# Artificial Intelligence
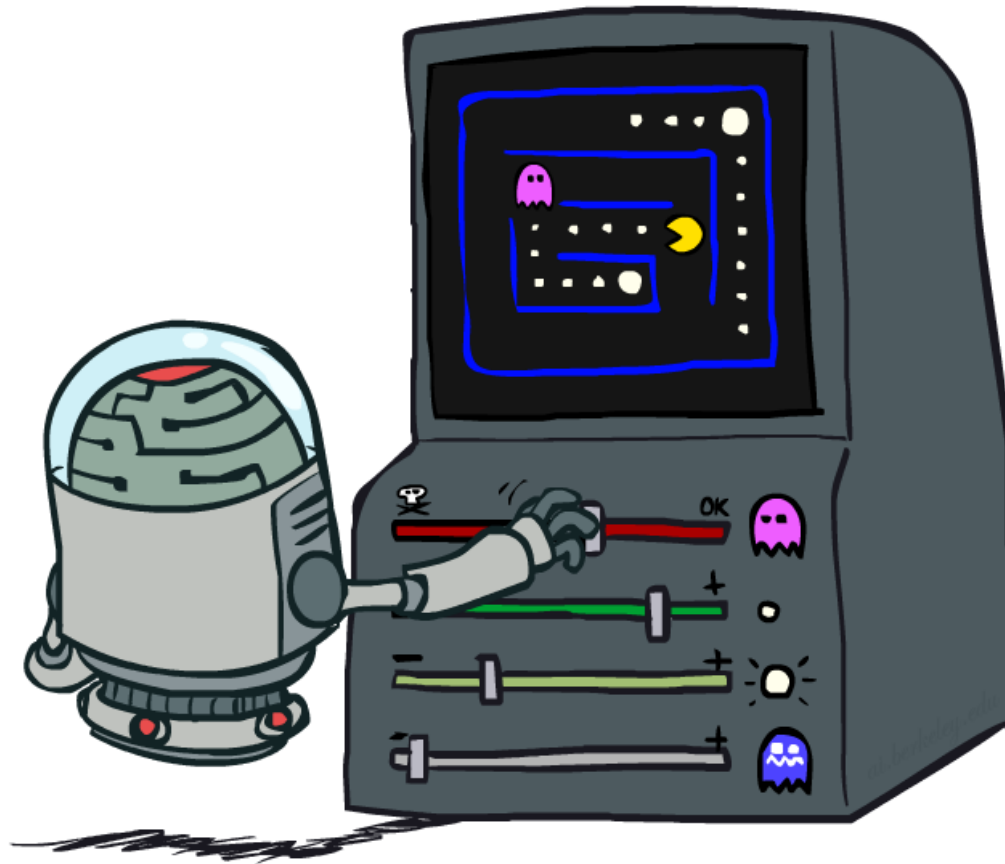## Reinforcement Learning



[These slides were created by Dan Klein, Pieter Abbeel, and Anca Dragan. http://ai.berkeley.edu.]

# Reinforcement Learning

o We still assume an MDP:

  o A set of states $s \in S$

  o A set of actions (per state) A

  o A model T(s,a,s')

  o A reward function R(s,a,s')

o Still looking for a policy $\pi(s)$

o New twist: don't know T or R, so must try out actions

o Big idea: Compute all averages over T using sample outcomes

# The Story So Far: MDPs and RL

## Known MDP: Offline Solution

| Goal | Technique |
|---|---|
| Compute V*, Q*, $\pi$* | Value / policy iteration |
| Evaluate a fixed policy $\pi$ | Policy evaluation |

## Unknown MDP: Model-Based

| Goal | Technique |
|---|---|
| Compute V*, Q*, $\pi$* | VI/PI on approx. MDP |
| Evaluate a fixed policy $\pi$ | PE on approx. MDP |

## Unknown MDP: Model-Free

| Goal | Technique |
|---|---|
| Compute V*, Q*, $\pi$* | Q-learning |
| Evaluate a fixed policy $\pi$ | Value Learning |

# Analogy: Expected Age

Goal: Compute expected age of UNR students

**Known P(A)**

$$E[A] = \sum_a P(a) \cdot a \qquad = 0.35 \times 20 + \dots$$

Without P(A), instead collect samples $[a_1, a_2, \dots a_N]$

**Unknown P(A): "Model Based"**

Why does this work? Because eventually you learn the right model.

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$

$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

**Unknown P(A): "Model Free"**

$$E[A] \approx \frac{1}{N} \sum_i a_i$$

Why does this work? Because samples appear with the right frequencies.

# Sample-Based Policy Evaluation?

○ We want to improve our estimate of V (utility) by computing these averages:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

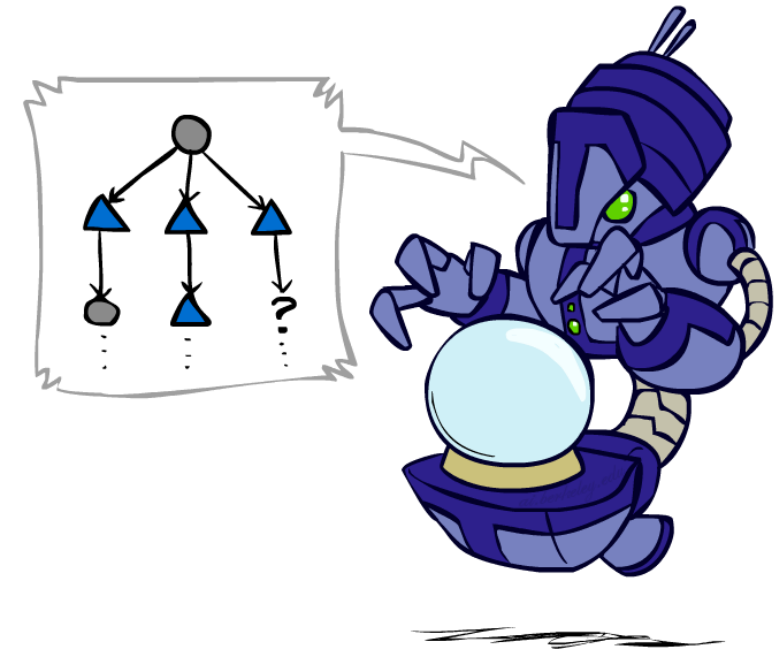○ Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_1 = R(s, \pi(s), s_1') + \gamma V_k^{\pi}(s_1')$$

$$sample_2 = R(s, \pi(s), s_2') + \gamma V_k^{\pi}(s_2')$$

$$\cdots$$

$$sample_n = R(s, \pi(s), s_n') + \gamma V_k^{\pi}(s_n')$$

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_i sample_i$$
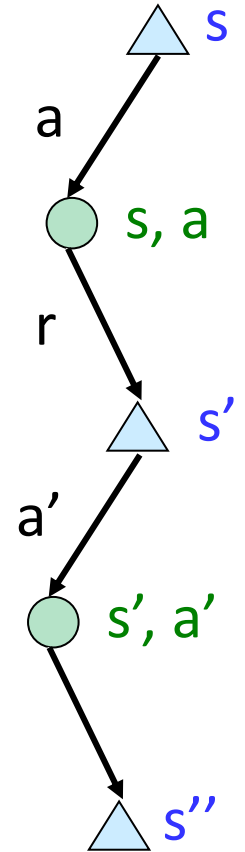
# Model-Free Learning

o Model-free (temporal difference) learning

   o Experience world through episodes

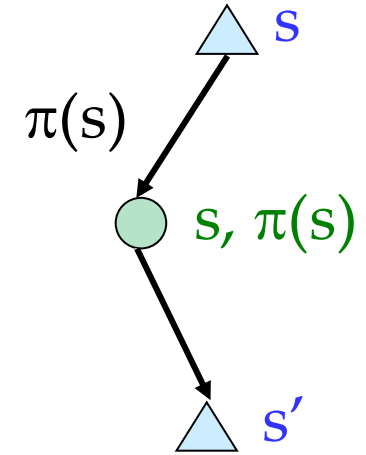$$(s, a, r, s', a', r', s'', a'', r'', s'''' \dots)$$

   o Update estimates each transition $(s, a, r, s')$

   o Over time, updates will mimic Bellman updates

# Temporal Difference Learning

o Temporal difference learning of values
  o Policy still fixed, still doing evaluation!
  o Move values toward value of whatever successor occurs: running average

$\pi(s)$

$s$

$s, \pi(s)$

$s'$

Sample of V(s): $sample = R(s, \pi(s), s') + \gamma V^{\pi}(s')$

Update to V(s): $V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + (\alpha)sample$

# Example: Temporal Difference Learning

## States



*Assume: γ = 1, α = 1/2*

## Observed Transitions

B, east, C, -2

C, east, D, -2



$$V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha \left[ R(s, \pi(s), s') + \gamma V^{\pi}(s') \right]$$

# Approximating Values through Samples

- Policy Evaluation:

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

- Value Iteration:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s')\right]$$

- Q-Value Iteration:

