▶ This in turn guarantees that the policy  $\pi$  is optimal, which is designated  $\pi^*$ .

# THE BELLMAN OPTIMALITY EQUATIONS

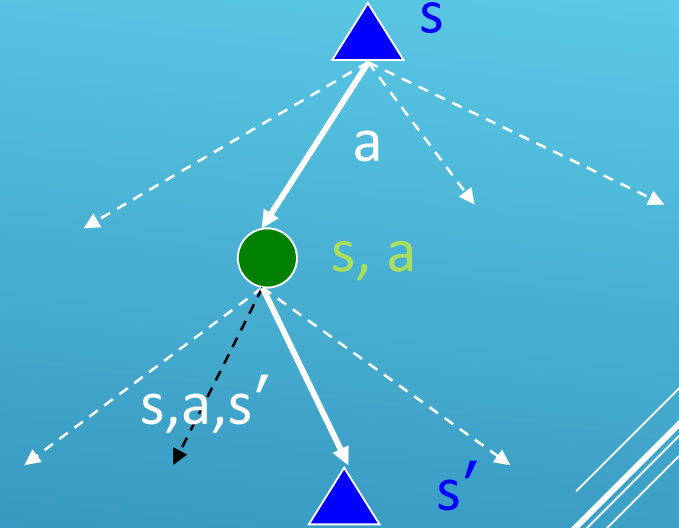
## State-Value Equation:

$$\boxed{V^*(s)} = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \boxed{V^*(s')}]$$

## Action-Value Equation:

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$

Simplifying Assumption: rewards are fixed for  $(s, a, s')$



# MUTUAL RECURSION

State-Value Equation:

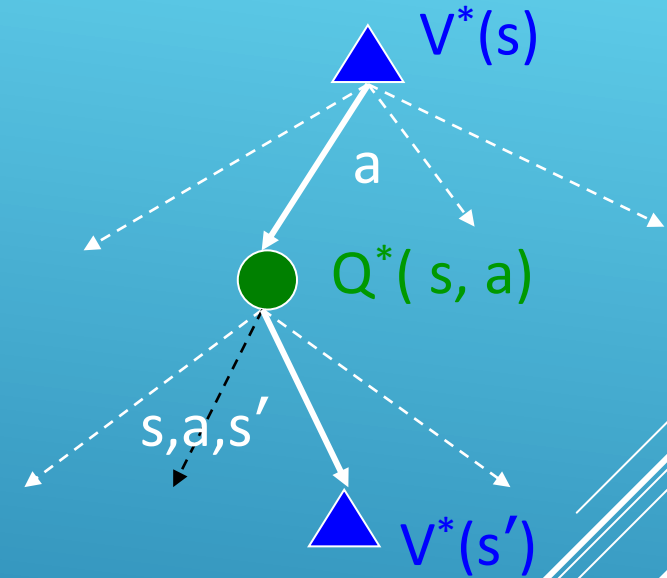
$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a Q^*(s, a)$$

Action-Value Equation:

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

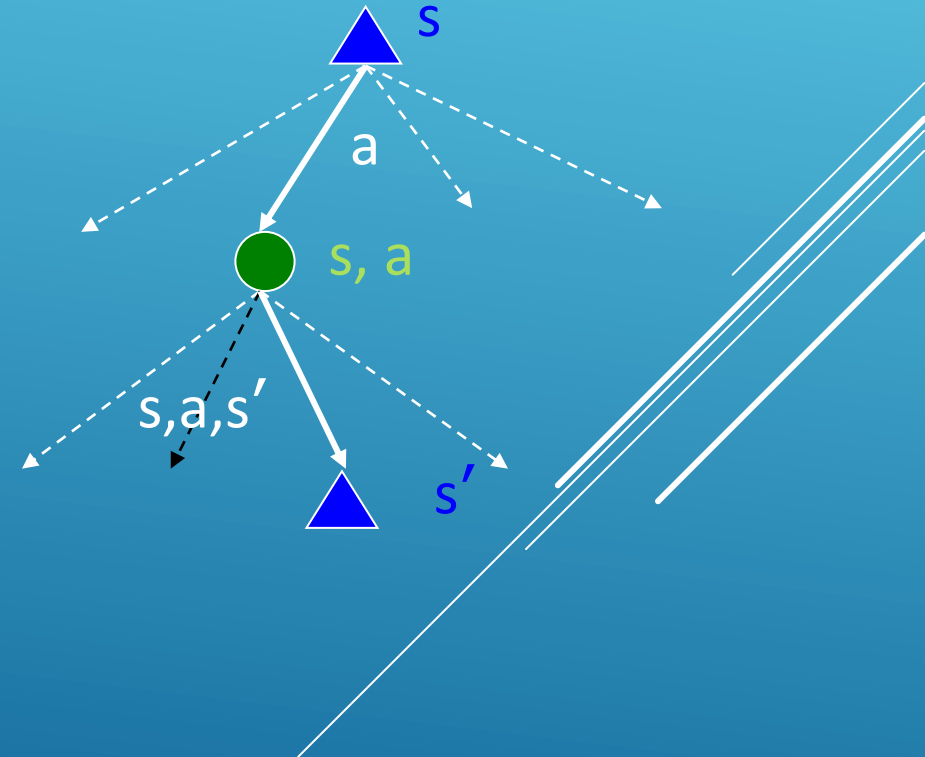




# THE OPTIMAL STATE-VALUE EQUATION $V^*$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

- ▶ Focusing on different Bellman Equations gives different algorithms
- ▶ The  $V^*$  equation gives rise to these dynamic programming algorithms previously discussed:
  - ▶ Value Iteration
  - ▶ Policy Iteration



# THE OPTIMAL VALUE UTILITY EQUATION $Q^*$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$

The  $Q^*$  equation gives rise to the Q-Learning algorithm.

# FROM EQUATION TO UPDATE RULE

- ▶ What to do about  $T(s,a,s')$  and  $R(s,a,s')$ , since we don't have these functions?

