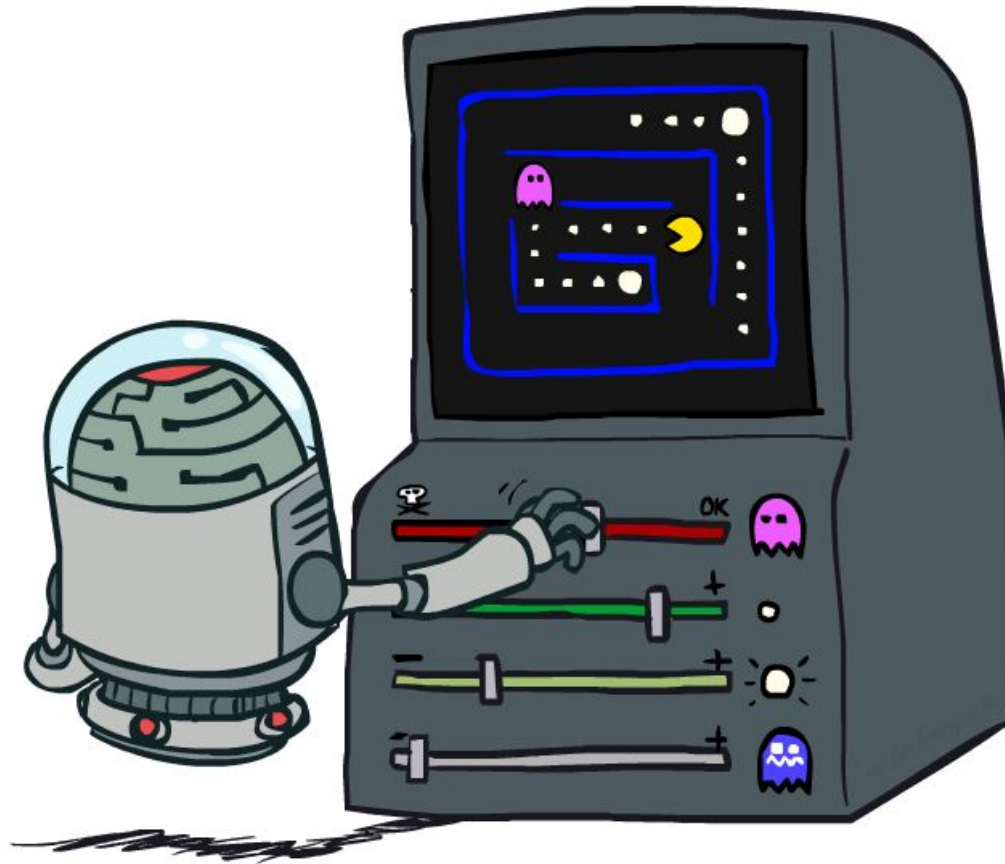# CS 188: Artificial Intelligence
## Reinforcement Learning II

Instructor: Stuart Russell and Dawn Song, University of California, Berkeley

# Recap: Reinforcement Learning

- Still assume a Markov decision process (MDP):
  - A set of states s $\in$ S
  - A set of actions (per state) A
  - A model T(s,a,s')
  - A reward function R(s,a,s')
- Still looking for a policy $\pi$(s)

- New twist: don't know T or R
  - I.e. we don't know which states are good or what the actions do
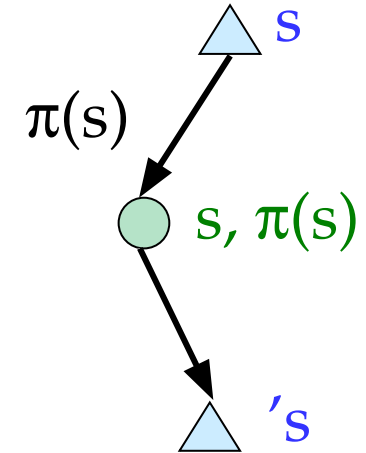  - Must actually try actions and states out to learn

# Recap: Reinforcement Learning

o Passive reinforcement learning:

   o A passive learning agent has a fixed policy that determines its behavior

o Model-based learning:

   o Learn an approximate MDP model based on experiences

o Model-free learning:

   o Do not learn an explicit MDP model

# Recap: Temporal Difference Learning

o Big idea: learn from every experience!
  o Update V(s) each time we experience a transition (s, a, s', r)
  o Likely outcomes s' will contribute updates more often