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

# Q-Learning

○ Q-Learning: sample-based Q-value iteration

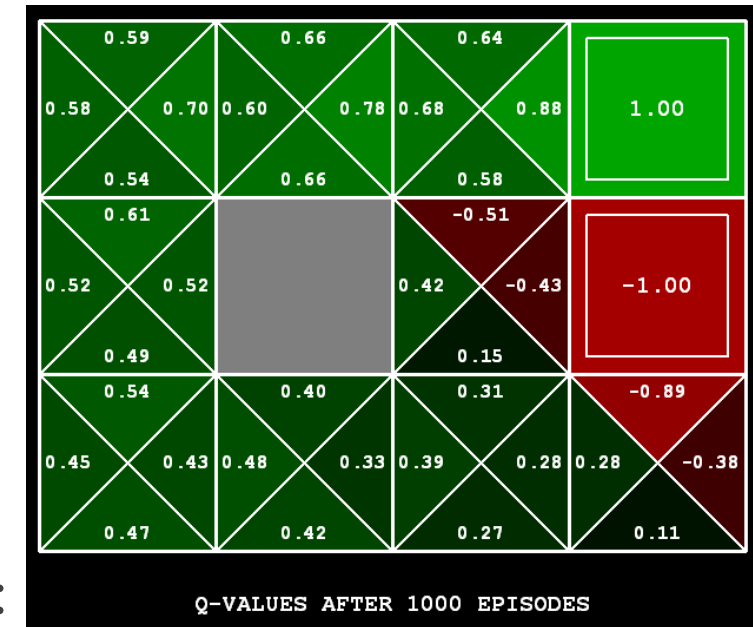$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

○ Learn Q(s,a) values as you go

  ○ Receive a sample (s,a,s',r)
  ○ Consider your old estimate: $Q(s,a)$
  ○ Consider your new sample estimate:

$$sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$   no longer policy evaluation!
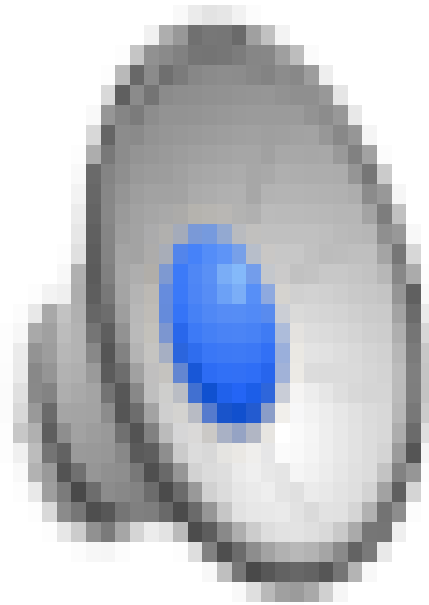
  ○ Incorporate the new estimate into a running average:

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha) \left[sample\right]$$
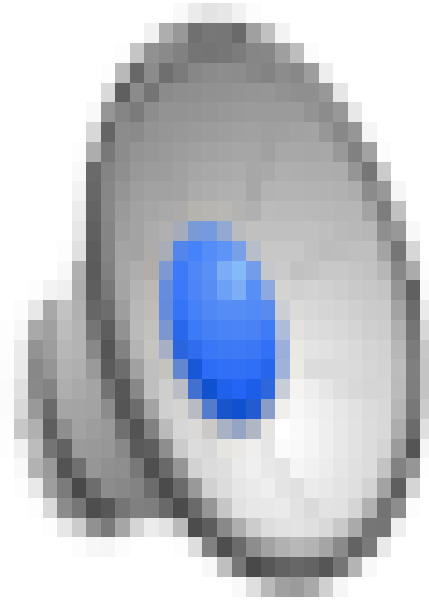


Q-VALUES AFTER 1000 EPISODES

[Demo: Q-learning – gridworld (L10D2)]
[Demo: Q-learning – crawler (L10D3)]
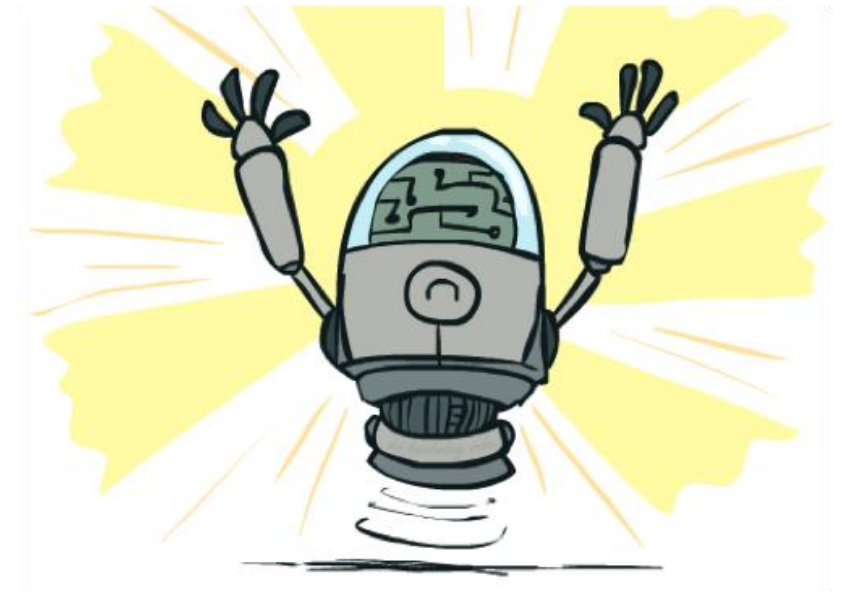
# Video of Demo Q-Learning -- Gridworld

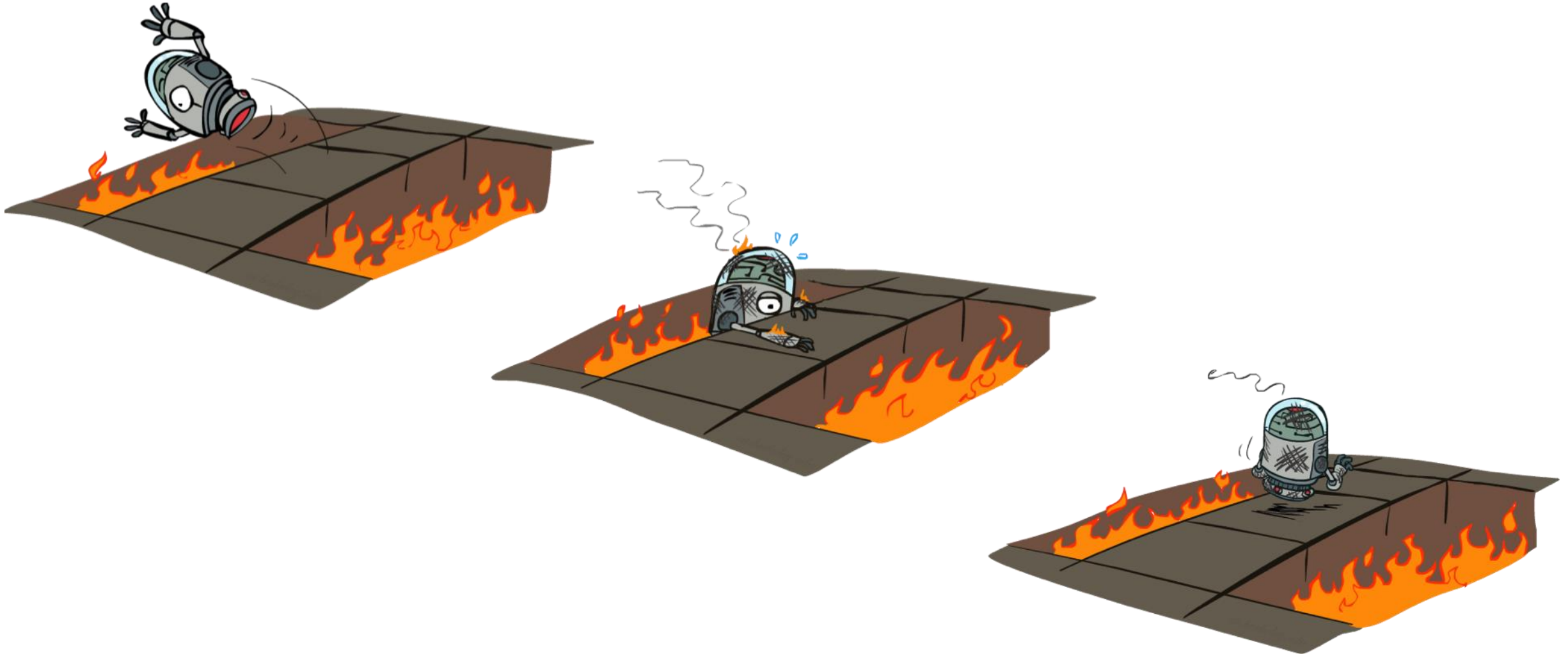# Video of Demo Q-Learning -- Crawler

# Q-Learning Properties

o Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!

o This is called <span style="color:red">off-policy learning</span>

o Caveats:
   o You have to explore enough
   o You have to eventually make the learning rate small enough
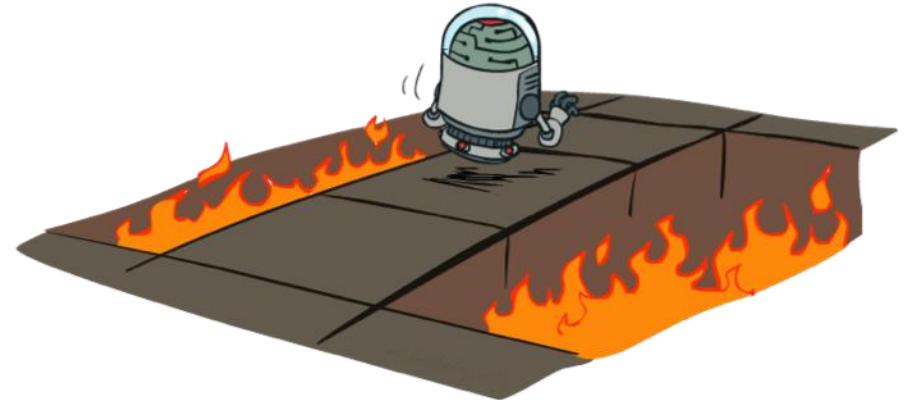   o … but not decrease it too quickly
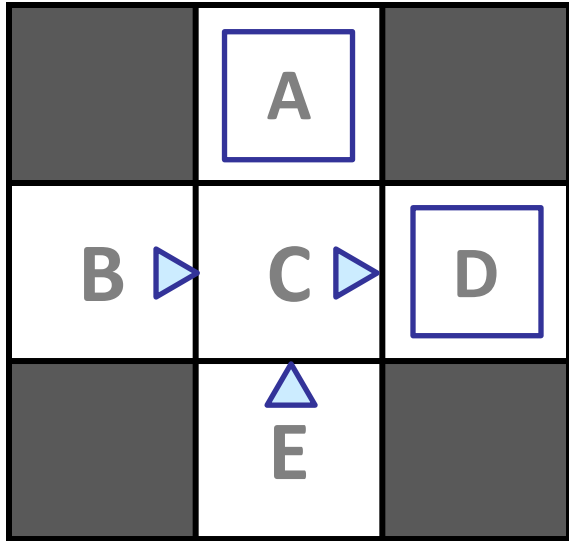
# Active Reinforcement Learning