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$

- ▶ Use sampling to learn  $Q(s,a)$  values as you go
  - ▶ Receive a sample transition:  $(s, a, r, s')$
  - ▶ Consider your old estimate:  $Q(s, a)$
  - ▶ Consider your new sample estimate:  $r + \gamma \max_{a'} Q_k(s', a')$
  - ▶ Incorporate the new estimate into a running average based on the learning rate  $\alpha$ :

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q_k(s', a') - Q(s, a) \right]$$

## Q-LEARNING UPDATE RULE (2)

- ▶ On transitioning from state  $s$  to state  $s'$  on action  $a$ , and receiving reward  $r$ , update:

$$Q_{k+1}(s, a) \leftarrow Q_k(s, a) + \alpha \left[ \overset{\text{target}}{r + \gamma \max_{a'} Q_k(s', a')} - \overset{\text{current}}{Q_k(s, a)} \right]$$

- ▶  $\alpha$  is the **learning rate**
  - ▶ A large  $\alpha$  results in quicker learning but may not converge.
  - ▶  $\alpha$  is often decreased as learning goes on.
- ▶  $\gamma$  is the **discount rate** i.e., discounts future rewards.

# CHOOSING AN ACTION: EXPLORATION VS EXPLOITATION

- How should an agent choose an action? An obvious answer is simply to follow the current policy. However, this is often not the best way to improve your model.
- **Exploit:** use your current model to maximize the expected utility now.
- **Explore:** choose an action that will help you improve your model.

# E-GREEDY METHOD

- With probability  $1 - \epsilon$ :

$A_t = \operatorname{argmax} Q_t(s, a)$ , select action with maximum value.

- With probability  $\epsilon$ :

$A_t =$  select with equal probability an action at state  $s$  from all possible actions.

# TABULAR Q-LEARNING ALGORITHM

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -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

# TRANSITION/REWARD PARAMETER TABLES

For every state  $s$ , there is a transition table  $\mathbf{T}$  and a reward table  $\mathbf{R}$ :

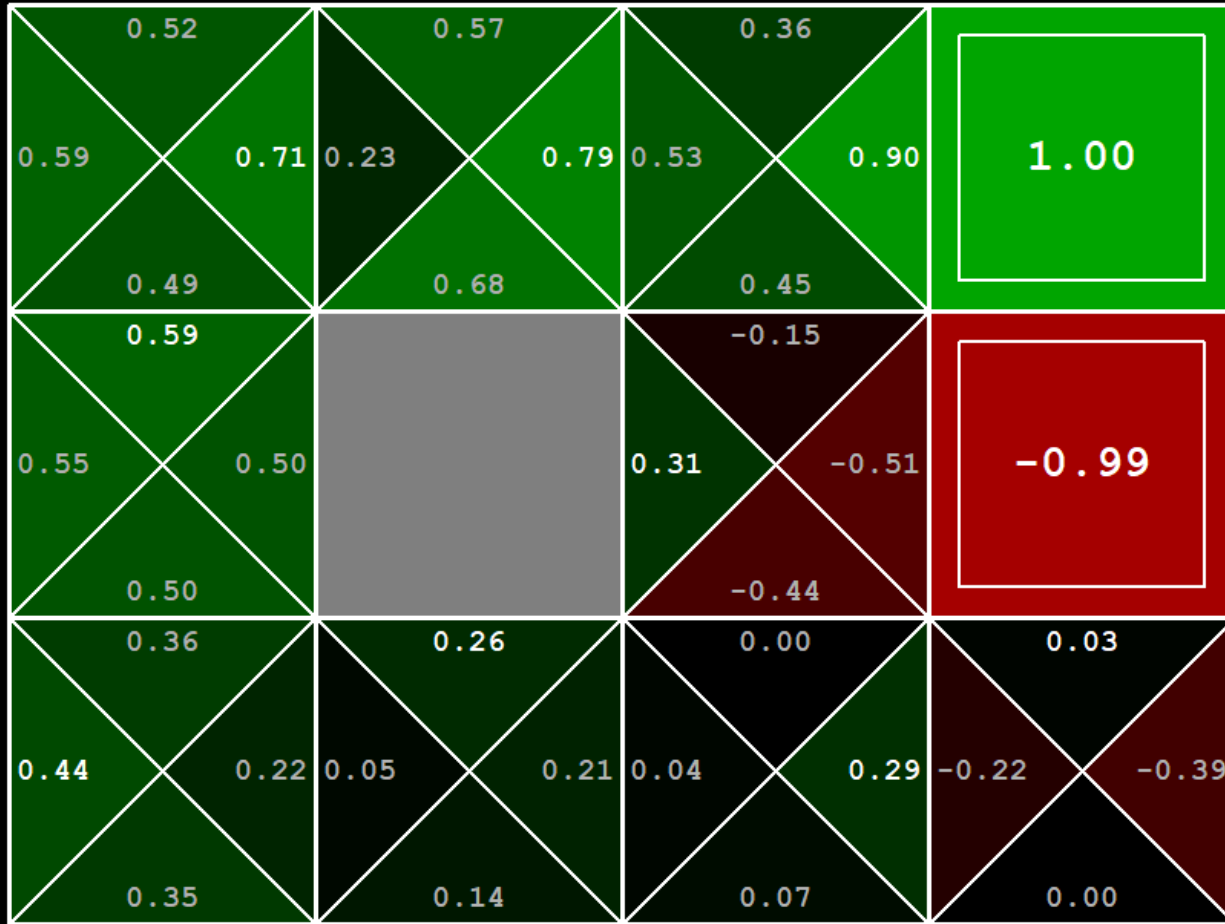
State  $s'$

Action  $a$

$\hat{T}(s, a0, s0)$ $\hat{R}(s, a0, s0)$	$\hat{T}(s, a0, s1)$ $\hat{R}(s, a0, s1)$	$\hat{T}(s, a0, s2)$ $\hat{R}(s, a0, s2)$	$\hat{T}(s, a0, s3)$ $\hat{R}(s, a0, s3)$
$\hat{T}(s, a1, s0)$ $\hat{R}(s, a1, s0)$	$\hat{T}(s, a1, s1)$ $\hat{R}(s, a1, s1)$	$\hat{T}(s, a1, s2)$ $\hat{R}(s, a1, s2)$	$\hat{T}(s, a1, s3)$ $\hat{R}(s, a1, s3)$
$\hat{T}(s, a2, s0)$ $\hat{R}(s, a2, s0)$	$\hat{T}(s, a2, s1)$ $\hat{R}(s, a2, s1)$	$\hat{T}(s, a2, s2)$ $\hat{R}(s, a2, s2)$	$\hat{T}(s, a2, s3)$ $\hat{R}(s, a2, s3)$



