

Q-LEARNING (YET AGAIN!)

Scott O'Hara

Metrowest Developers Machine Learning Group

12/12/2018

REFERENCES

The material for this talk is primarily drawn from the slides, notes and lectures of these courses:

University of California, Berkeley CS188:

- ▶ *CS188 – Introduction to Artificial Intelligence*, Profs. Dan Klein, Pieter Abbeel, et al.
<http://ai.berkeley.edu/home.html>


David Silver, DeepMind:

- ▶ *Introduction to Reinforcement Learning*
<http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>

CS181 course at Harvard University:

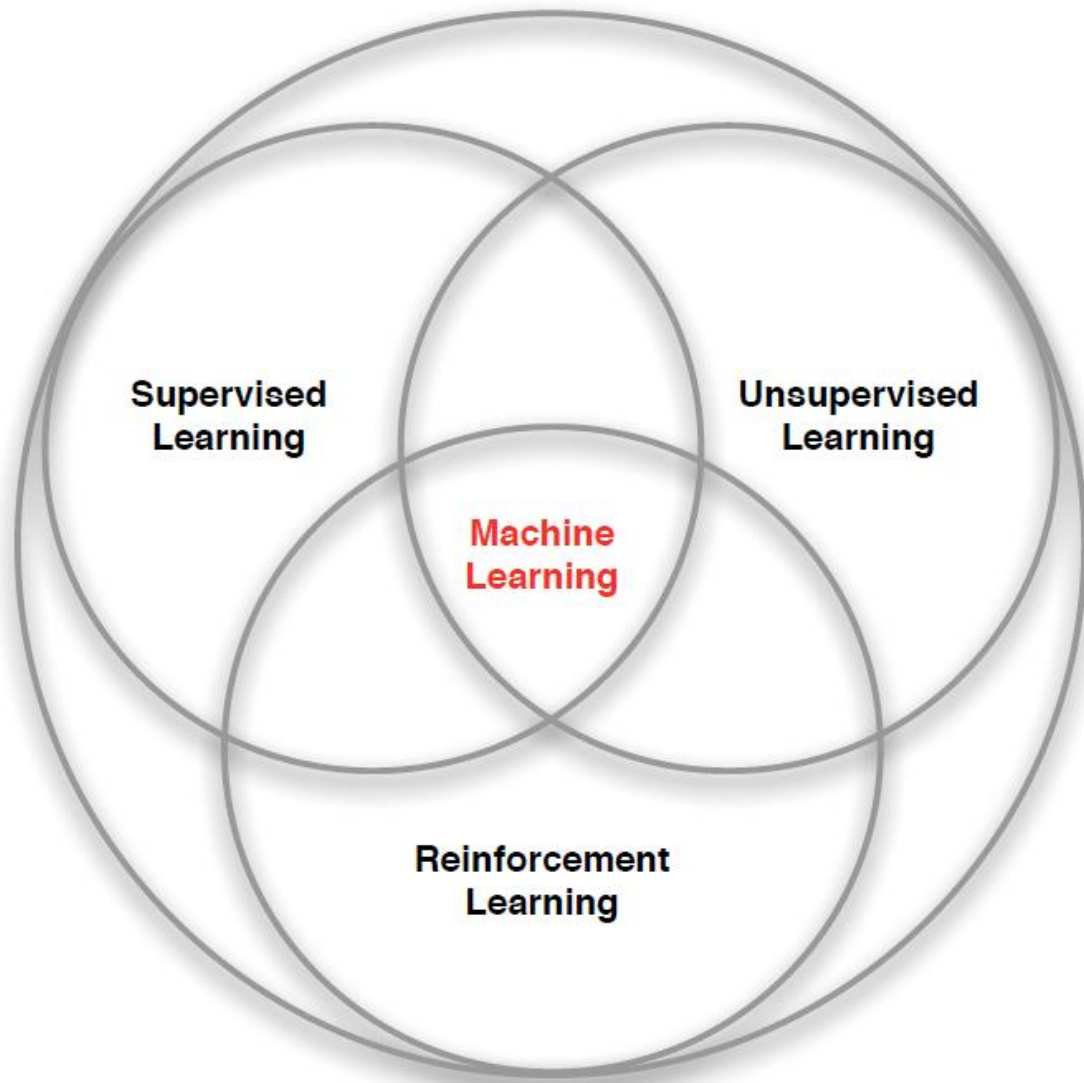
- ▶ *CS181 Intelligent Machines: Perception, Learning and Uncertainty*, Sarah Finney, Spring 2009
- ▶ *CS181 Intelligent Machines: Perception, Learning and Uncertainty*, Prof. David C Brooks, Spring 2011
- ▶ *CS181 – Machine Learning*, Prof. Ryan P. Adams, Spring 2014. <https://github.com/wihl/cs181-spring2014>
- ▶ *CS181 – Machine Learning*, Prof. David Parkes, Spring 2017. <https://harvard-ml-courses.github.io/cs181-web-2017/>

Q-LEARNING: TALK OUTLINE

- Reinforcement Learning
 - What is Q-Learning
 - Q-Learning Demos
 - Approximate Q-Learning
- 
- Several thin, parallel white lines of varying lengths and orientations are positioned in the bottom right corner of the slide, creating a modern, abstract graphic element.

REINFORCEMENT LEARNING





BRANCHES OF MACHINE LEARNING

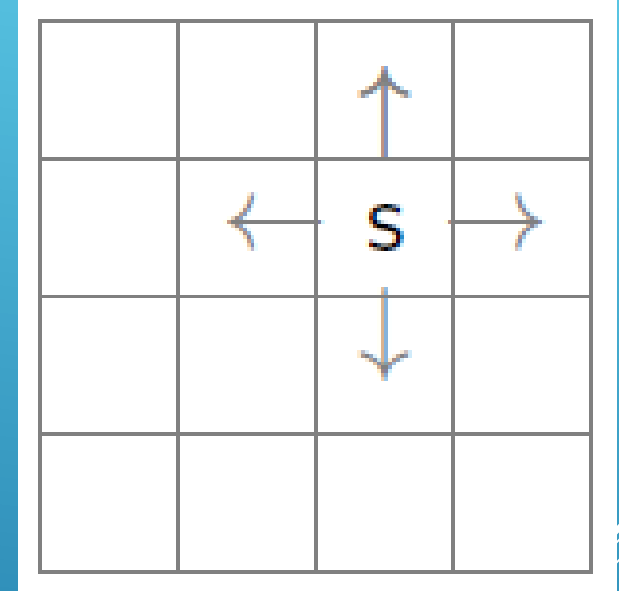
Credit: adapted from lecture slides David Silver, DeepMind, "Introduction to Reinforcement Learning"

REINFORCEMENT LEARNING: THE BASIC IDEA

- Select an action.
- If action leads to reward, reinforce that action.
- If action leads to punishment, avoid that action.
- Basically, a computational form of Behaviorism (Pavlov, B. F. Skinner).

THE LEARNING FRAMEWORK

- Learning is performed **online**, learn as we interact with the world
- In contrast with supervised learning, there are no training or test sets. The reward is accumulated over interactions with the environment.
- Data is not fixed, more information is acquired as you go.
- The training distribution can be influenced by action decisions.



EXAMPLES OF REINFORCEMENT LEARNING

- Fly stunt maneuvers in a helicopter
- Defeat the world champion at Backgammon
- Manage an investment portfolio
- Control a power station
- Make a humanoid robot walk
- Play Atari games better than humans