# Model-Free Learning

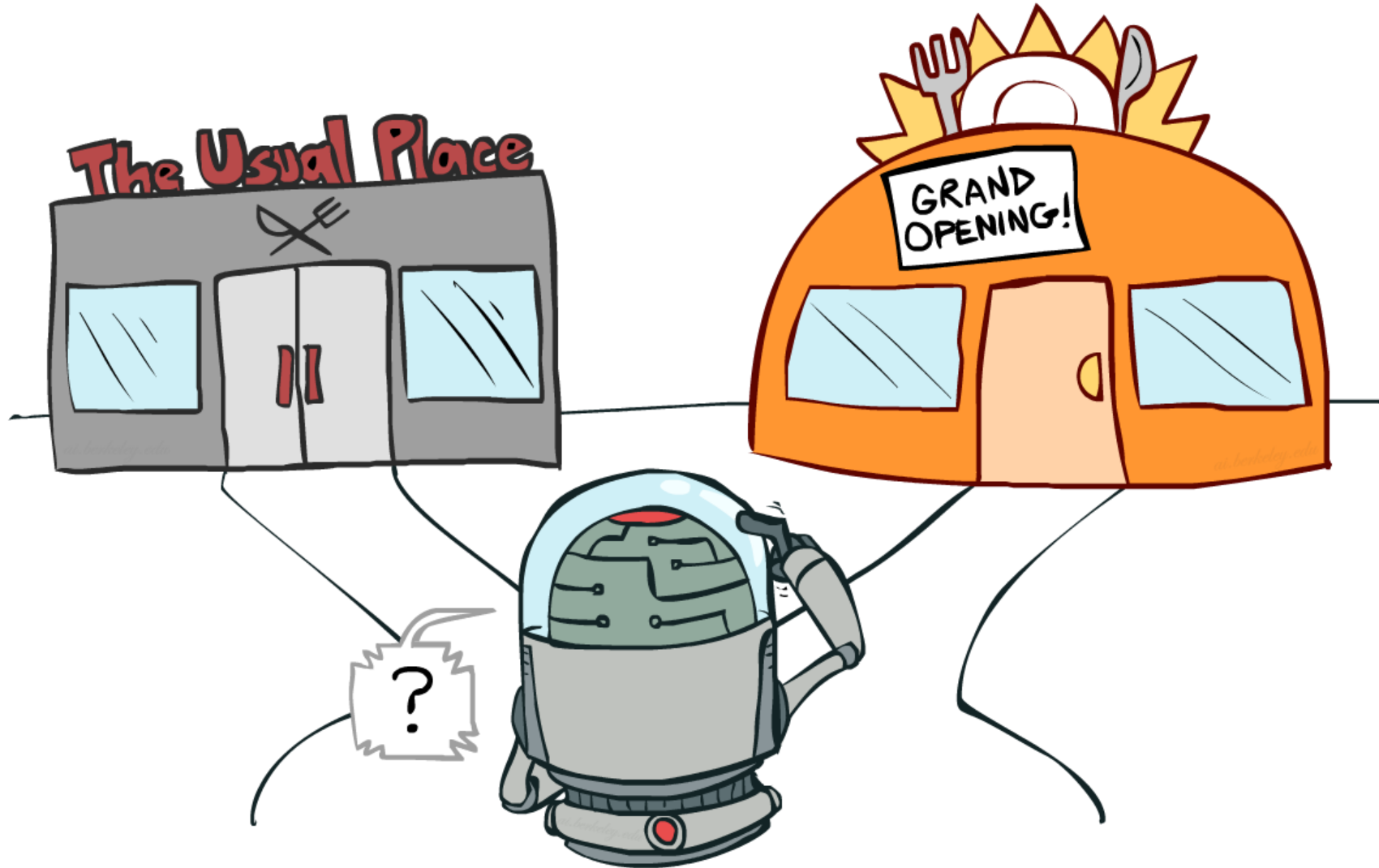o   act according to current optimal (based on Q-Values)
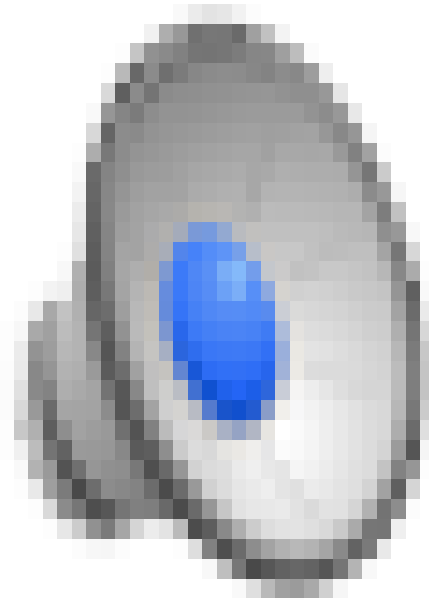o   but also explore…

# Model-Based Learning

act according to current optimal
also explore!

# Exploration vs. Exploitation

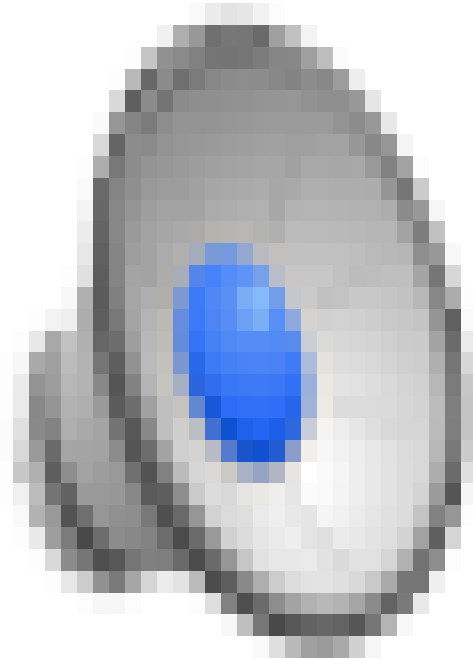# Video of Demo Q-learning – Manual Exploration – Bridge Grid

# How to Explore?

o **Several schemes for forcing exploration**
  - o Simplest: random actions (ε-greedy)
    - o Every time step, flip a coin
    - o With (small) probability ε, act randomly
    - o With (large) probability 1-ε, act on current policy

  - o Problems with random actions?
    - o You do eventually explore the space, but keep thrashing around once learning is done
    - o One solution: lower ε over time
    - o Another solution: exploration functions

# Video of Demo Q-learning – Epsilon-Greedy – Crawler

# Exploration Functions

o **When to explore?**
  - o Random actions: explore a fixed amount
  - o Better idea: explore areas whose badness is not (yet) established, eventually stop exploring

o **Exploration function**
  - o Takes a value estimate u and a visit count n, and returns an optimistic utility, e.g. $f(u, n) = u + k/n$
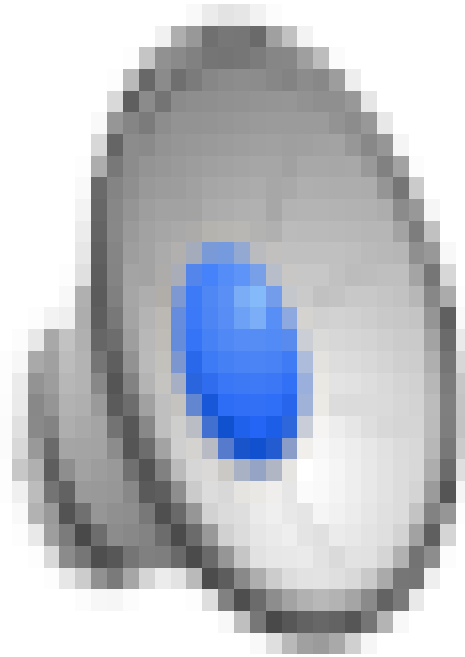
  Regular Q-Update: $\quad Q(s, a) \leftarrow_\alpha R(s, a, s') + \gamma \max_{a'} Q(s', a')$

  Modified Q-Update: $Q(s, a) \leftarrow_\alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$