# Q-LEARNING EXAMPLE: GRIDWORLD



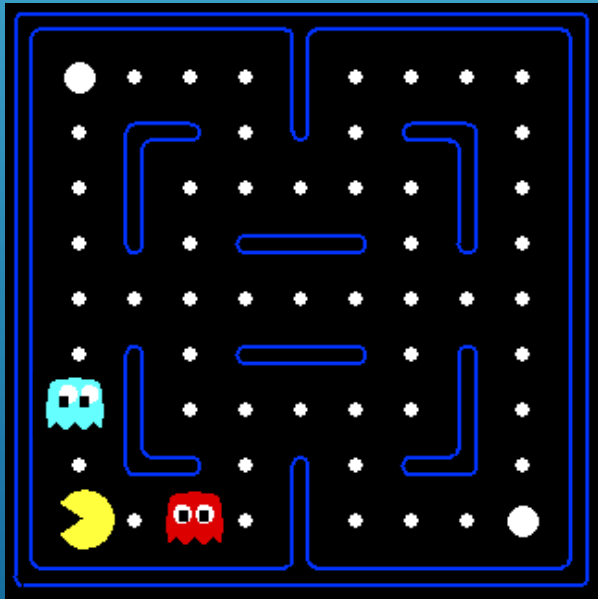
Q-VALUES AFTER 40 EPISODES

# APPROXIMATE REINFORCEMENT LEARNING

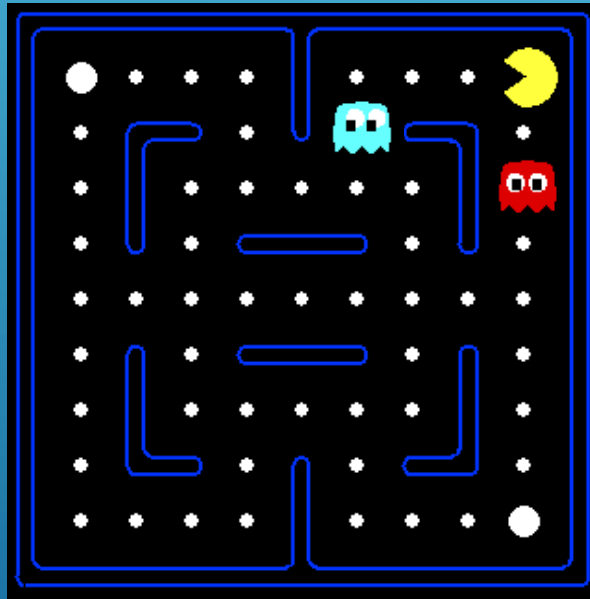
- Problems with tabular reinforcement learning:
  - Large state spaces need large amounts of memory.
  - Large state spaces require a long time to fill with accurate values.
- Key question: How can experience with a limited subset of the state space be usefully generalized to produce a good approximation over a much larger subset?
- Generalization can be accomplished with function approximation, which is an instance of supervised learning.
- In theory, any of the supervised learning techniques can be used as a function approximator, though some approaches are better than others.

# EXAMPLE: PACMAN

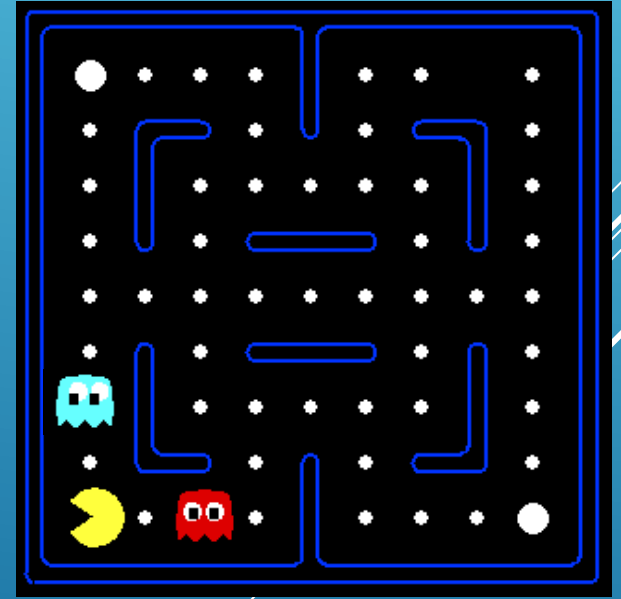
Let's say we discover through experience that this state is bad:

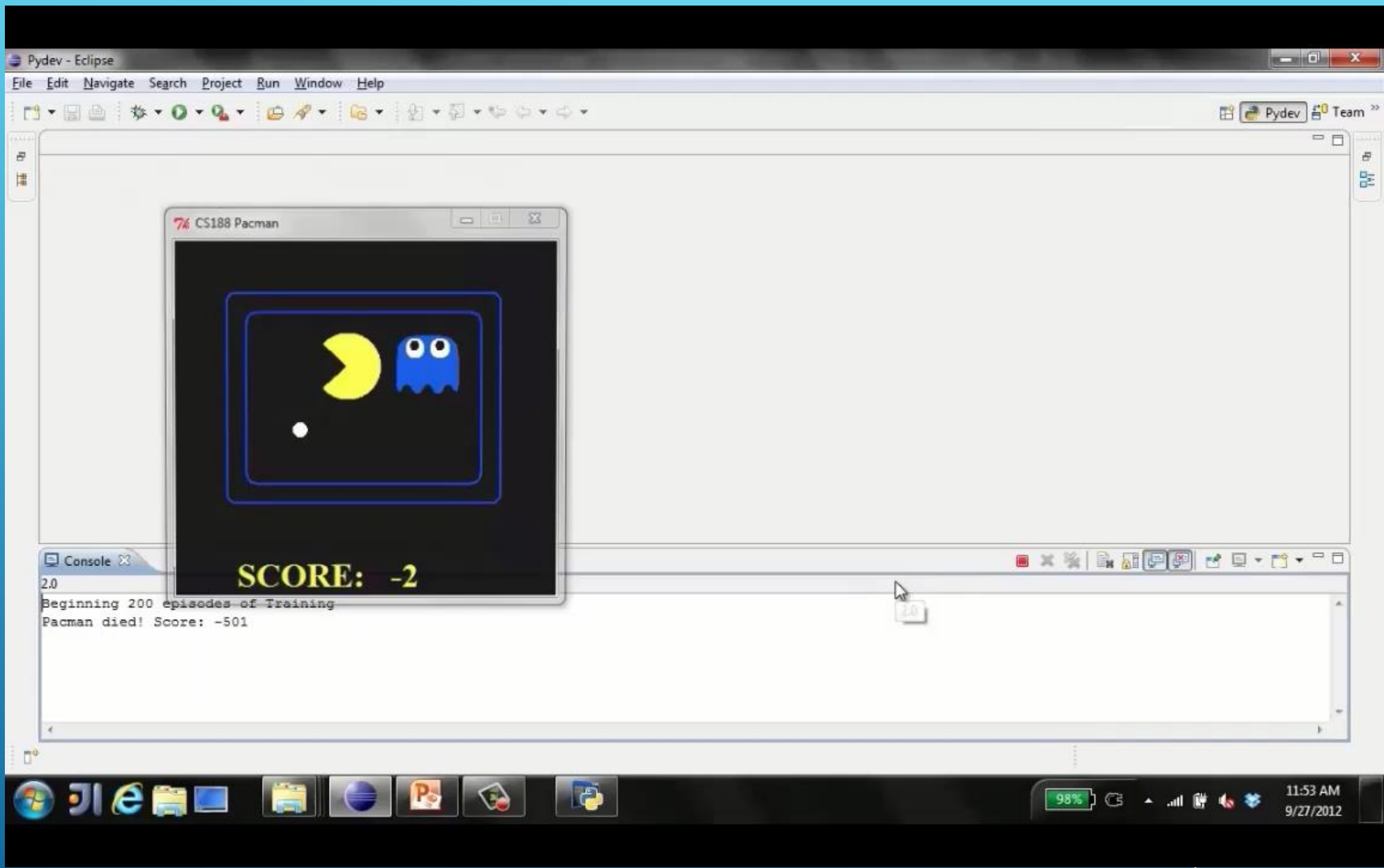


In naïve q-learning, we know nothing about this state:

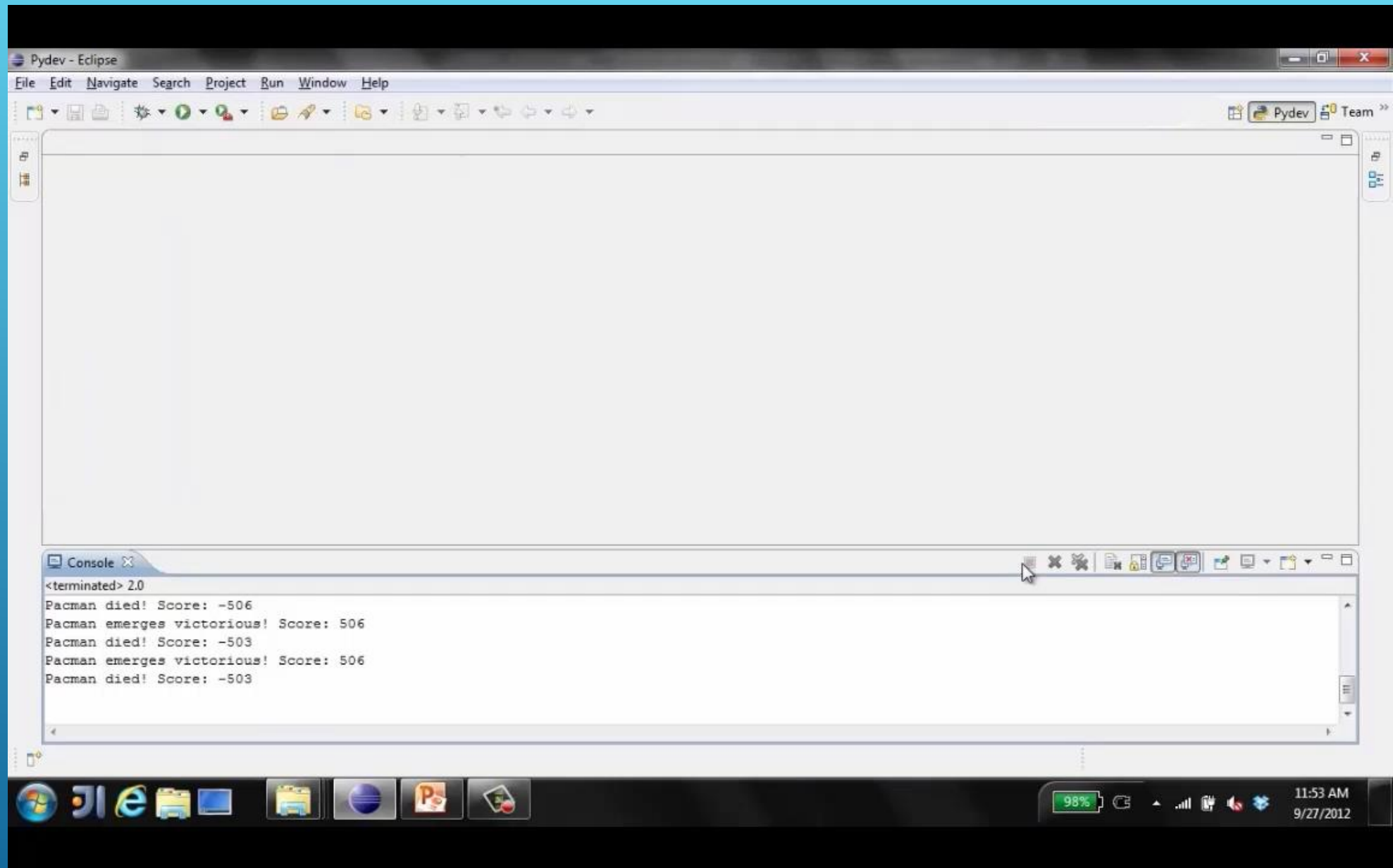


Or even this one!