EXAMPLES OF REINFORCEMENT LEARNING

RL Course by David Silver – Lecture 1: Introduction to Reinforcement Learning

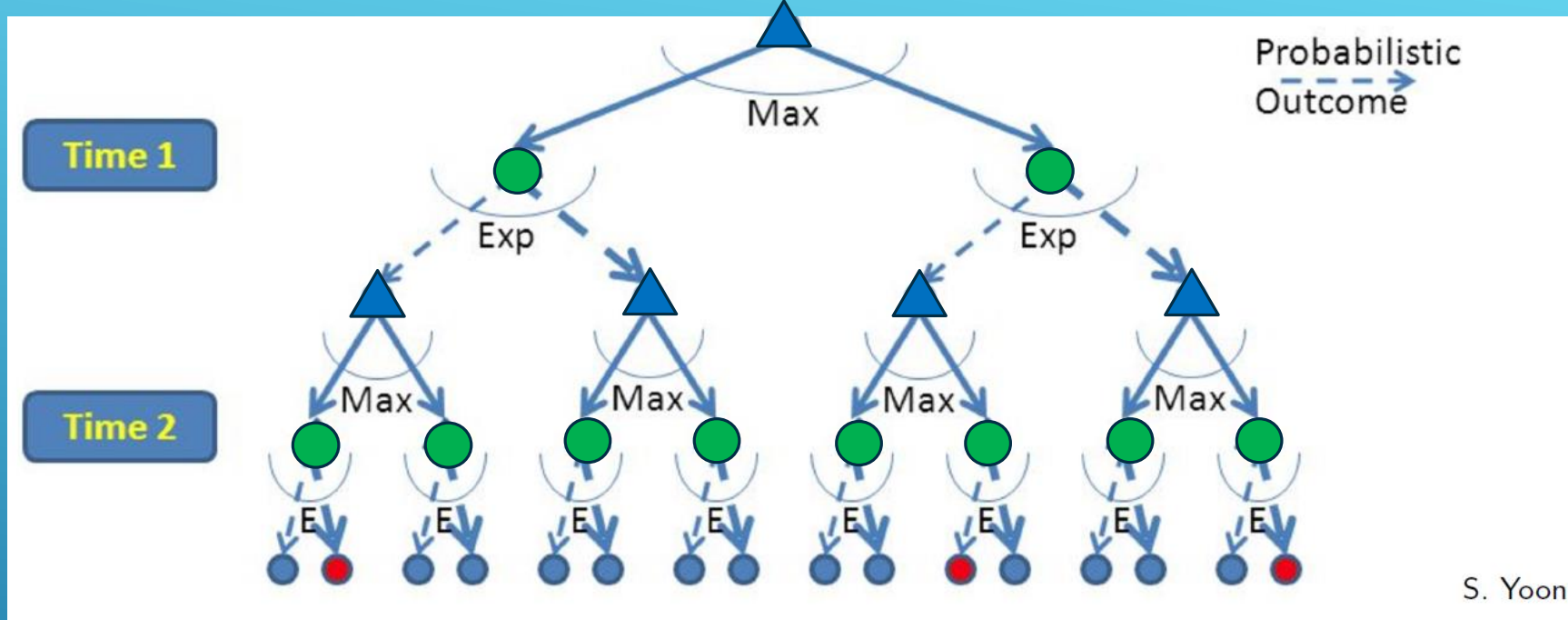
<https://www.youtube.com/watch?v=2pWv7GOvuf0>

12:25 – 22.00

Q-LEARNING



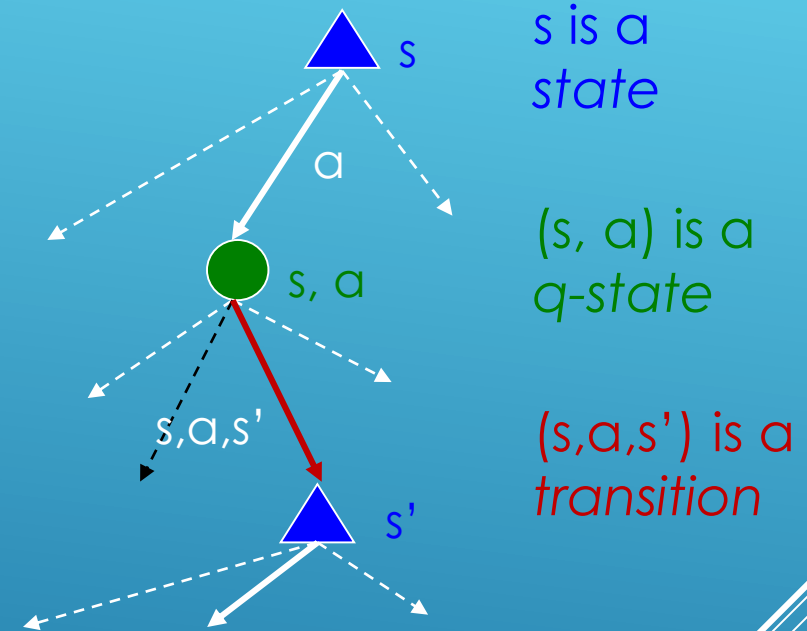
RL IS LIKE A GAME AGAINST NATURE



- 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)

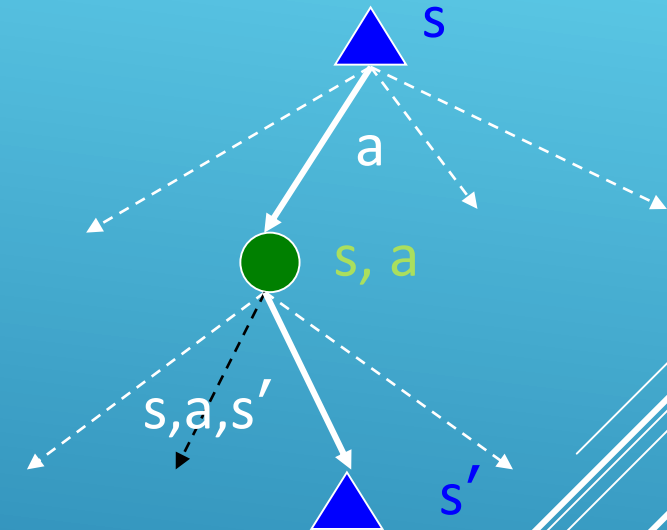
QUANTITIES TO OPTIMIZE

- The value (utility) of a **state s** :
 $V(s)$ = expected utility starting in s and acting optimally
- The value (utility) of a **q-state (s,a)** :
 $Q(s,a)$ = expected utility starting out having taken action a from state s and (thereafter) acting optimally
- The optimal policy:
 $\pi(s)$ = optimal action from state s



THE BELLMAN EQUATIONS

- ▶ The Bellman Equations define a relationship, which when satisfied guarantees that $V(\mathbf{s})$ and $Q(\mathbf{s}, \mathbf{a})$ are optimal for each state and action.
- ▶ This in turn guarantees that the policy π^* is optimal.
- ▶ There is one equation $V^*(\mathbf{s})$ for each state \mathbf{s} .
- ▶ There is one equation $Q^*(\mathbf{s}, \mathbf{a})$ for each state \mathbf{s} and action \mathbf{a} .

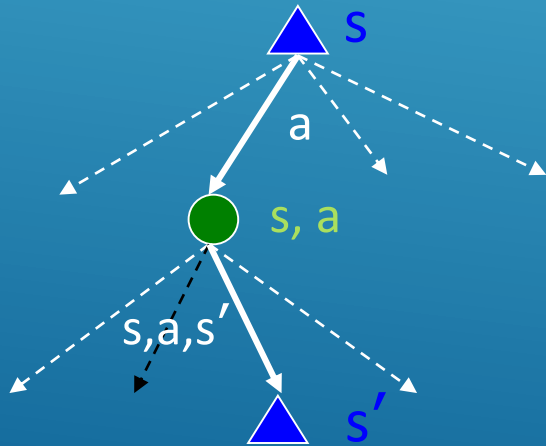


THE BELLMAN EQUATIONS

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

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

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



- **States:** s_1, \dots, s_n
- **Actions:** a_1, \dots, a_m
- **Reward Function:**
 $R(s, a, s') \in R$
- **Transition model:**
 $T(s, a, s') = P(s' | s, a)$
- **Discount factor:** $\gamma \in [0, 1]$

THE OPTIMAL VALUE UTILITY EQUATION V^*

- ▶ Focusing on different Bellman Equations Gives Different Algorithms
- ▶ The V^* equation gives rise to these algorithms previously discussed:
 - ▶ Expectimax
 - ▶ Value Iteration
 - ▶ Policy Iteration

$$V^*(s) = \max_a Q^*(s, a) \quad [1]$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')] \quad [2]$$



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

THE OPTIMAL VALUE UTILITY EQUATION Q^*

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

$$V^*(s) = \max_a Q^*(s, a) \quad [1]$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')] \quad [2]$$