  - o Note: this propagates the "bonus" back to states that lead to unknown states as well!
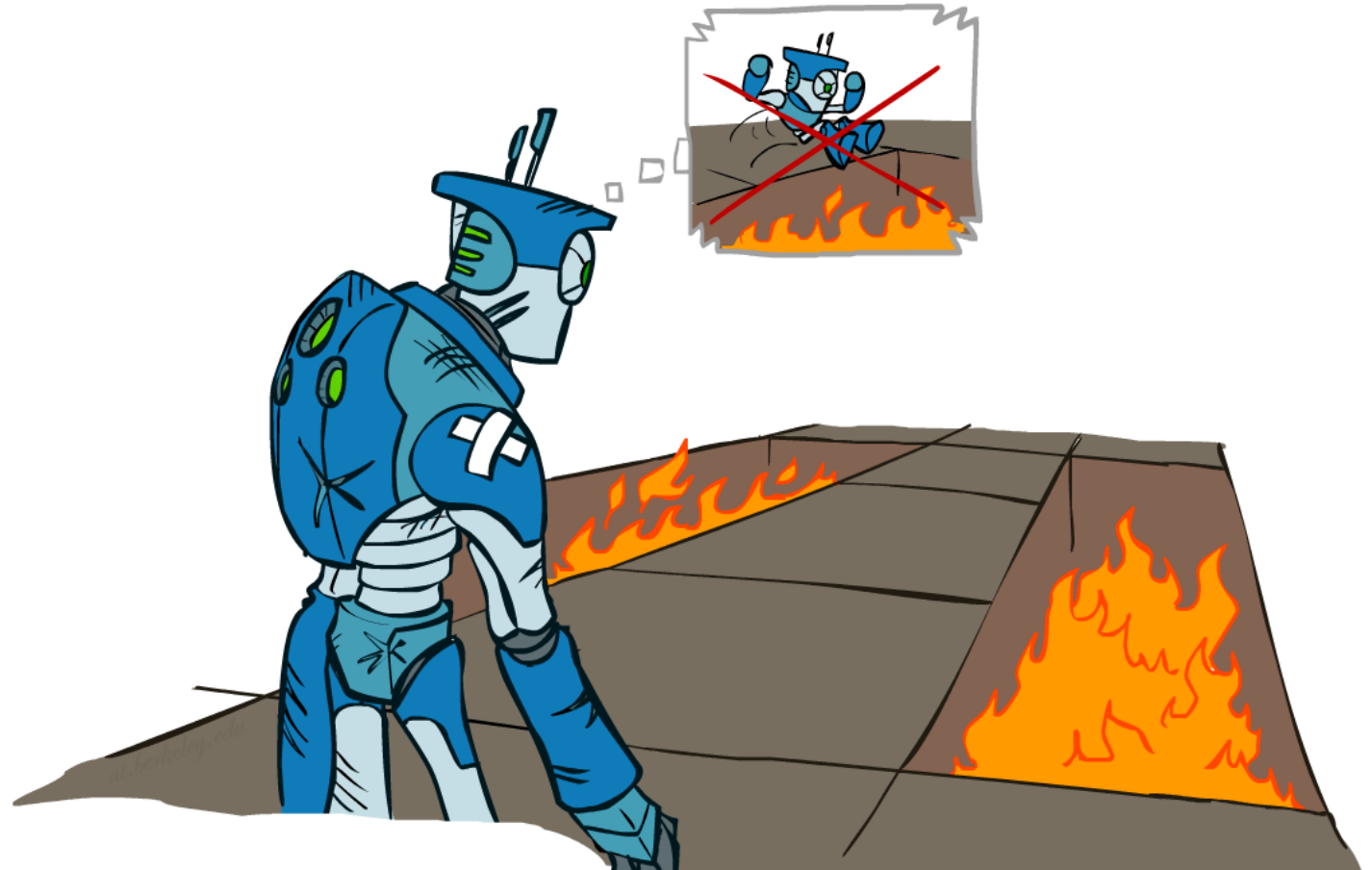
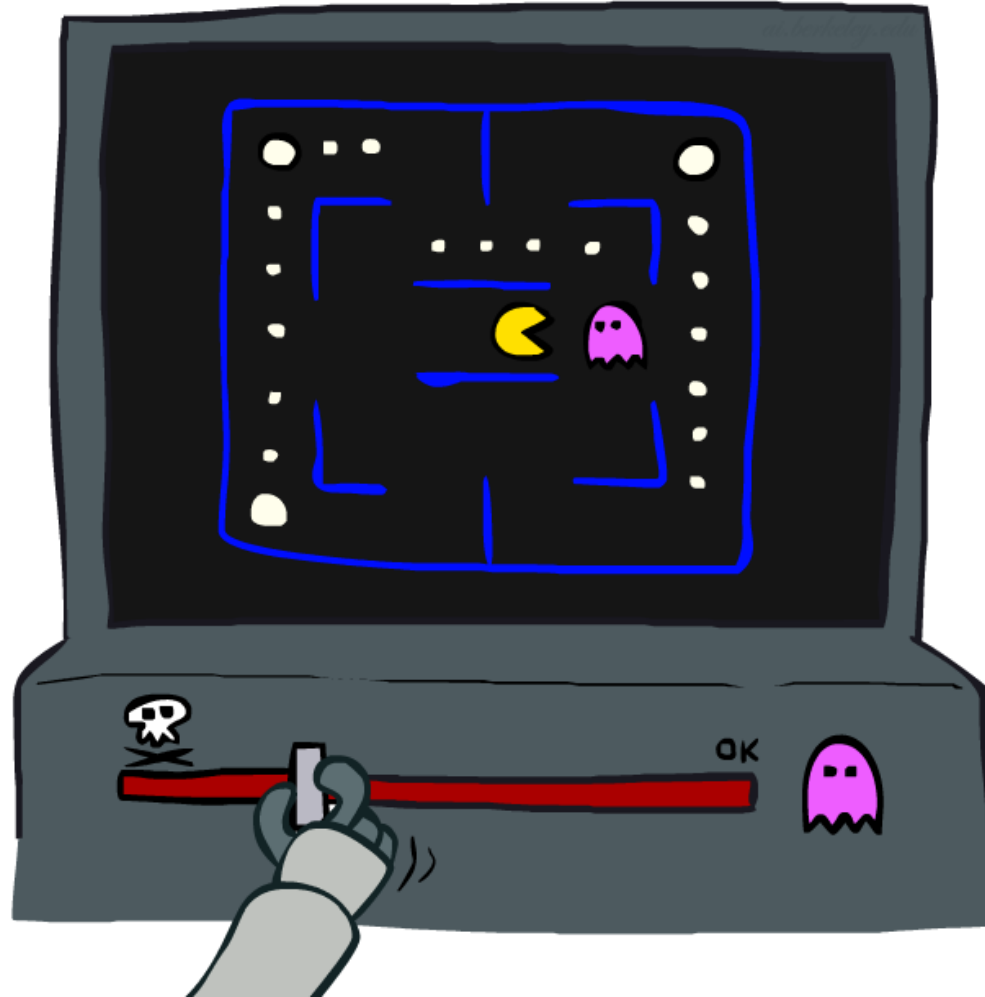# Video of Demo Q-learning – Exploration Function – Crawler

# Regret

o Even if you learn the optimal policy, you still make mistakes along the way!

o Regret is a measure of your total mistake cost: the difference between your (expected) rewards, including youthful suboptimality, and optimal (expected) rewards

o Minimizing regret goes beyond learning to be optimal – it requires optimally learning to be optimal

o Example: random exploration and exploration functions both end up optimal, but random exploration has higher regret
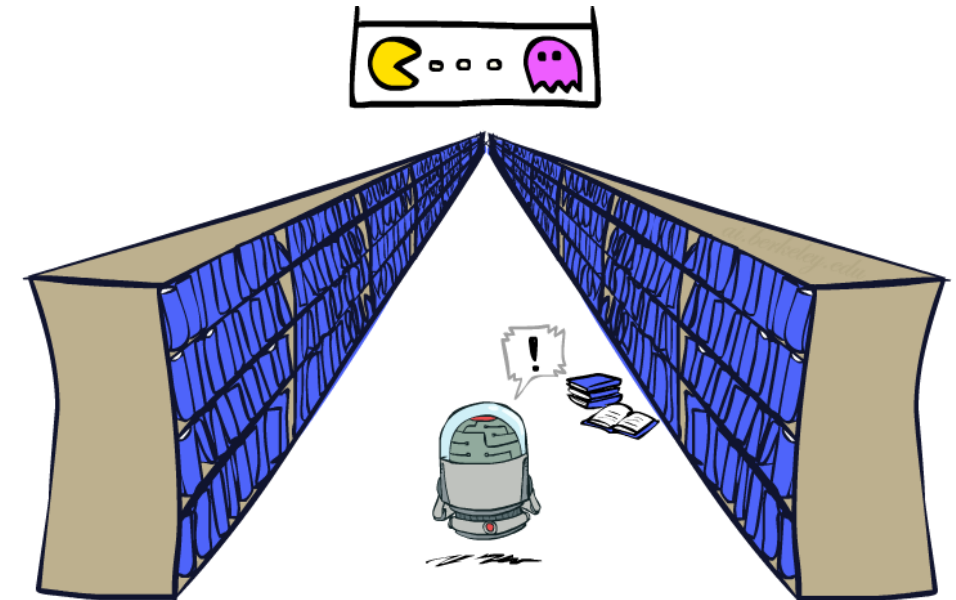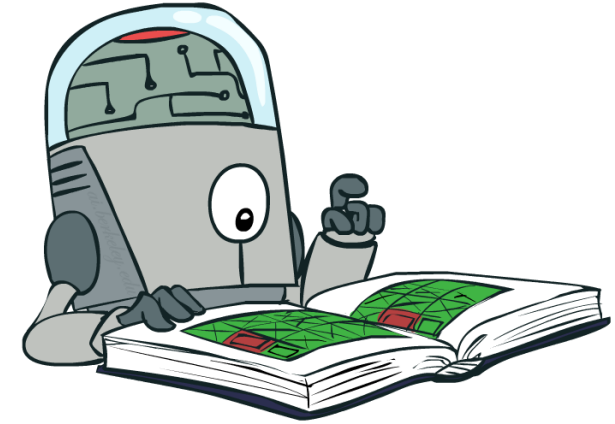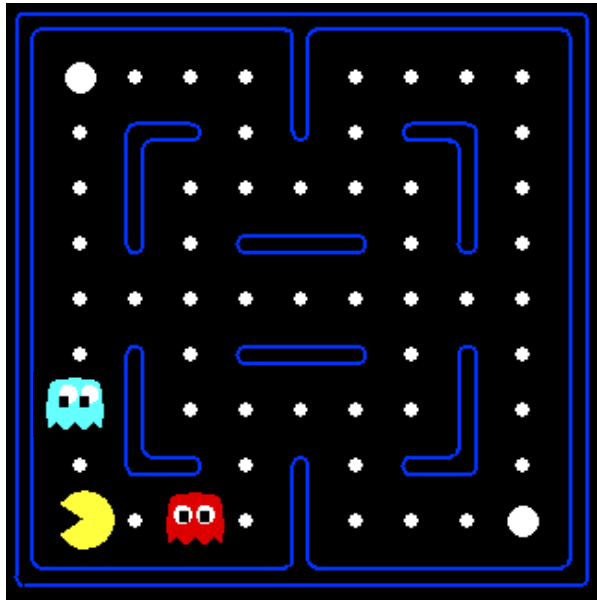
# Approximate Q-Learning