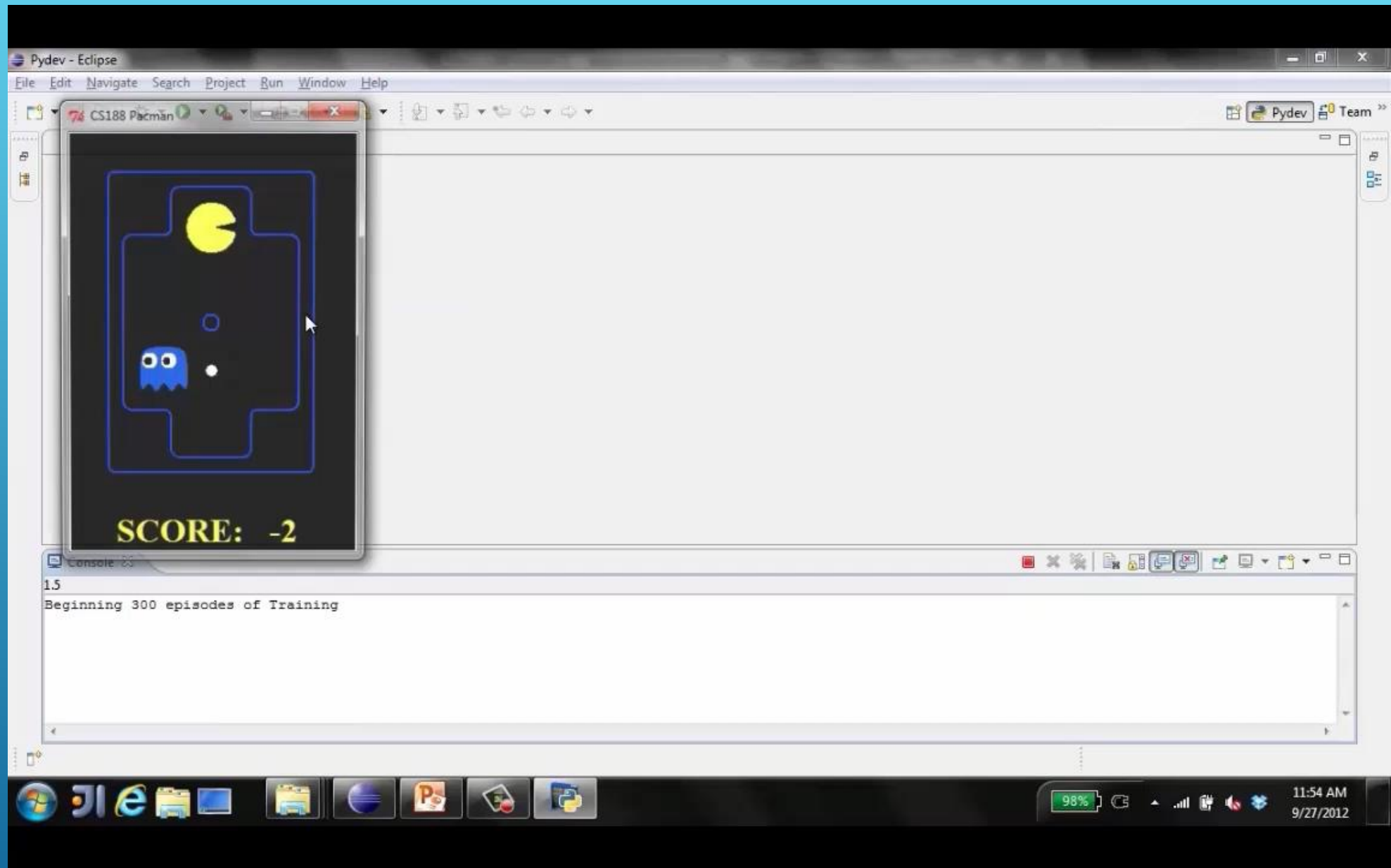




# TINY PACMAN DEMO 1



## TINY PACMAN DEMO 2

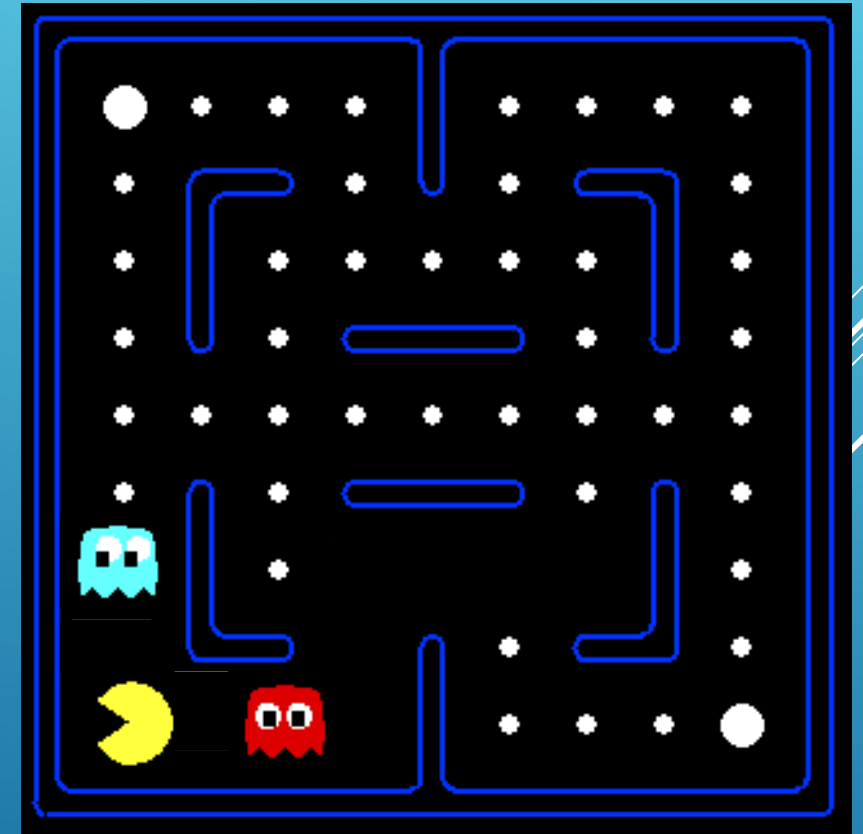


## TINY PACMAN DEMO 3

# FEATURE-BASED REPRESENTATIONS

Solution: describe a state using a vector of features (properties)