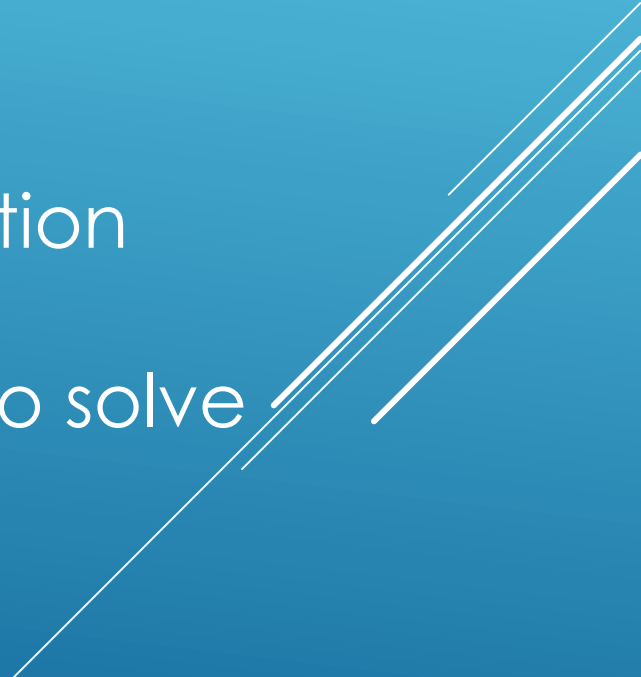


$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$$

MODEL-BASED VS. MODEL-FREE

- Model-Based RL
 - Explicitly learn:
 - the transition model $T(s,a,s')$
 - The reward model $R(s,a,s')$
 - Use Value Iteration, etc. to find optimal policy π^* .
- Model-Free RL
 - Don't learn $T(s,a,s')$ and $R(s,a,s')$.
 - Q-Learning – learn $Q(s,a)$ directly.

MODEL-BASED RL PROS AND CONS

- **Pros:**
 - Makes maximal use of experience.
 - Solves model optimally, given enough experience.
 - **Cons:**
 - Requires computationally expensive solution procedure.
 - Requires the model to be small enough to solve
- 
- A series of three parallel white diagonal lines in the bottom right corner of the slide, extending from the middle of the right edge towards the bottom left.

MODEL-FREE RL PROS AND CONS

- **Pros:**

- Solution procedure is much more efficient.
- Can handle much larger models.

- **Cons:**

- Learns more slowly.
- Does not learn as much as a model-based RL in a single training episode. (“Leaves information on the table”).

Q-LEARNING

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

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

- ▶ 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: $Q(s,a) = r + \gamma \max_{a'} Q_k(s', a')$

- ▶ Incorporate the new estimate into a running average based on the learning rate α :

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q_k(s', a') \right]$$

Q-LEARNING UPDATE RULE

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

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q_k(s', a') \right]$$

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

Q-LEARNING ALGORITHM

For each state s and action a :

$$Q(s, a) \leftarrow 0$$

Begin in state s :

Repeat:

Given s and the set of actions A_s associated with s :

- ▶ **CHOOSE ACTION** $a \in A_s \leftarrow \text{HOW?}$
- ▶ Apply action a .
- ▶ Receive reward r and new state s' .
- ▶ With transition $\langle s, a, r, s' \rangle$, update $Q(s, a)$

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q_k(s', a') \right]$$

$$s \leftarrow s'$$

CHOOSING AN ACTION: EXPLORATION VS EXPLOITATION

- **Exploit:** maximize the expected utility now.
- **Explore:** choose an action that will help you improve your model.
- **How to Exploit?** use the current policy.
 - e.g., for Q-Learning in state s , choose action with largest $Q(s,a)$
- **How to Explore?**
 - choose an action randomly
 - choose an action you haven't chosen yet
 - choose an action that will take you to an unexplored state.

EXPLORATION STRATEGY: ϵ -GREEDY

- Explore with probability ϵ . Exploit with probability $1 - \epsilon$.
- Weaknesses:
 - Does not exploit when learning has converged.
- Uses:
 - appropriate if the world is changing.

EXPLORATION STRATEGY: BOLTZMANN

- In state s , choose action a with probability p :

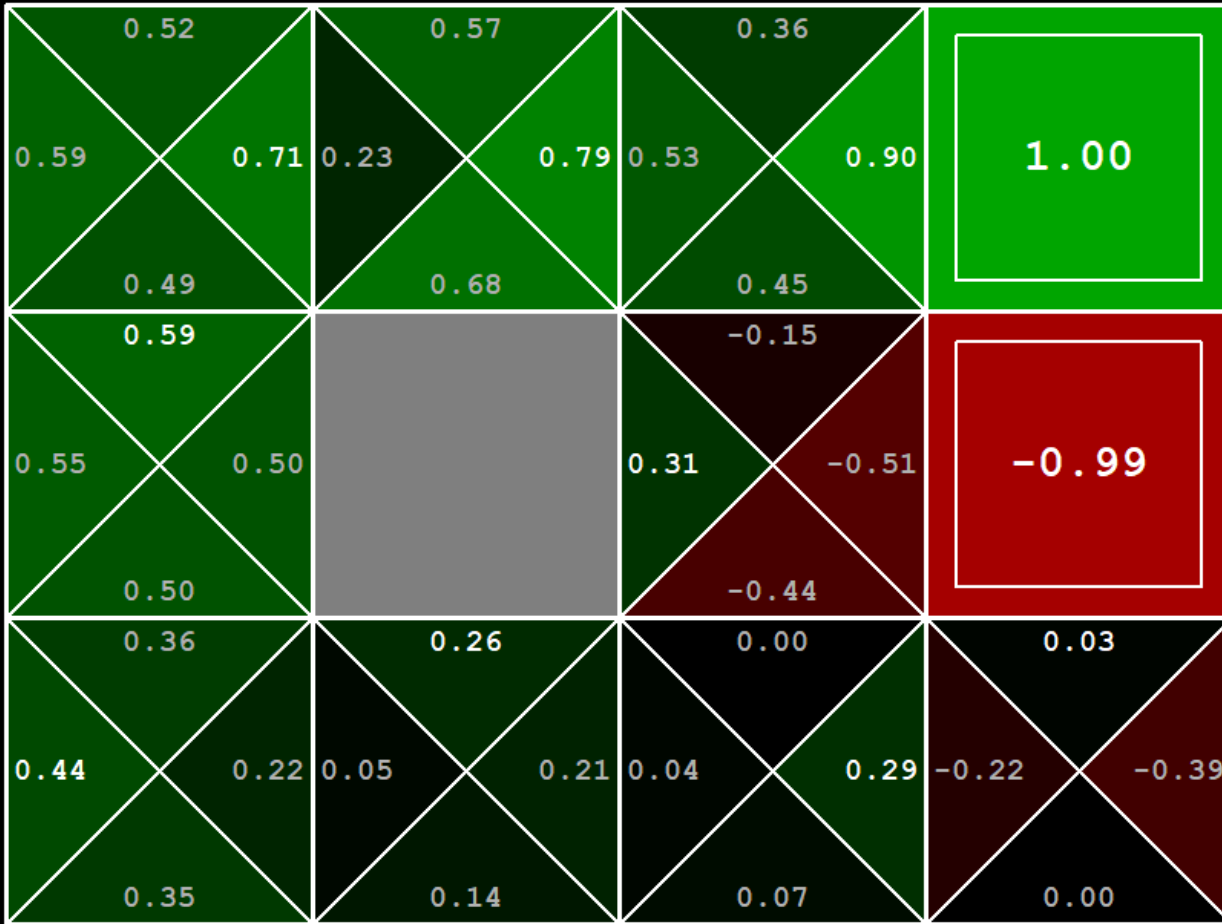
$$p = \frac{e^{\frac{Q(s,a)}{t}}}{\sum_{a'} e^{\frac{Q(s,a')}{t}}}$$

- Simulated annealing: t is a “temperature” which “cools” over time.
- High temperature means more exploration
 - $e^{\frac{Q(s,a)}{t}} \rightarrow 1$, implies actions taken at about equal probability.
- As t cools, exploration is reduced.
 - $e^{\frac{Q(s,a)}{t}} \rightarrow \infty$, largest $Q(s,a)$ becomes most probable.