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

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

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

Same update:   $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$
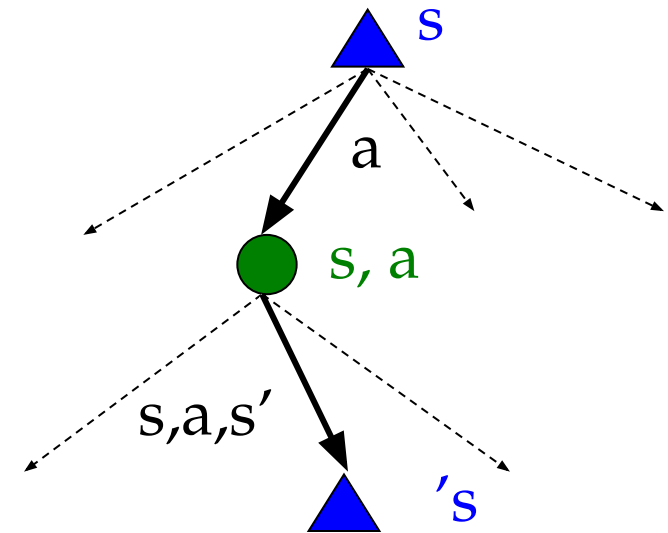
s

$\pi(s)$

s, $\pi(s)$

s'

# Recap: Problems with TD Value Learning

○ TD value leaning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages

○ However, if we want to turn values into a (new) policy, we're sunk:

$$\pi(s) = \arg\max_a Q(s, a)$$

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

○ Idea: learn Q-values, not values
○ Makes action selection model-free too!

s

a

s, a

s,a,s'

s'

# Detour: Q-Value Iteration

o Value iteration: find successive (depth-limited) values
  o Start with $V_0(s) = 0$, which we know is right
  o Given $V_k$, calculate the depth k+1 values for all states:

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

o But Q-values are more useful, so compute them instead
  o Start with $Q_0(s,a) = 0$
  o Given $Q_k$, calculate the depth k+1 q-values for all q-states:

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

# Q-Learning

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

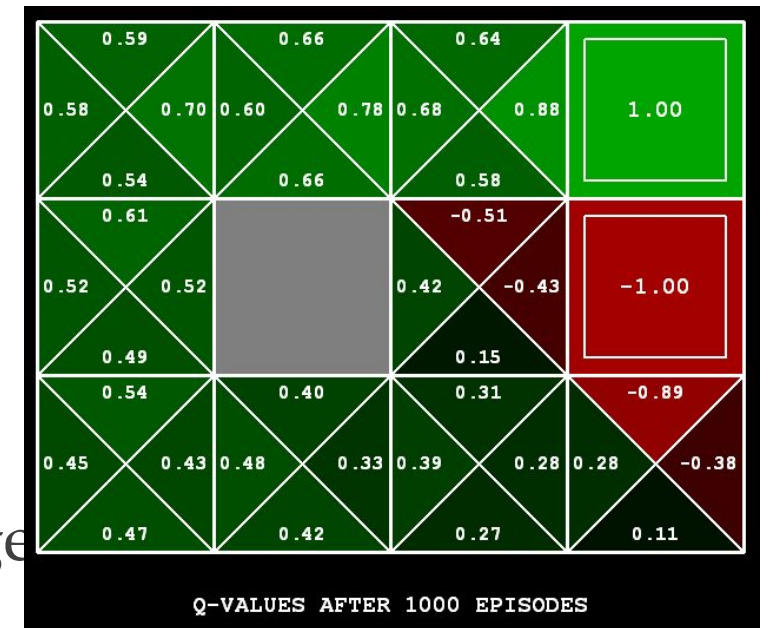$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

○ Learn Q(s,a) values as you go



Q-VALUES AFTER 1000 EPISODES

   ○ Receive a sample (s,a,s′,r)