- ▶ Features are functions from states to real numbers (often 0/1) that capture important properties of the state
- ▶ Example features:
  - ▶ Distance to closest ghost
  - ▶ Distance to closest dot
  - ▶ Number of ghosts
  - ▶  $1 / (\text{dist to dot})^2$
  - ▶ Is Pacman in a tunnel? (0/1)
  - ▶ Is it the exact configuration on this slide?
  - ▶ ..... etc.
- ▶ Can also describe a q-state  $(s, a)$  with features (e.g. action moves closer to food)



# LINEAR VALUE FUNCTIONS

- ▶ Using a feature representation, we can write a  $q$  function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- ▶ **Advantage:** our experience is summed up in a few powerful numbers
- ▶ **Disadvantage:** states may share features but actually be very different in value!



# Q-LEARNING WITH LINEAR Q-FUNCTIONS

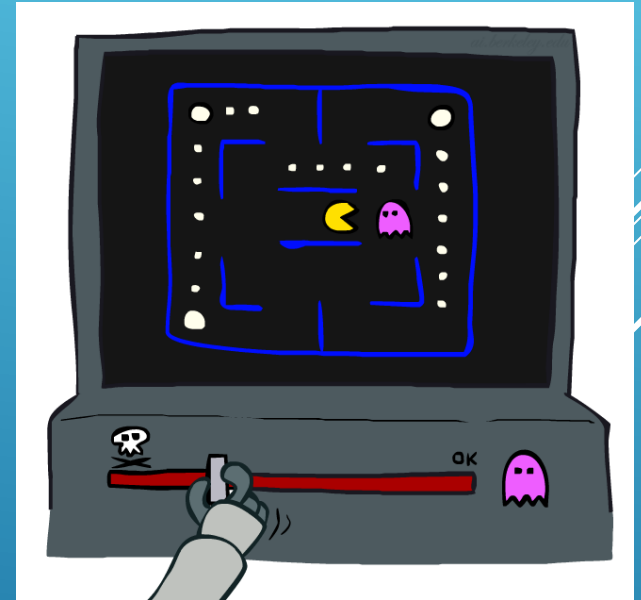
$$Q(s, a, \bar{w}) = w_1 x_1(s, a) + w_2 x_2(s, a) + \dots + w_n x_n(s, a)$$

## Exact Q's:

- ▶ transition =  $(s, a, r, s')$
- ▶ difference =  $\left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$
- ▶  $Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$

## Approximate Q's:

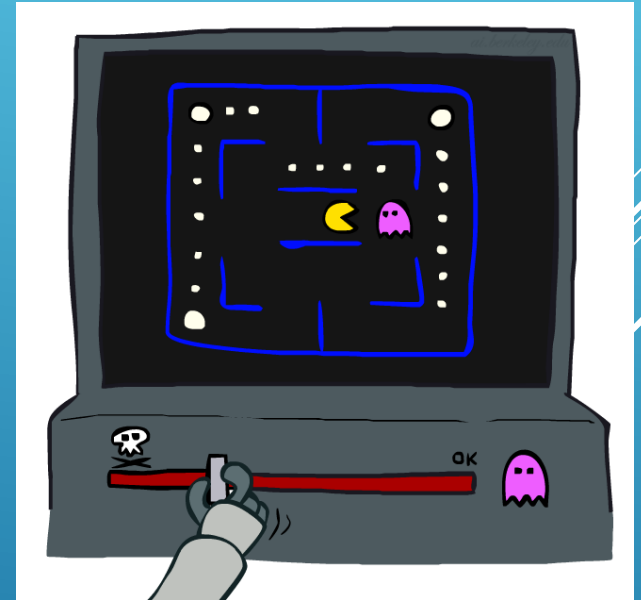
- ▶  $Q(s, a, \bar{w}) = \bar{w}^T \bar{x}$
- ▶  $\nabla Q(s, a, \bar{w}) = \bar{x}$  for linear functions
- ▶  $w_i \leftarrow w_i + \alpha [\text{difference}] \frac{\partial x_i(s, a)}{\partial w_i}$
- ▶  $w_i \leftarrow w_i + \alpha [\text{difference}] x_i$



# INTUITIVE INTERPRETATION

$$\mathbf{w}_i \leftarrow \mathbf{w}_i + \alpha[\text{difference}] \frac{\partial x_i(s, a)}{\partial \mathbf{w}_i}$$

- ▶ Adjust weights of active features.
- ▶ If something good happens, increase the features that were on i.e., prefer similar states more that have that state's features.
- ▶ If something unexpectedly bad happens, blame the features that were on i.e., prefer similar states less that have that state's features.



# Q-LEARNING UPDATE RULES: TABULAR VS. LINEAR ACTION-VALUE FUNCTION

**Tabular Update Rule:**

$$Q_{k+1}(s, a) \leftarrow Q_k(s, a) + \alpha \left[ \overset{\text{target}}{r + \gamma \max_{a'} Q_k(s', a')} - \overset{\text{current}}{Q_k(s, a)} \right]$$

**Linear Action-Value Update Rule:**

$$\bar{w} \leftarrow \bar{w} + \alpha \left[ \overset{\text{target}}{R + \gamma \max_{a'} \hat{q}(S', A', \bar{w})} - \overset{\text{current}}{\hat{q}(S, A, \bar{w})} \right] \nabla \hat{q}(S, A, \bar{w})$$

# APPROXIMATE Q-LEARNING ALGORITHM

## Q-Learning ~~Episodic Sarsa~~ with function approximation

### Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size  $\alpha > 0$ , small  $\varepsilon > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\varepsilon$ -greedy)

    Loop for each step of episode:

        Take action  $A$ , observe  $R, S'$

        If  $S'$  is terminal:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

            Go to next episode

        Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \cancel{\gamma \hat{q}(S', A', \mathbf{w})} - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