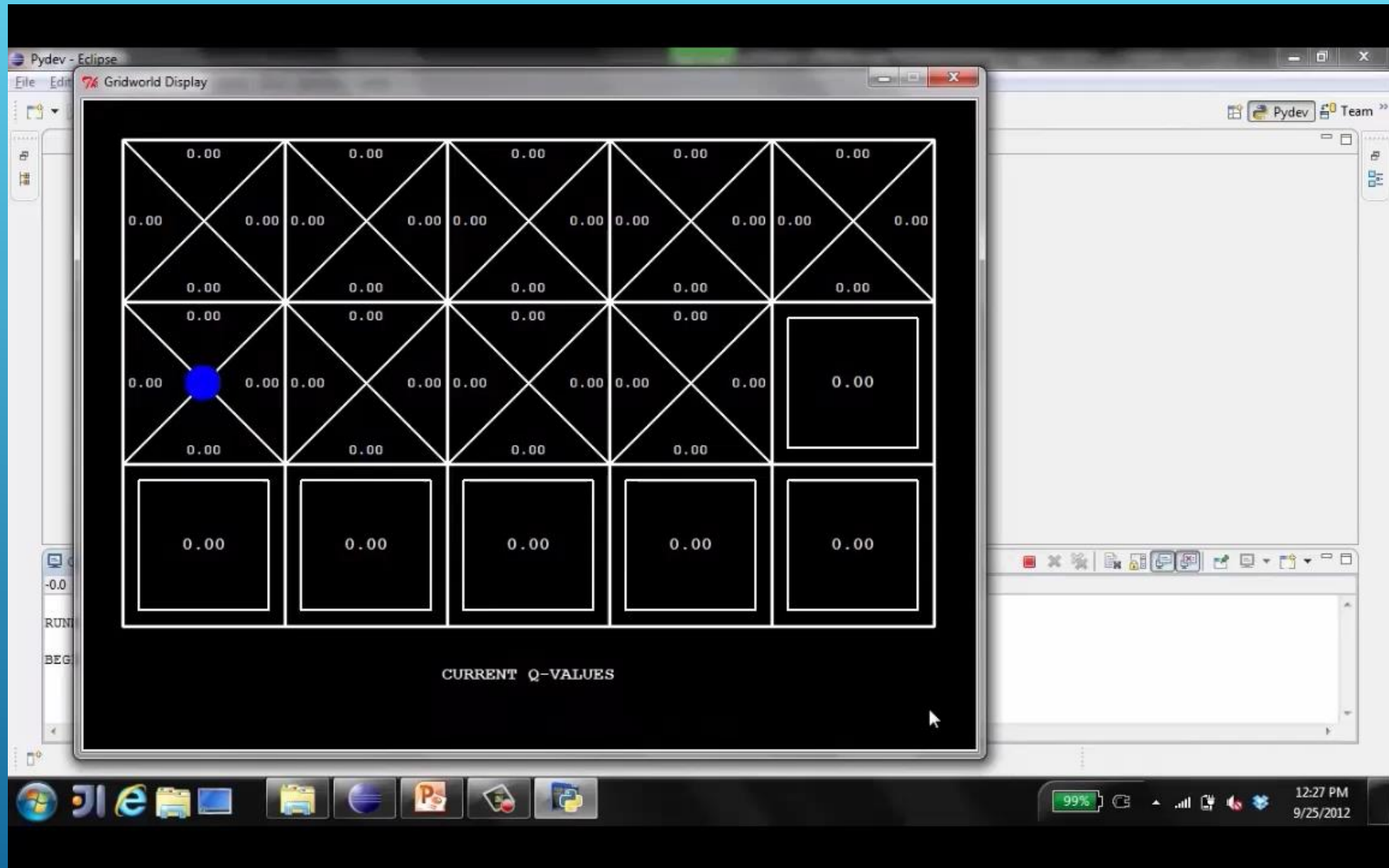
Q-LEARNING DEMOS



Q-LEARNING EXAMPLE: GRIDWORLD

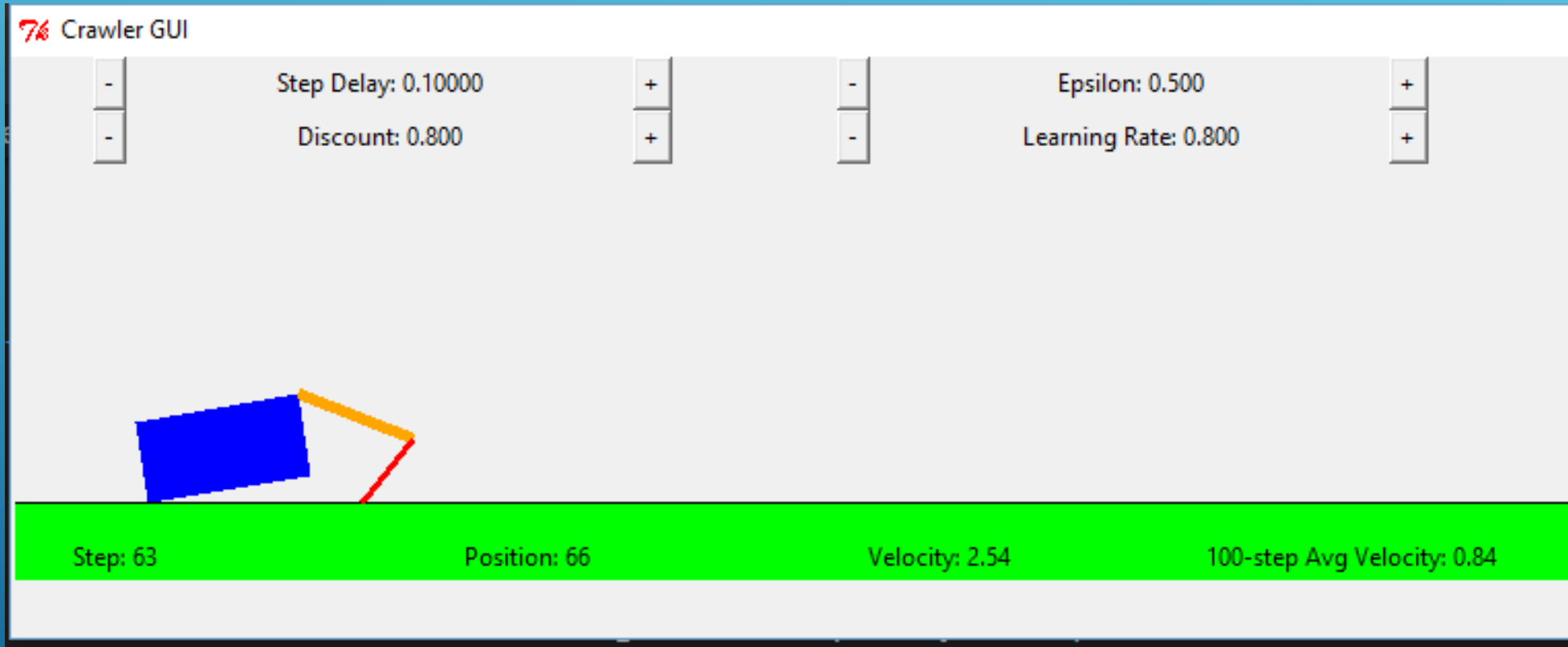


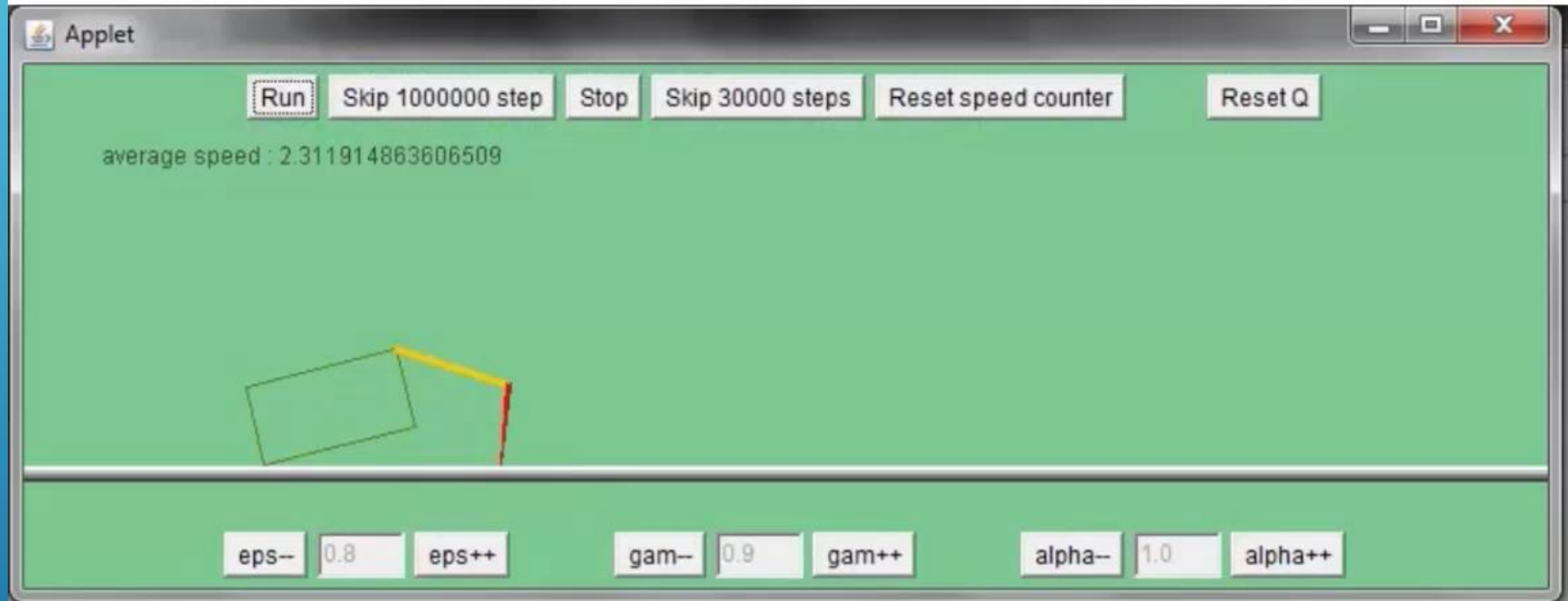
Q-VALUES AFTER 40 EPISODES



VIDEO OF Q-LEARNING DEMO -- GRIDWORLD

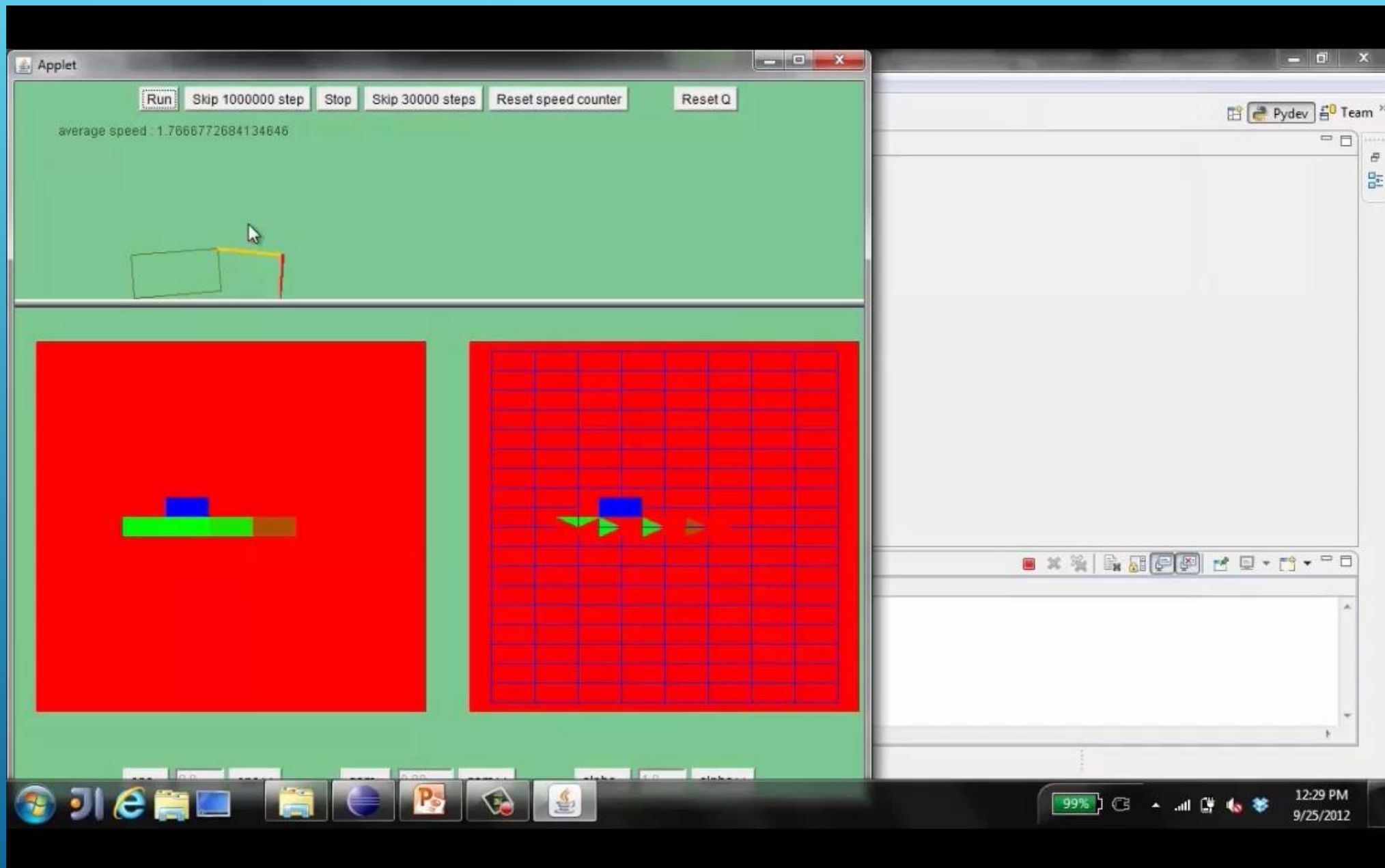
Q-LEARNING EXAMPLE: CRAWLER ROBOT





VIDEO 1 OF Q-LEARNING DEMO -- CRAWLER

Credit: CS188 – Intro to AI, UC Berkeley, ai.berkeley.edu



VIDEO 2 OF Q-LEARNING DEMO -- CRAWLER

Credit: CS188 – Intro to AI, UC Berkeley, ai.berkeley.edu

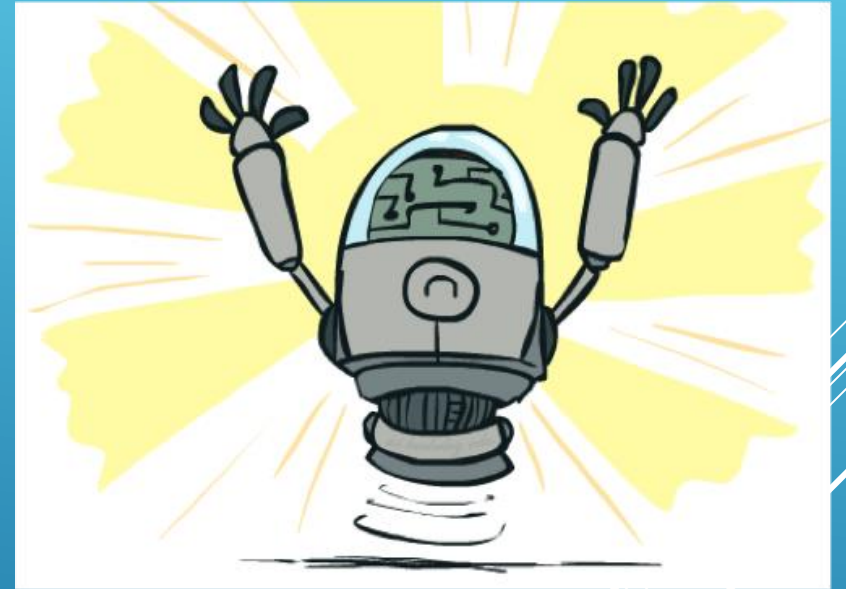
Q-LEARNING EXAMPLE: DISCOUNT EFFECT

Update rule: $Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right]$