   ○ Consider your old estimat $Q(s, a)$

   ○ 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

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

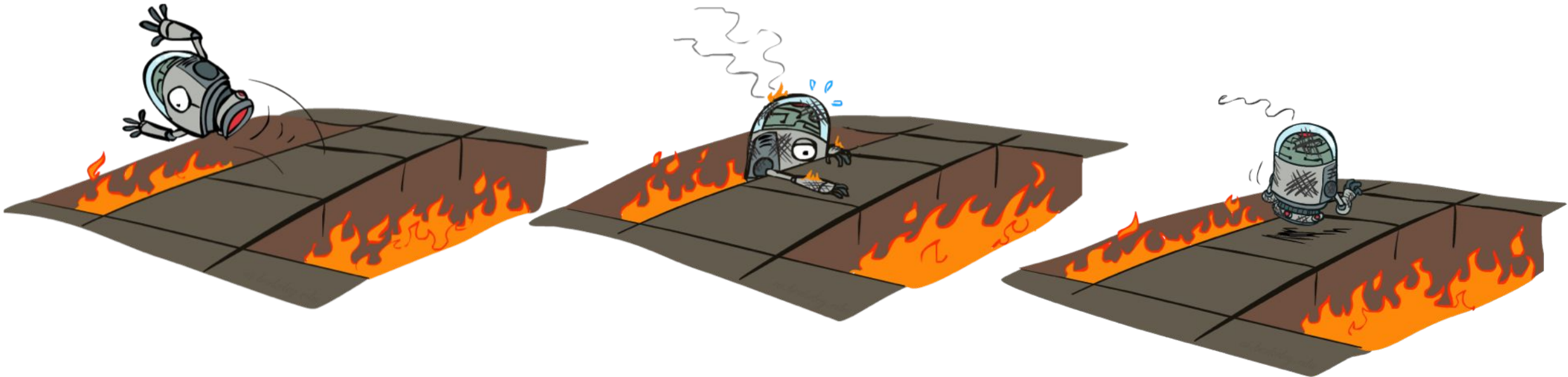# Video of Demo Q-Learning -- Gridworld

# Video of Demo Q-Learning -- Crawler
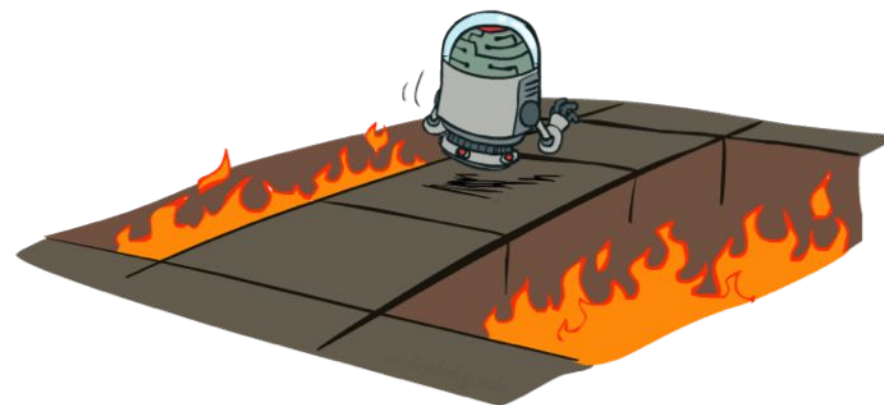
# Active Reinforcement Learning

o Passive reinforcement learning:

    o A passive learning agent has a fixed policy that determines its behavior

o Active reinforcement learning:

    o An active learning agent gets to decide what actions to take

# Q-Learning:
## act according to current optimal (and also explore...)

o Full reinforcement learning: optimal policies (like value iteration)
  - o You don't know the transitions T(s,a,s')
  - o You don't know the rewards R(s,a,s')
  - o You choose the actions now
  - o Goal: learn the optimal policy / values