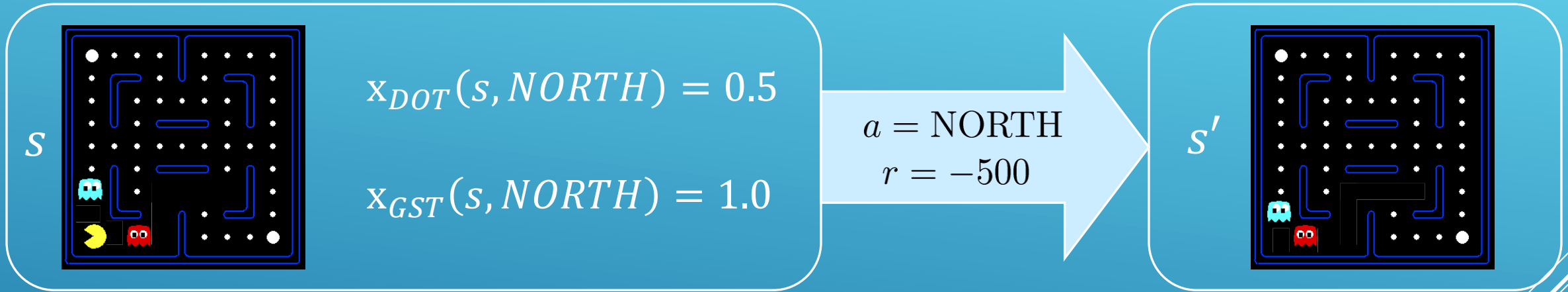
$S \leftarrow S'$

$A \leftarrow A'$

$\gamma \max_a \hat{q}(S', a, \mathbf{w})$

# EXAMPLE: Q-PACMAN

$$Q(s, a) = 4.0 x_{DOT}(s, a) - 1.0 x_{GST}(s, a)$$



$$x_{DOT}(s, \text{NORTH}) = 0.5$$

$$x_{GST}(s, \text{NORTH}) = 1.0$$

$$Q(s, \text{NORTH}) = +1$$

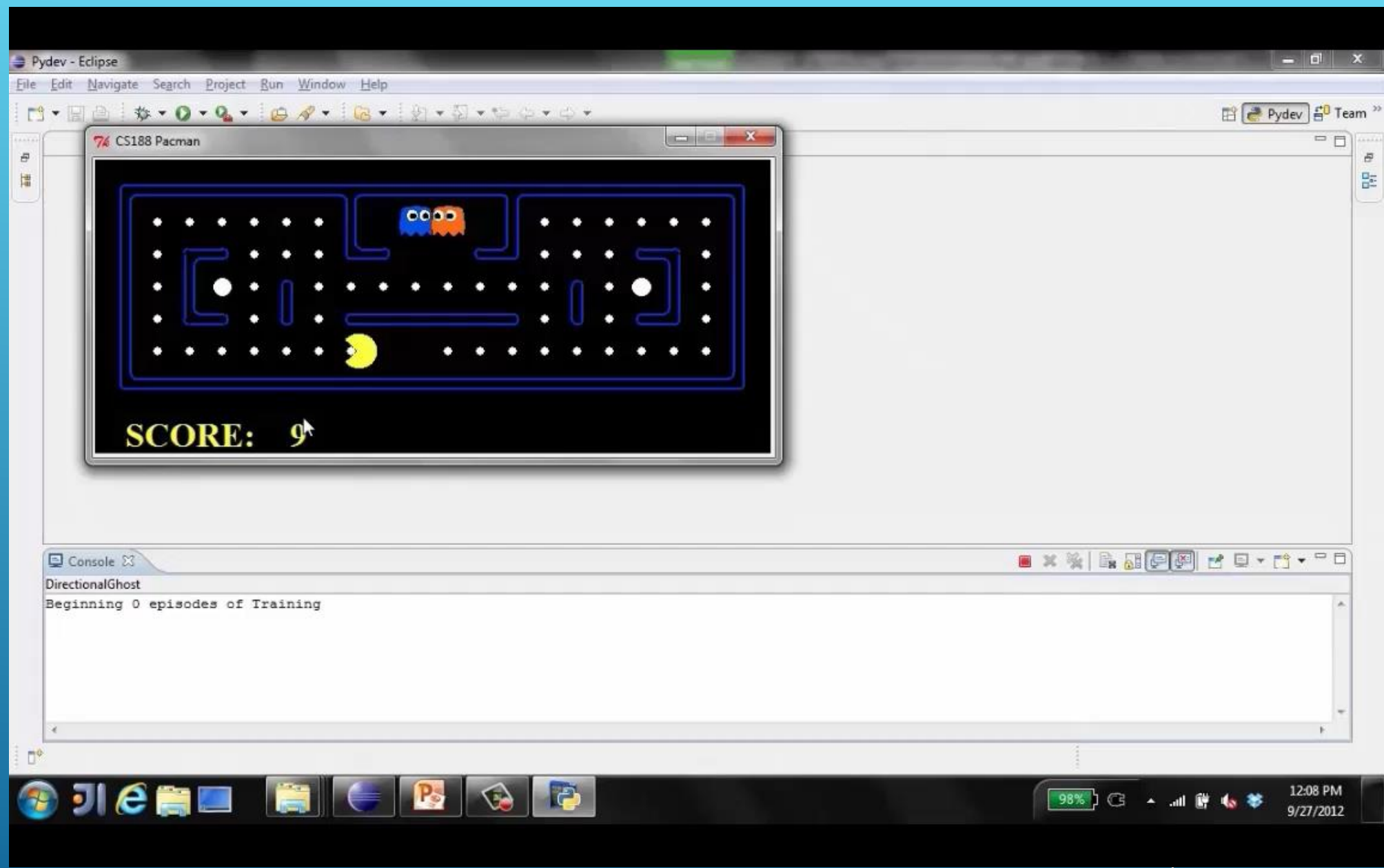
$$r + \gamma \max_{a'} Q(s', a') = -500 - 1 = -501$$

$$Q(s', \cdot) = 0$$

difference = -501

$$\begin{aligned} w_{DOT} &\leftarrow 4.0 + \alpha[-501]0.5 \\ w_{GST} &\leftarrow -1.0 + \alpha[-501]1.0 \end{aligned}$$