Description	Training Steps	<u>Discount</u> (γ)	Learning Rate (α)	Avg Velocity
Default	~100K	0.8	0.8	~1.73
Low γ	~100K	0.5	0.8	0.0
High γ	~100K	0.919	0.8	~3.33
Low α	~100K	0.8	0.2	~1.73
High α	~100K	0.8	0.9	~1.73
High γ , Low α	~100K	0.919	0.2	~3.33

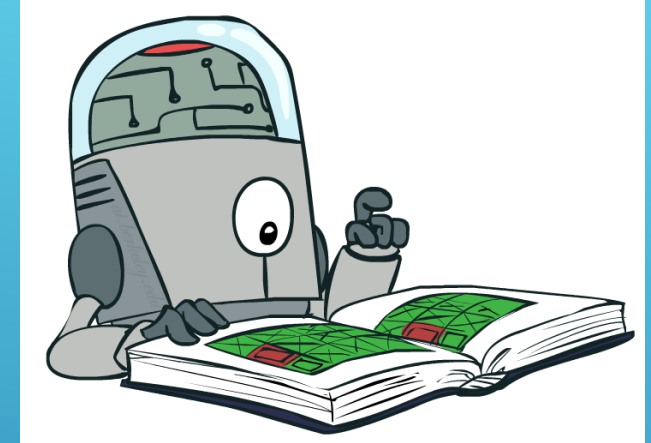
Q-LEARNING PROPERTIES

- ▶ Amazing result: Q-learning converges to optimal policy -- even if you're acting sub-optimally!
- ▶ This is called **off-policy learning**
- ▶ Caveats:
 - ▶ You have to explore enough
 - ▶ You have to eventually make the learning rate small enough
 - ▶ ... but not decrease it too quickly
 - ▶ Basically, in the limit, it doesn't matter how you select actions (!)