# Generalizing Across States

o Basic Q-Learning keeps a table of all q-values

o In realistic situations, we cannot possibly learn about every single state!
  o Too many states to visit them all in training
  o Too many states to hold the q-tables in memory

o Instead, we want to generalize:
  o Learn about some small number of training states from experience
  o Generalize that experience to new, similar situations
  o This is a fundamental idea in machine learning, and we'll see it over and over again
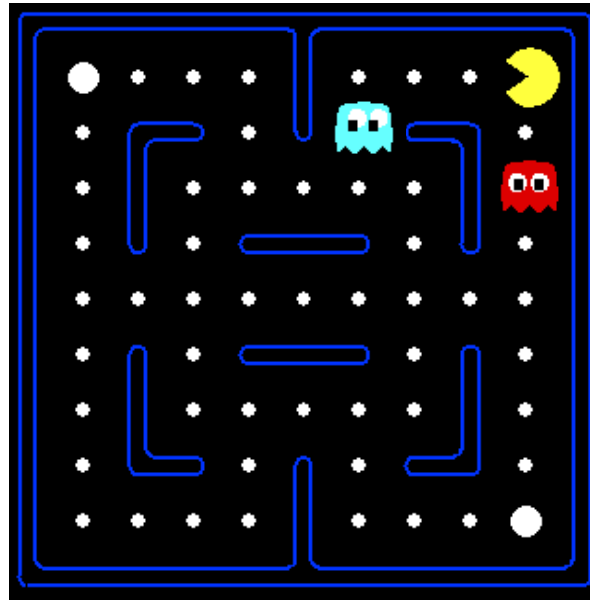
[demo – RL pacman]

# Example: Pacman

Let's say we discover through experience that this state is bad:
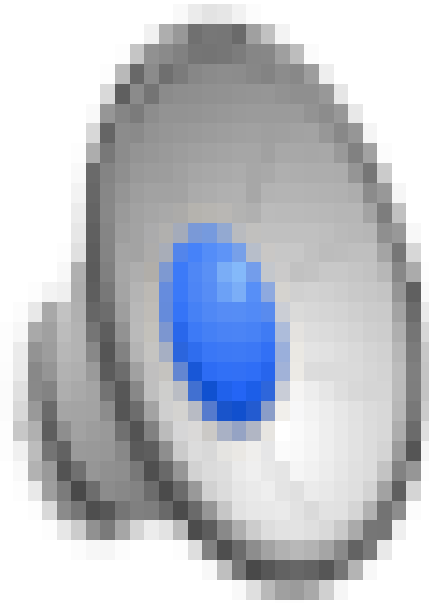
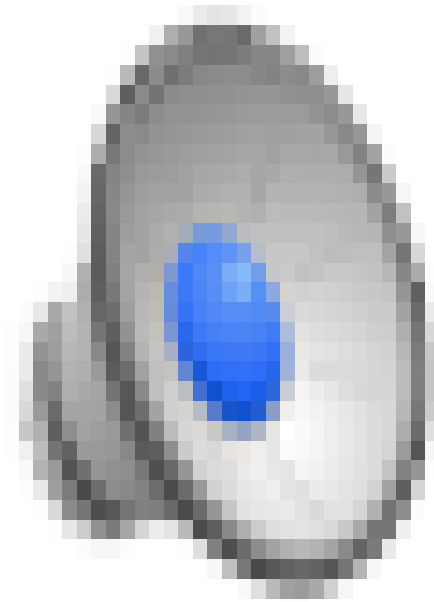In naïve q-learning, we know nothing about this state:

Or even this one!

# Video of Demo Q-Learning Pacman – Tiny – Watch All

# Video of Demo Q-Learning Pacman – Tiny – Silent Train

# Video of Demo Q-Learning Pacman – Tricky – Watch All

# Feature-Based Representations

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

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

# Linear Value Functions

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

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

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

o Advantage: our experience is summed up in a few powerful numbers

o Disadvantage: states may share features but actually be very different in value!

# Error (least square)

What is the partial derivative of function

$$g(w_1, w_2) = \frac{1}{2}(y - (w_1 f_1(x) + w_2 f_2(x)))^2$$

w.r.t. $w_1$, i.e., $\frac{\partial g(w_1, w_2)}{\partial w_1}$?

Assume $y$ is a constant and $f$ is a known function that maps a vector in $\mathbb{R}^n$ to a scalar

A: $f_1(x)$

B: $y - (w_1 f_1(x) + w_2 f_2(x))$

C: $w_1 f_1(x) + w_2 f_2(x) - y$

D: $(w_1 f_1(x) + w_2 f_2(x) - y) f_1(x)$

# Error (least square)

What is the partial derivative of function

$$g(w_1, w_2) = \frac{1}{2}(y - (w_1 f_1(x) + w_2 f_2(x)))^2$$

w.r.t. $w_1$, i.e., $\frac{\partial g(w_1, w_2)}{\partial w_1}$?

Assume $y$ is a constant and $f$ is a known function that maps a vector in $\mathbb{R}^n$ to a scalar

A: $f_1(x)$

B: $y - (w_1 f_1(x) + w_2 f_2(x))$

C: $w_1 f_1(x) + w_2 f_2(x) - y$

D: $(w_1 f_1(x) + w_2 f_2(x) - y) f_1(x)$

Let $h(w_1, w_2) = y - (w_1 f_1(x) + w_2 f_2(x))$

Then $g(w_1, w_2) = \frac{1}{2}\left(h(w_1, w_2)\right)^2$

$$\frac{\partial g(w_1, w_2)}{\partial w_1} = \frac{\partial g(w_1, w_2)}{\partial h(w_1, w_2)} \frac{\partial h(w_1, w_2)}{\partial w_1}$$

$$= \frac{1}{2} \times 2 \times h(w_1, w_2) \times (-f_1(x))$$

$$= -h(w_1, w_2) f_1(x)$$

34

# Updating a linear value function

Original Q-learning: Update Q values directly (stored in a table)

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

Latest sample      Previous estimate

Difference

Can be viewed as trying to reduce prediction error at s, a:

$$Q(s,a) \leftarrow Q(s,a) - \alpha \nabla Error \qquad Error = \frac{1}{2}\left(sample - Q(s,a)\right)^2$$

Approximate Q-Learning with Linear Q-Value Function:

$$Q_w(s,a) = w_1 f_1(s,a) + \dots + w_n f_n(s,a)$$

Update weights to reduce prediction error at s, a:

$$w_i \leftarrow w_i - \alpha \frac{\partial Error(w_1, w_2, \dots, w_n)}{\partial w_i} \qquad Error(w) = \frac{1}{2}\left(sample - Q_w(s,a)\right)^2$$

35

# Updating a linear value function

$$Q_w(s, a) = w_1 f_1(s, a) + \ldots + w_n f_n(s, a)$$

$$w_i \leftarrow w_i - \alpha \frac{\partial Error(w_1, w_2, \ldots, w_n)}{\partial w_i} \qquad Error(w) = \frac{1}{2} \left(sample - Q_w(s, a)\right)^2$$

$$\frac{\partial Error(w)}{\partial w_i} = (Q_w(s, a) - sample) \frac{\partial Q_w(s, a)}{\partial w_i}$$
$$= (Q_w(s, a) - sample) f_i(s, a)$$

Final Update Rule for Approximate Q-Learning with Linear Q-Value Function:

$$\boxed{w_i} \leftarrow w_i + \alpha \left(r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a)\right) \boxed{f_i(s, a)}$$

Original Q-Learning Update Rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$
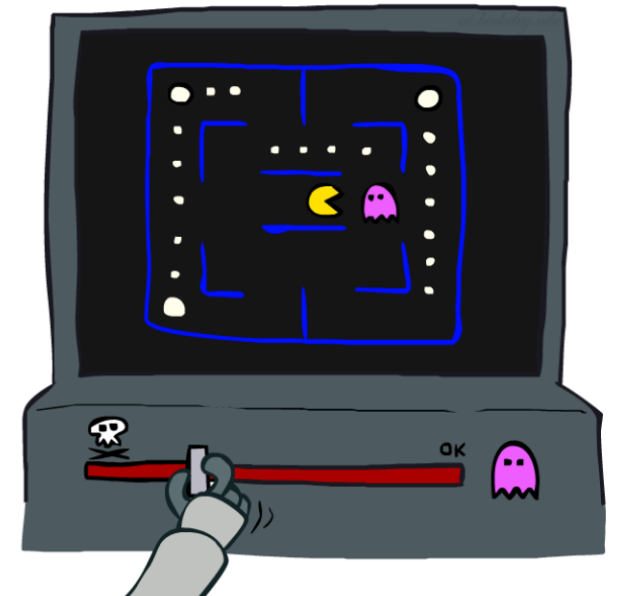
# Approximate Q-Learning

Update Rule for Approximate Q-Learning with Linear Q-Value Function:

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right) f_i(s, a)$$

Qualitative justification:

- Pleasant surprise: increase weights on +valued features, decrease on − ones
  - As a result, $Q_w$ increased for states with the same (similar) features too. Will now prefer all states with that state's features.
- Unpleasant surprise: decrease weights on +valued features, increase on − ones
  - Disprefer all states with that state's features

# Approximate Q-Learning

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

○ 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

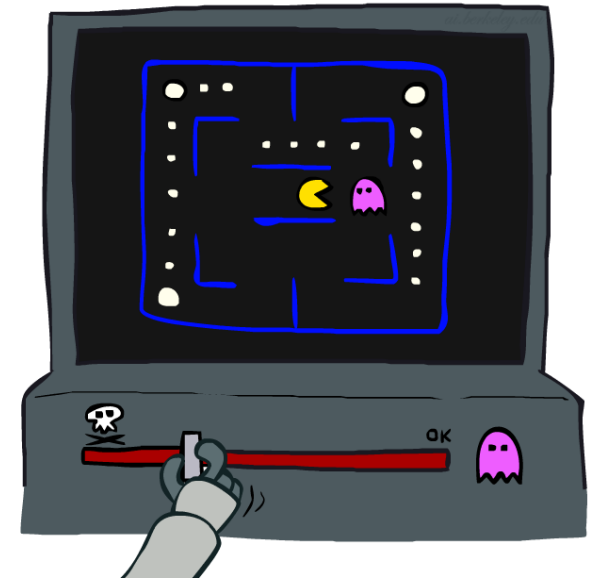$$w_i \leftarrow w_i + \alpha \,[\text{difference}]\, f_i(s,a)$$ Approximate Q's

○ Intuitive interpretation:
   ○ Adjust weights of active features
   ○ E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

○ Formal justification: least squares

# What if non-linear value function

Update Rule for Q-Learning:

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

Update Rule for Approximate Q-Learning with Linear Q-Value Function:

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s',a') - Q_w(s,a) \right) f_i(s,a)$$

Update Rule for Approximate Q-Learning with differentiable Q-function $Q_w(s,a)$:

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s',a') - Q_w(s,a) \right) \frac{\partial Q_w(s,a)}{\partial w_i}$$

$$\text{If } Q_w(s,a) = w_1 f_1(s,a) + \ldots + w_n f_n(s,a)$$

$$\frac{\partial Q_w(s,a)}{\partial w_i} = f_1(s,a)$$

# What if non-linear value function

Update Rule for Approximate Q-Learning with Q-function $Q_w(s, a)$:

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right) \frac{\partial Q_w(s, a)}{\partial w_i}$$

Why?

$$w_i \leftarrow w_i - \alpha \frac{\partial Error(w_1, w_2, \ldots, w_n)}{\partial w_i} \qquad Error(w) = \frac{1}{2} \left( sample - Q_w(s, a) \right)^2$$

$$\frac{\partial Error(w)}{\partial w_i} = (Q_w(s, a) - sample) \frac{\partial Q_w(s, a)}{\partial w_i}$$

$$w_i - \alpha \frac{\partial Error(w_1, w_2, \ldots, w_n)}{\partial w_i} = w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right) \frac{\partial Q_w(s, a)}{\partial w_i}$$

# What if non-linear value function

Update Rule for Approximate Q-Learning with Q-function $Q_w(s, a)$:

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right) \frac{\partial Q_w(s, a)}{\partial w_i}$$

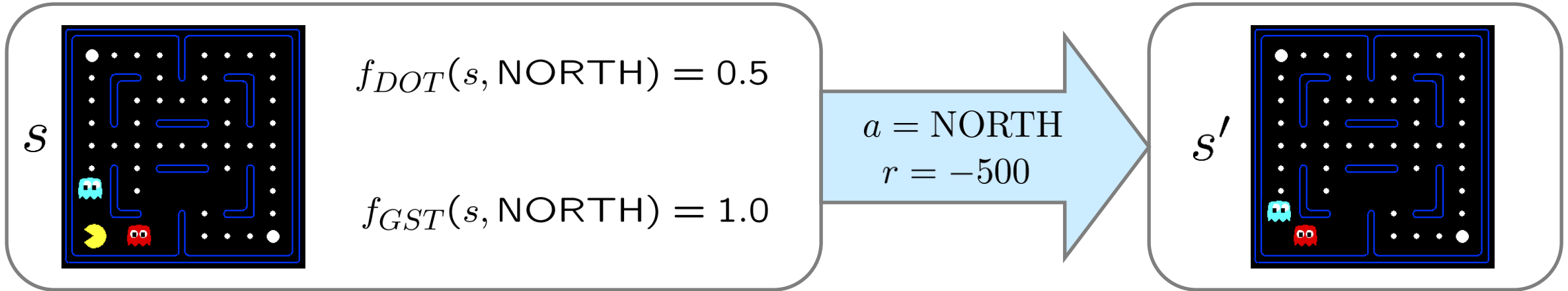Example: $Q_w(s, a) = \exp(w_1 f_1(s, a) + w_2 f_2(s, a) + \cdots + w_n f_n(s, a))$

$$\frac{\partial Q_w(s, a)}{\partial w_i} = \exp(w_1 f_1(s, a) + w_2 f_2(s, a) + \cdots + w_n f_n(s, a)) f_i(s, a)$$

$$= Q_w(s, a) f_i(s, a)$$

Update Rule:

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right) Q_w(s, a) f_i(s, a)$$
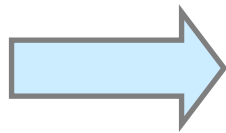
# Example: Q-Pacman

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

$s$

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

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

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

$s'$

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

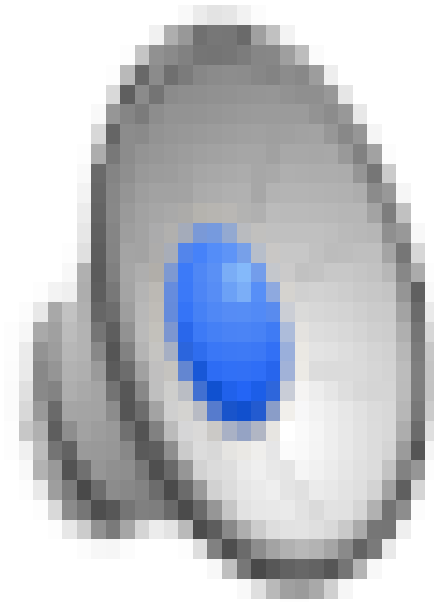$r + \gamma \max_{a'} Q(s', a') = -500 + 0$

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

difference $= -501$

$w_{DOT} \leftarrow 4.0 + \alpha[-501]0.5$
$w_{GST} \leftarrow -1.0 + \alpha[-501]1.0$
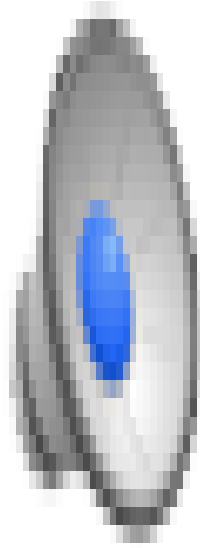
$$Q(s,a) = 3.0f_{DOT}(s,a) - 3.0f_{GST}(s,a)$$

[Demo: approximate Q-learning pacman (L11D10)]