$$Q(s, a) = 3.0 x_{DOT}(s, a) - 3.0 x_{GST}(s, a)$$



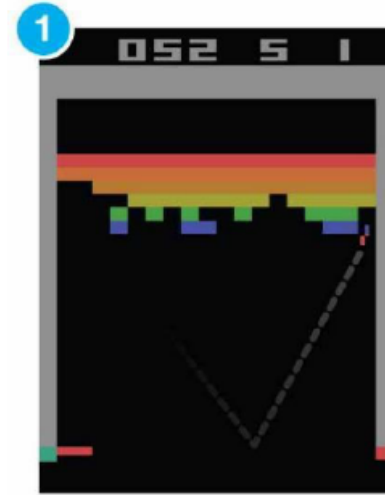
# APPROXIMATE Q-LEARNING DEMO -- PACMAN

# DEEP Q-LEARNING: APPROXIMATE Q-LEARNING WITH DEEP NEURAL NETWORKS



# Q-Learning: Representation Matters

- In practice, Tabular Q-Learning is impractical
  - Very limited states/actions
  - Cannot generalize to unobserved states



- Think about the **Breakout** game

- State: screen pixels
  - Image size:  $84 \times 84$  (resized)
  - Consecutive 4 images
  - Grayscale with 256 gray levels

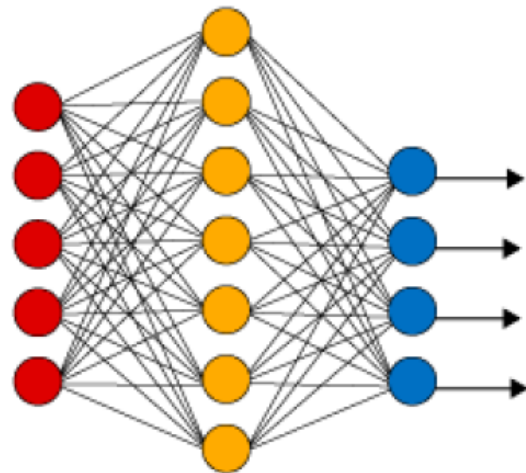
$256^{84 \times 84 \times 4}$  rows in the Q-table!

$= 10^{69,970} \gg 10^{82}$  atoms in the universe

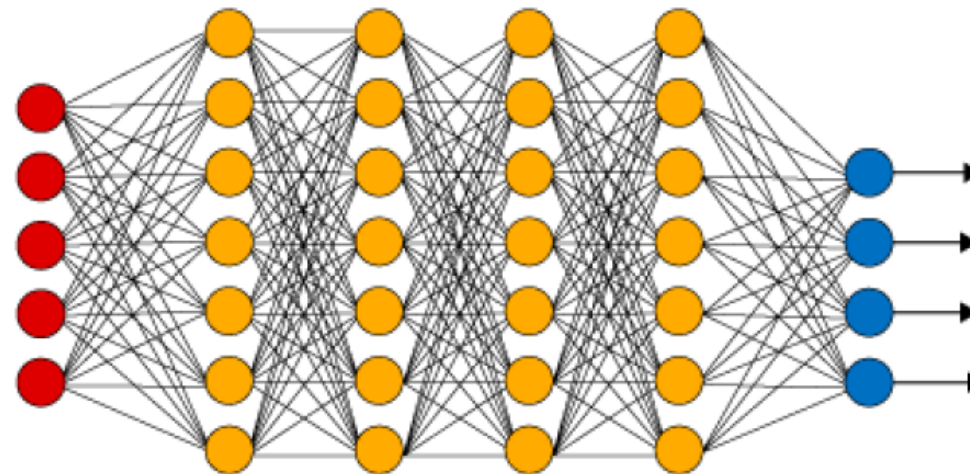


# Combing Neurons in Hidden Layers: The “Emergent” Power to Approximate

Simple Neural Network



Deep Learning Neural Network



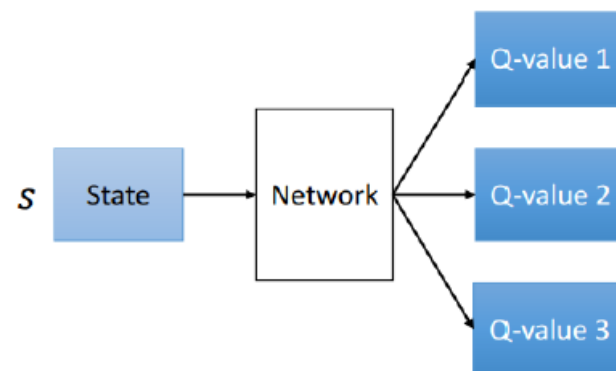
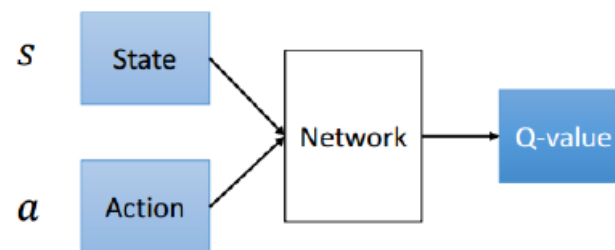
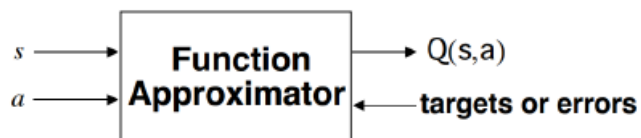
● Input Layer    ● Hidden Layer    ● Output Layer

**Universality:** For any arbitrary function  $f(x)$ , there exists a neural network that closely approximate it for any input  $x$

# DQN: Deep Q-Learning

Use a neural network to approximate the Q-function:

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



***EXTRA SLIDES***

EDITINGA



# Q-LEARNING UPDATE RULE: AN ALTERNATE INTERPRETATION

$$Q_{k+1}(s, a) \leftarrow Q_k(s, a) + \alpha \left[ \overset{\text{target}}{r + \gamma \max_{a'} Q_k(s', a')} - \overset{\text{current}}{Q_k(s, a)} \right]$$



$$Q_{k+1}(s, a) \leftarrow \overset{\text{old}}{(1 - \alpha) Q_k(s, a)} + \alpha \left[ \overset{\text{new}}{r + \gamma \max_{a'} Q_k(s', a')} \right]$$