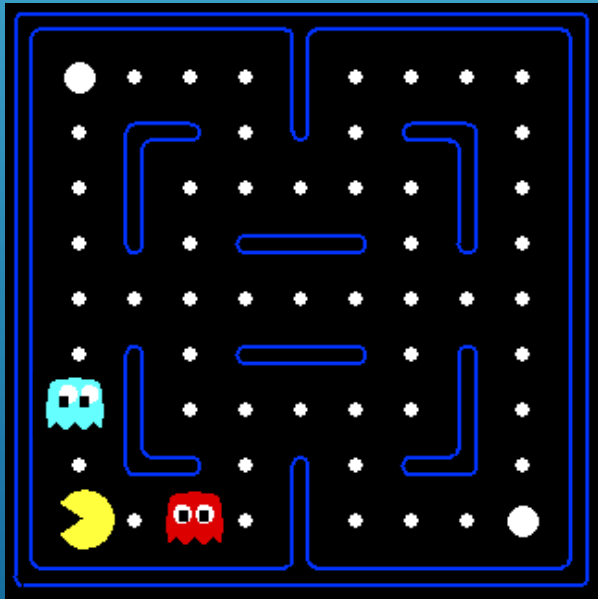
GENERALIZING ACROSS STATES

- ▶ Basic Q-Learning keeps a table of all q-values
- ▶ In realistic situations, we cannot possibly learn about every single state!
 - ▶ Too many states to visit them all in training
 - ▶ Too many states to hold the q-tables in memory
- ▶ Instead, we want to generalize:
 - ▶ Learn about some small number of training states from experience
 - ▶ Generalize that experience to new, similar situations

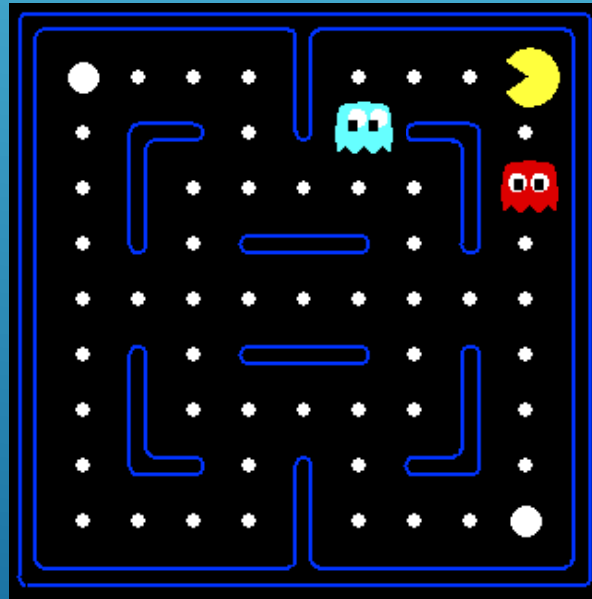


EXAMPLE: PACMAN

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

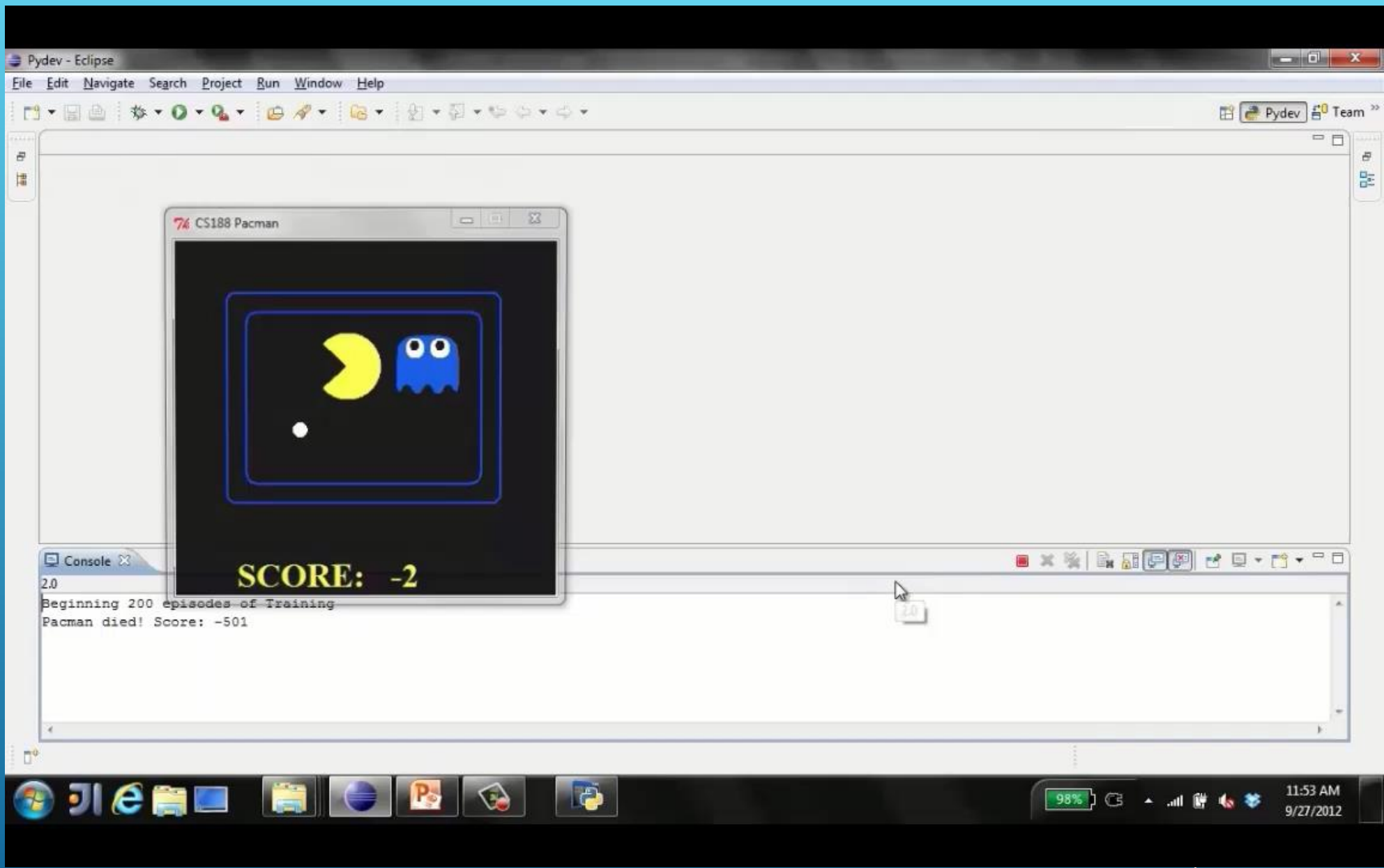


In naïve q-learning, we know nothing about this state:

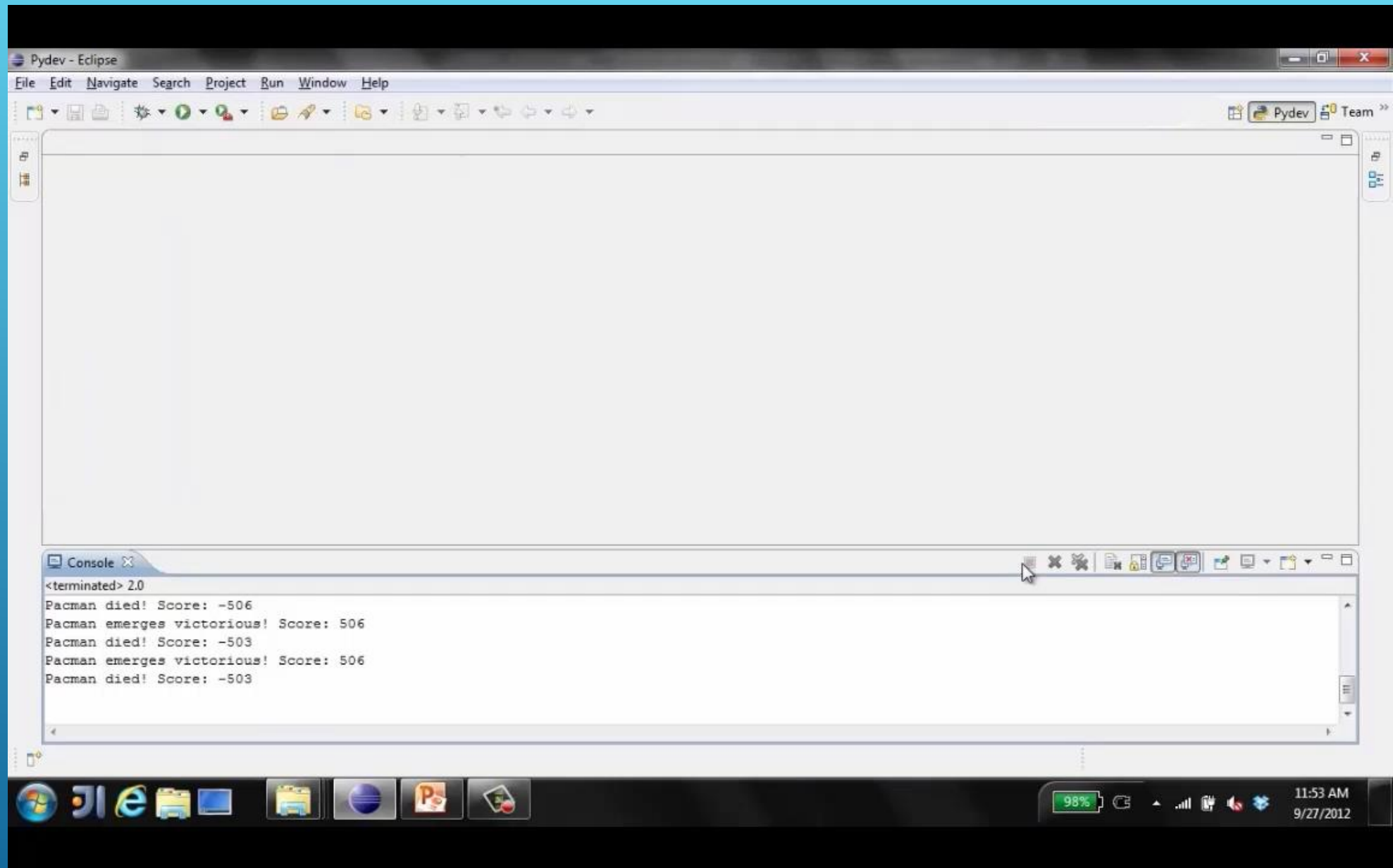


Or even this one!

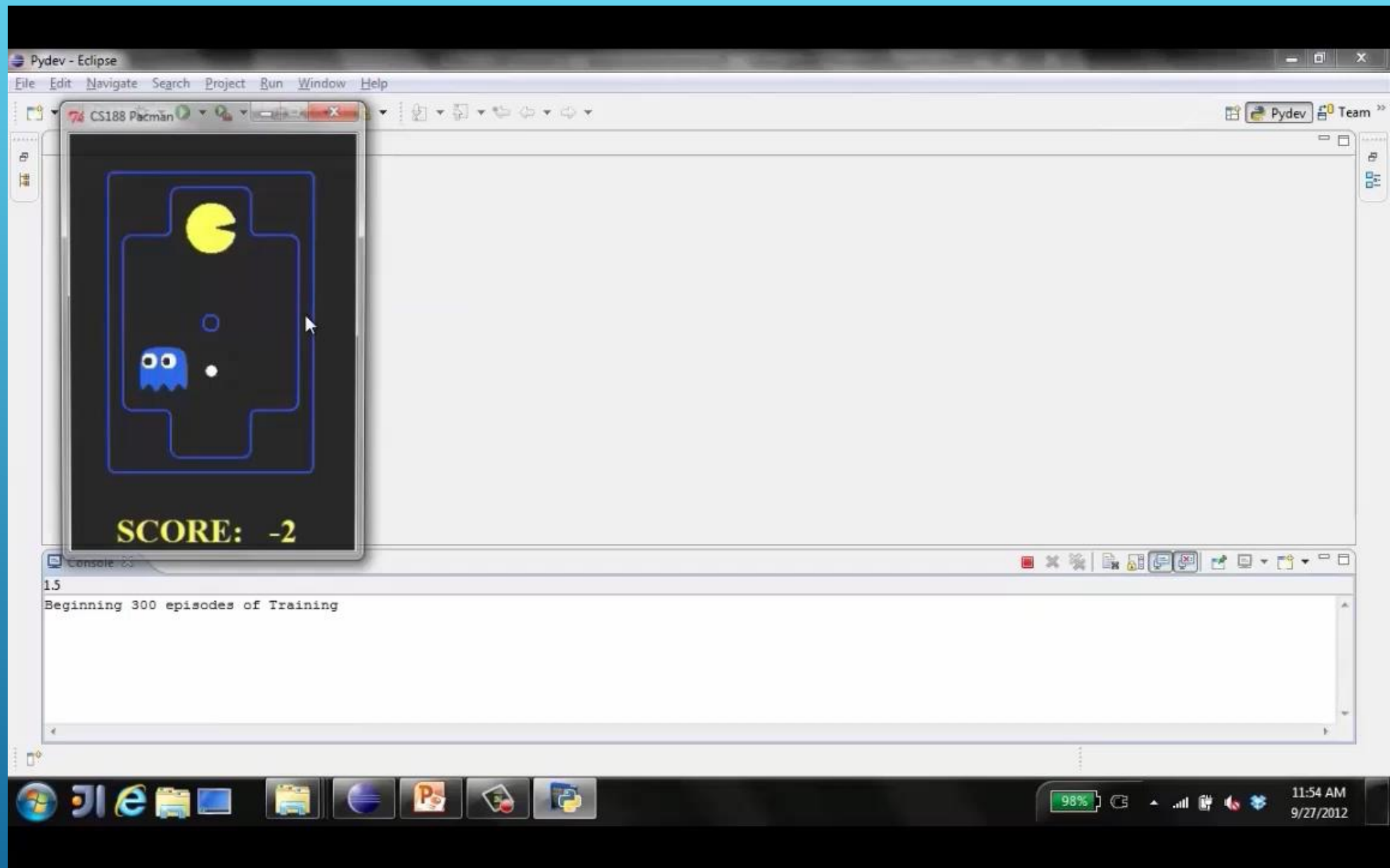




TINY PACMAN DEMO 1



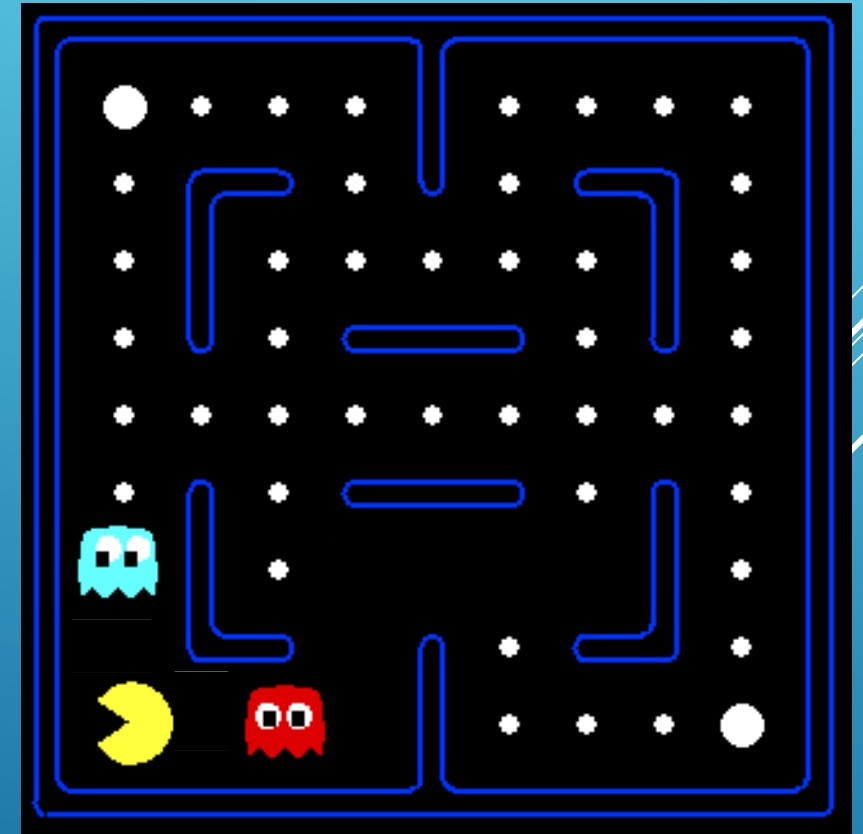
TINY PACMAN DEMO 2



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)
 - ▶ etc.
 - ▶ Is it the exact state on this slide?
 - ▶ 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!