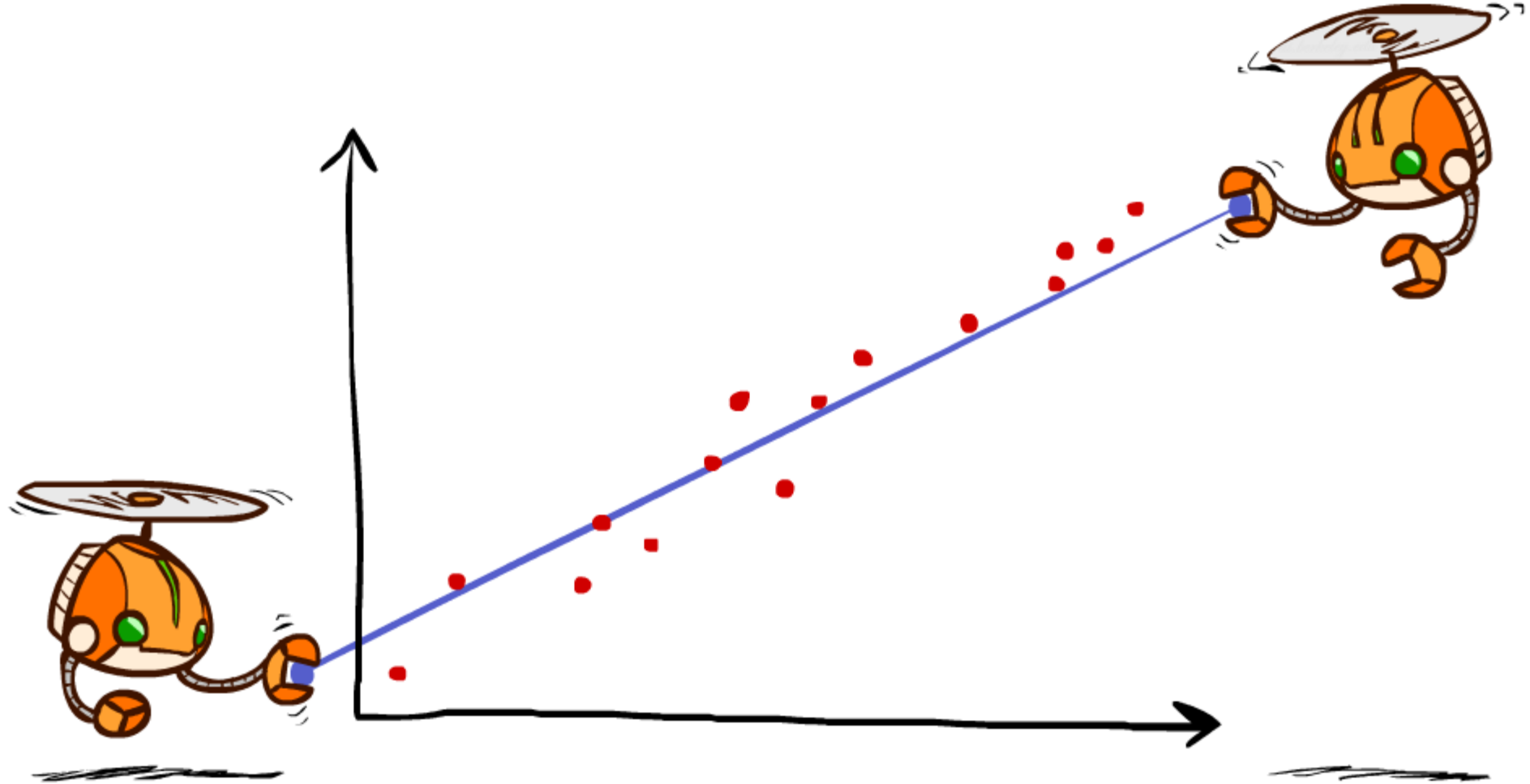# Video of Demo Approximate Q-Learning -- Pacman
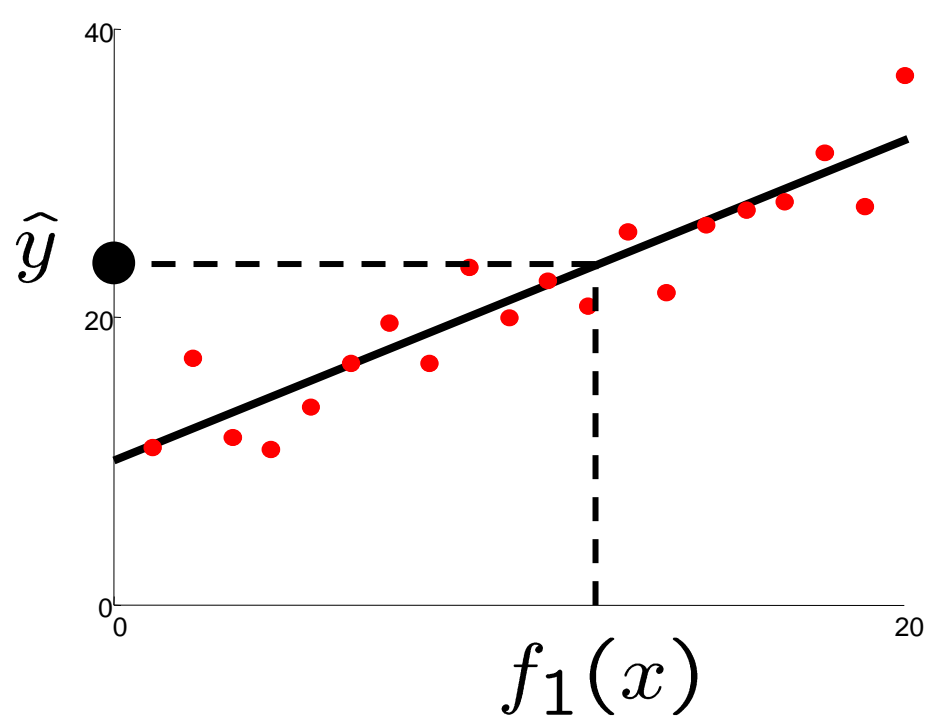
# DeepMind Atari (©Two Minute Lectures) approximate Q-learning with neural nets
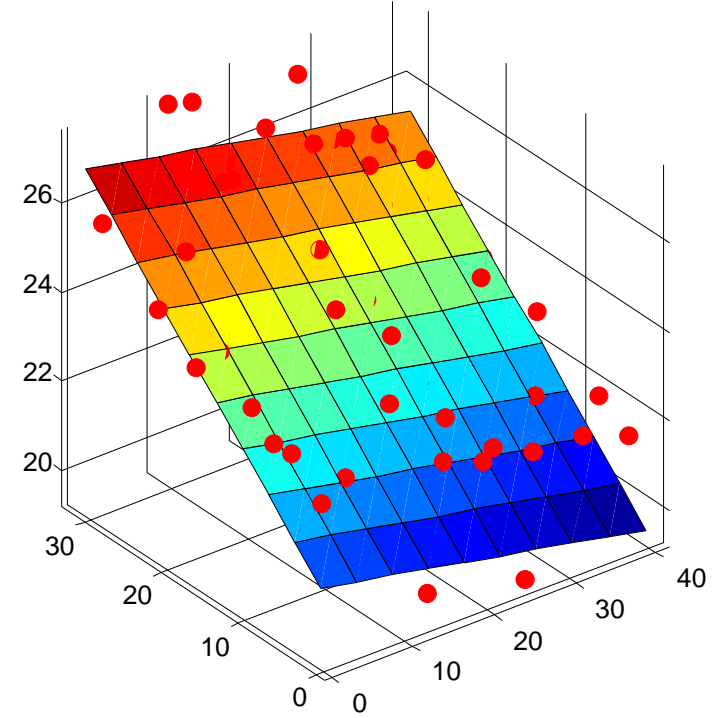
# Q-Learning and Least Squares

# Linear Approximation: Regression



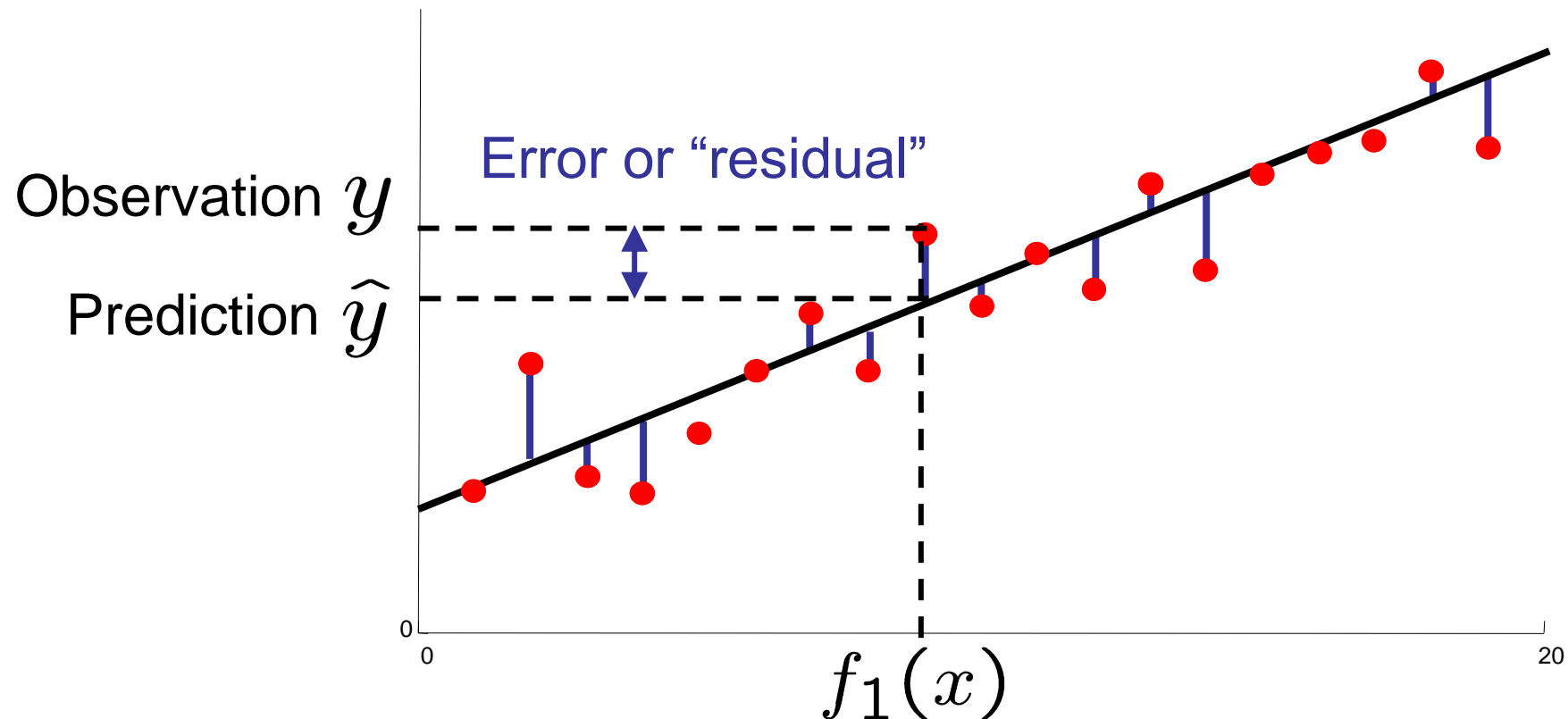Prediction:
$$\hat{y} = w_0 + w_1 f_1(x)$$

Prediction:
$$\hat{y}_i = w_0 + w_1 f_1(x) + w_2 f_2(x)$$
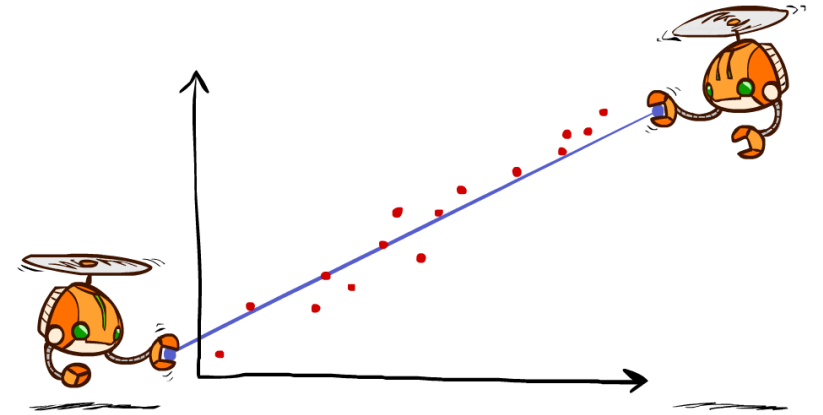
# Optimization: Least Squares

$$\text{total error} = \sum_i (y_i - \widehat{y}_i)^2 = \sum_i \left(y_i - \sum_k w_k f_k(x_i)\right)^2$$