APPROXIMATE Q-LEARNING

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

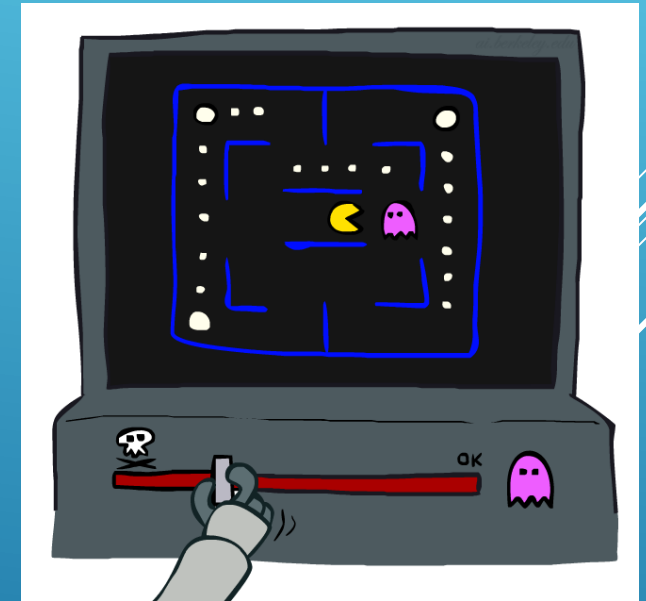
- ▶ Q-learning with linear Q-functions:

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

- ▶ Intuitive interpretation:

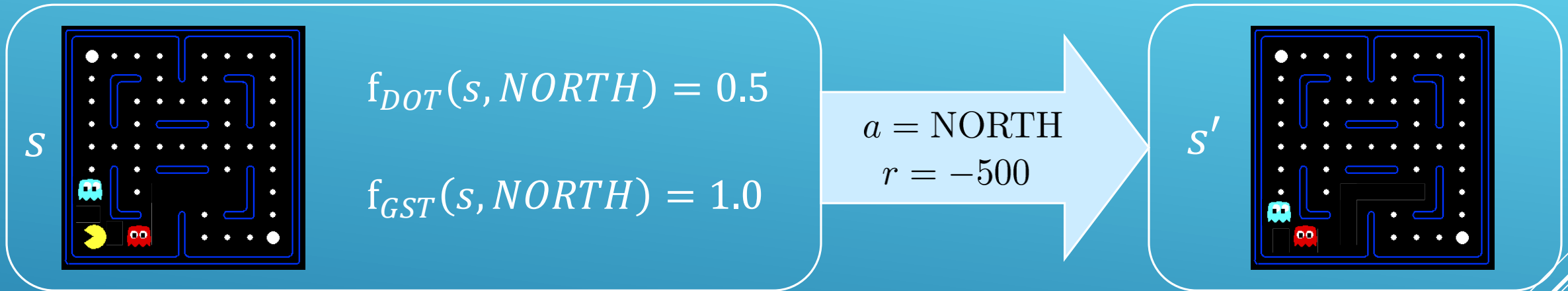
- ▶ Adjust weights of active features.
- ▶ E.g., if something unexpectedly bad happens, blame the features that were on i.e., prefer less all states with that state's features.

- ▶ Formal justification: online least squares



EXAMPLE: Q-PACMAN

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



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

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

$a = \text{NORTH}$
 $r = -500$

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

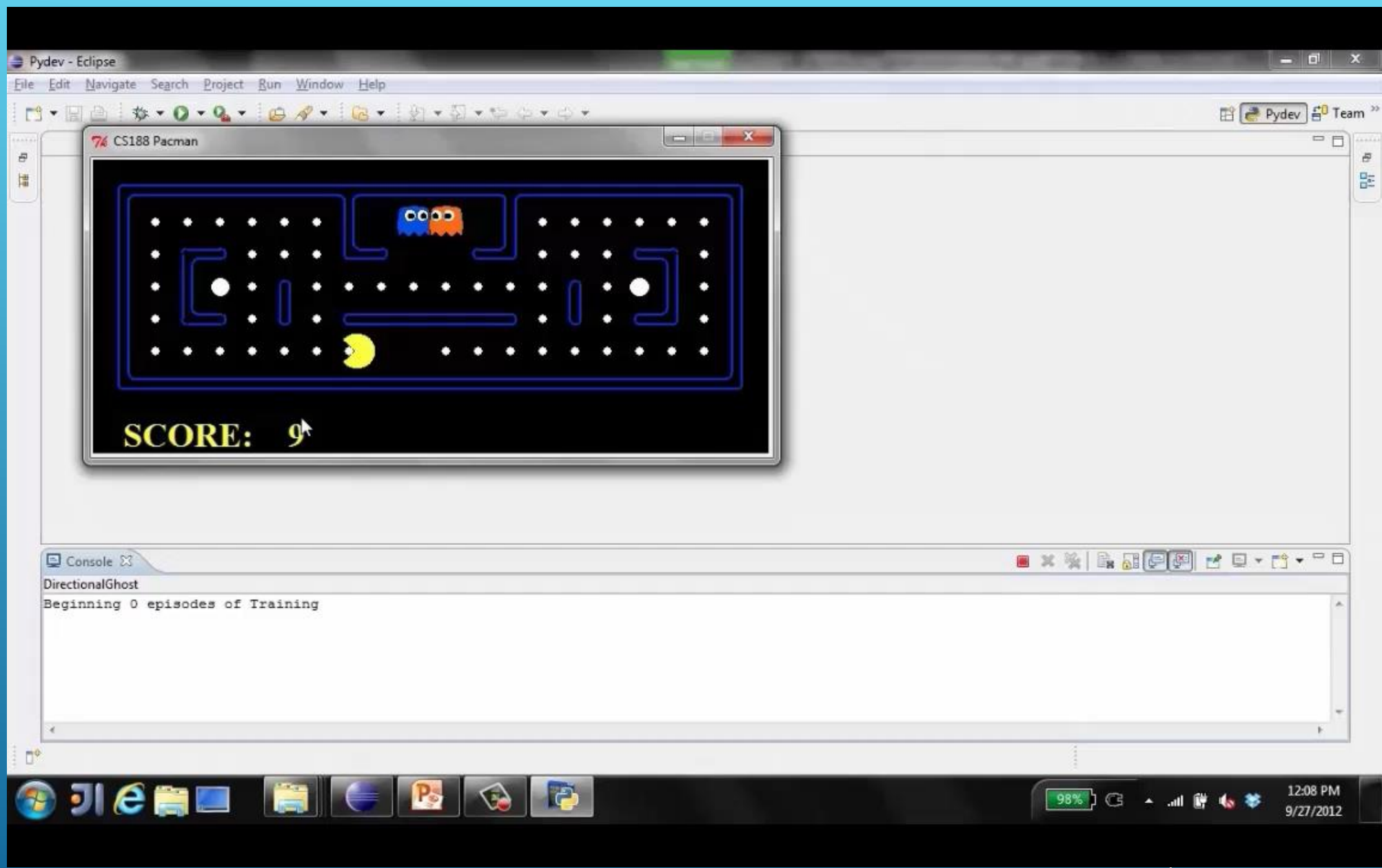
$$r + \gamma \max_{a'} Q(s', a') = -500 + 0$$

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

difference = -501

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

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



APPROXIMATE Q-LEARNING DEMO -- PACMAN

NEXT TIME: MORE RL

Possible Topics:

- Credit Assignment
- Online Learning vs Offline learning
 - SARSA vs Q-Learning
- Hybrid Approaches:
 - Temporal Difference (TD) Value Learning
 - Dyna-Q
- Partially Observable Markov Decision Processes (POMDPs)
- Deep Q-Learning

EXTRA SLIDES

MARKOV DECISION PROCESSES

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

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

- **Transition model:**

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

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