Observation $y$

Prediction $\widehat{y}$

Error or "residual"

$f_1(x)$

# Minimizing Error

Imagine we had only one point x, with features f(x), target value y, and weights w:

$$\text{error}(w) = \frac{1}{2}\left(y - \sum_k w_k f_k(x)\right)^2$$

$$\frac{\partial\ \text{error}(w)}{\partial w_m} = -\left(y - \sum_k w_k f_k(x)\right) f_m(x)$$

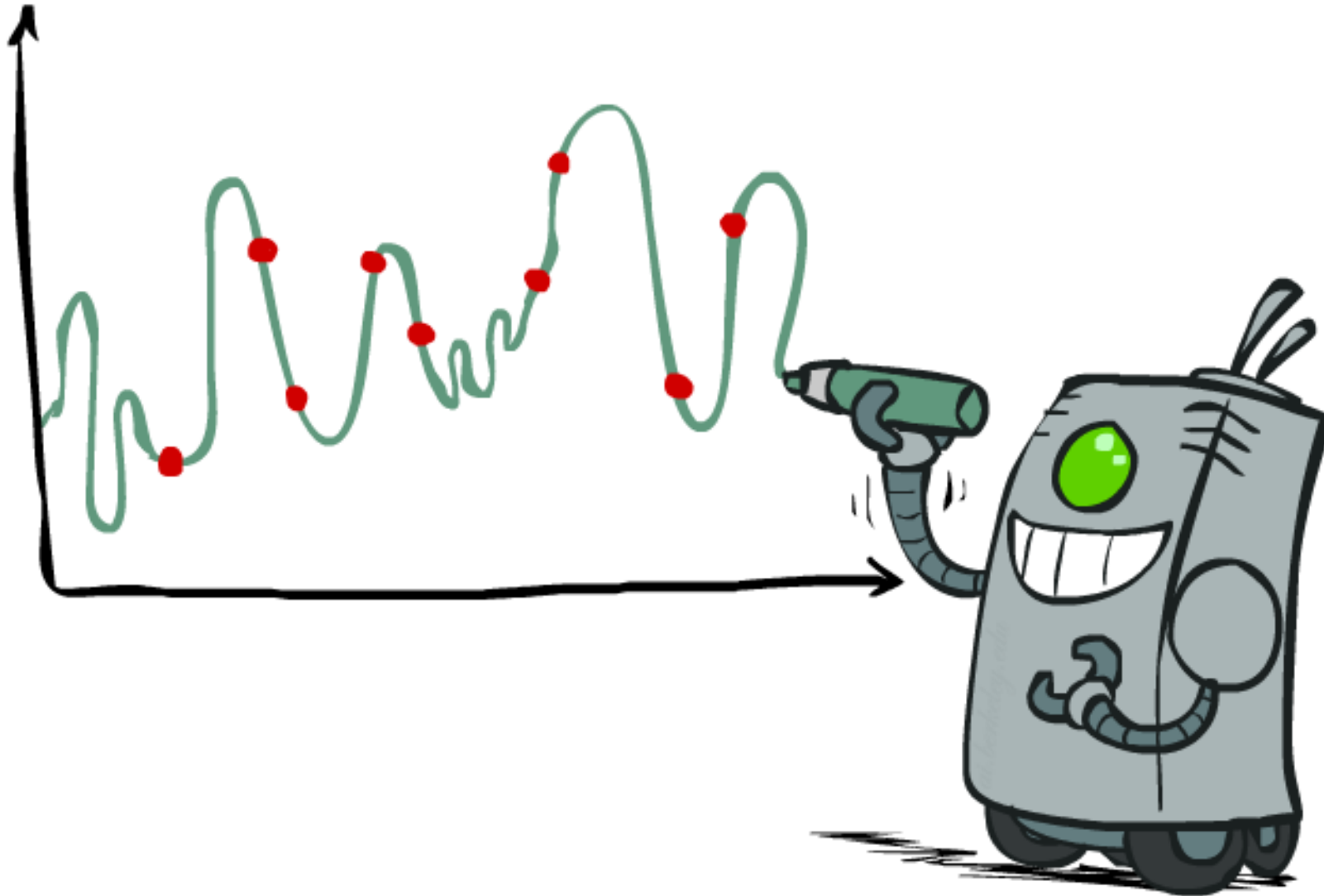$$w_m \leftarrow w_m + \alpha\left(y - \sum_k w_k f_k(x)\right) f_m(x)$$

Approximate q update explained:

$$w_m \leftarrow w_m + \alpha\left[r + \gamma \max_a Q(s', a') - Q(s, a)\right] f_m(s, a)$$
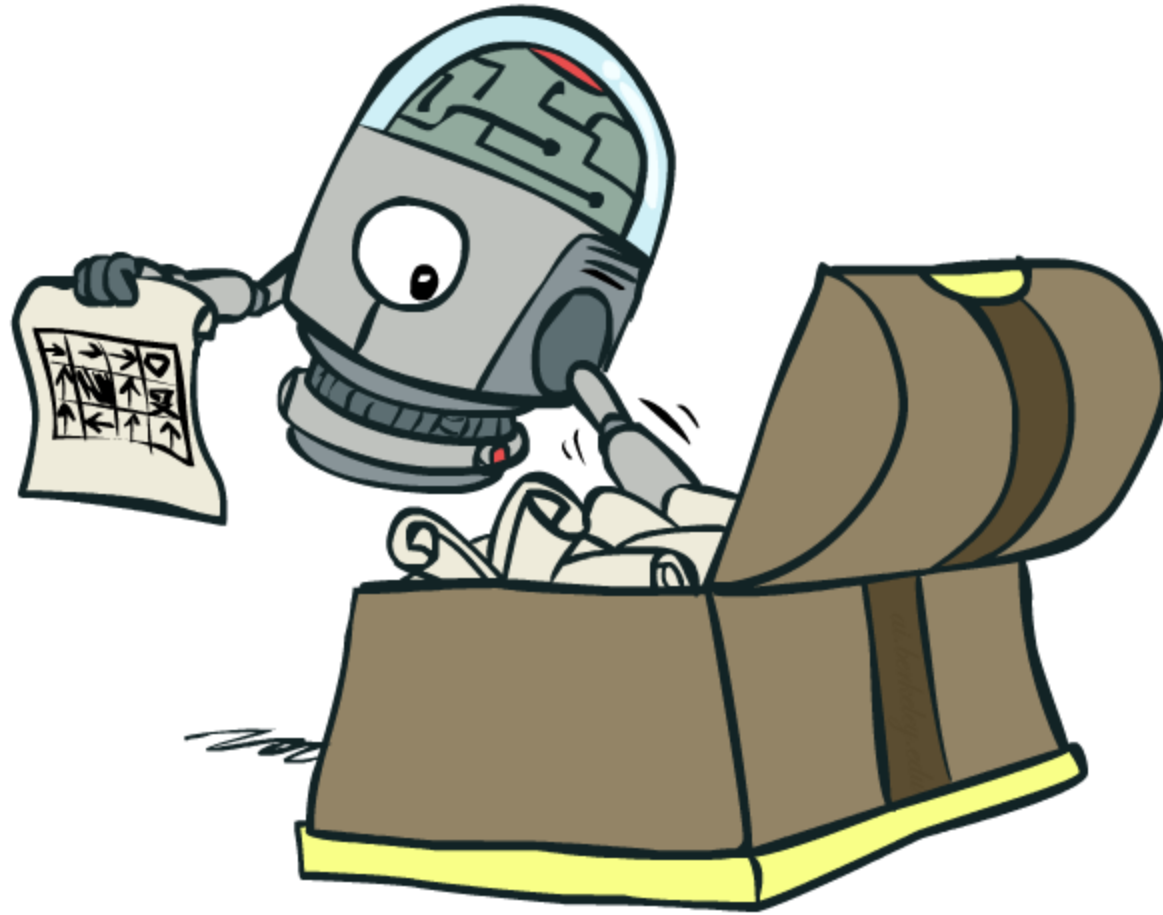
"target"      "prediction"

# Overfitting: Why Limiting Capacity Can Help

# Policy Search

# Policy Search

o Simplest policy search:

    o Start with an initial linear value function or Q-function

    o Nudge each feature weight up and down and see if your policy is better than before

o Problems:

    o How do we tell the policy got better?

    o Need to run many sample episodes!

    o If there are a lot of features, this can be impractical

o Better methods exploit lookahead structure, sample wisely, change multiple parameters…

# The Story So Far: MDPs and RL

## Known MDP: Offline Solution

| Goal | Technique |
|---|---|
| Compute V*, Q*, π* | Value / policy iteration |
| Evaluate a fixed policy π | Policy evaluation |

## Unknown MDP: Model-Based

| Goal | *use features to generalize* | Technique |
|---|---|---|
| Compute V*, Q*, π* | | VI/PI on approx. MDP |
| Evaluate a fixed policy π | | PE on approx. MDP |

## Unknown MDP: Model-Free

| Goal | *use features to generalize* | Technique |
|---|---|---|
| Compute V*, Q*, π* | | Q-learning |
| Evaluate a fixed policy π | | Value Learning |