o In this case:
  - o Learner makes choices!
  - o Fundamental tradeoff: exploration vs. exploitation
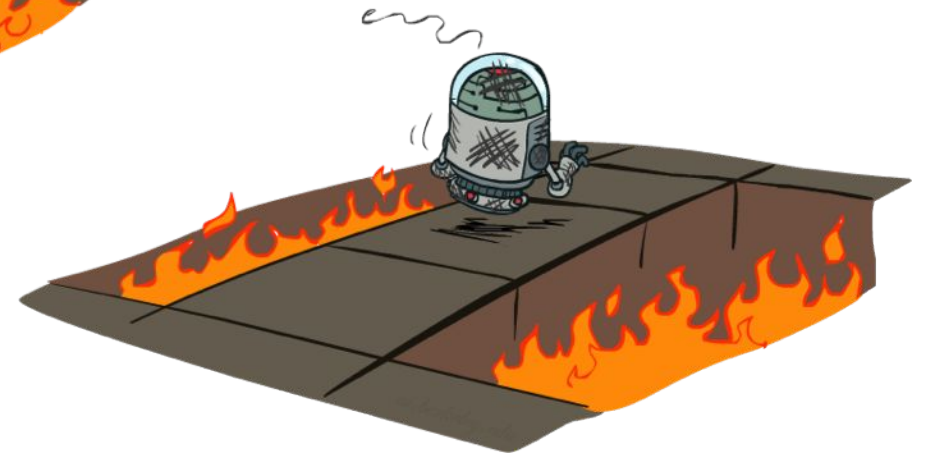  - o This is NOT offline planning! You actually take actions in the world and find out what happens...
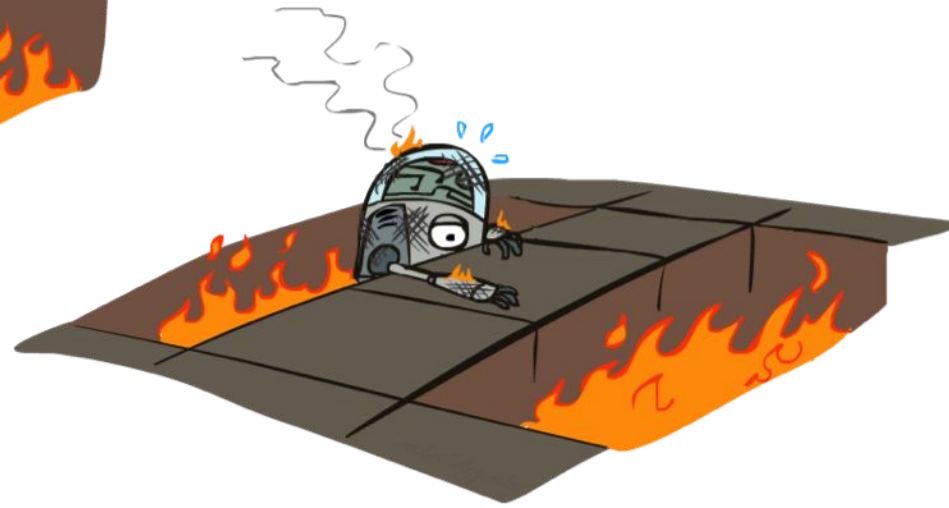
# Q-Learning Properties

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

o This is called off-policy learning

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
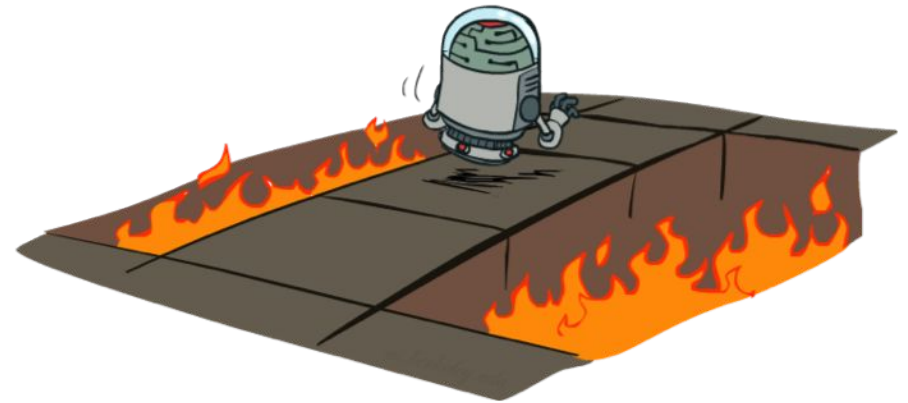  o Basically, in the limit, it doesn't matter how you select actions (!)
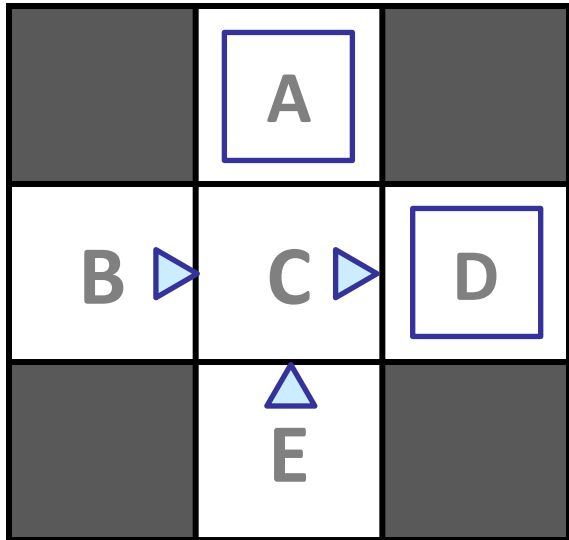
# Active Reinforcement Learning

# Model-Free Learning

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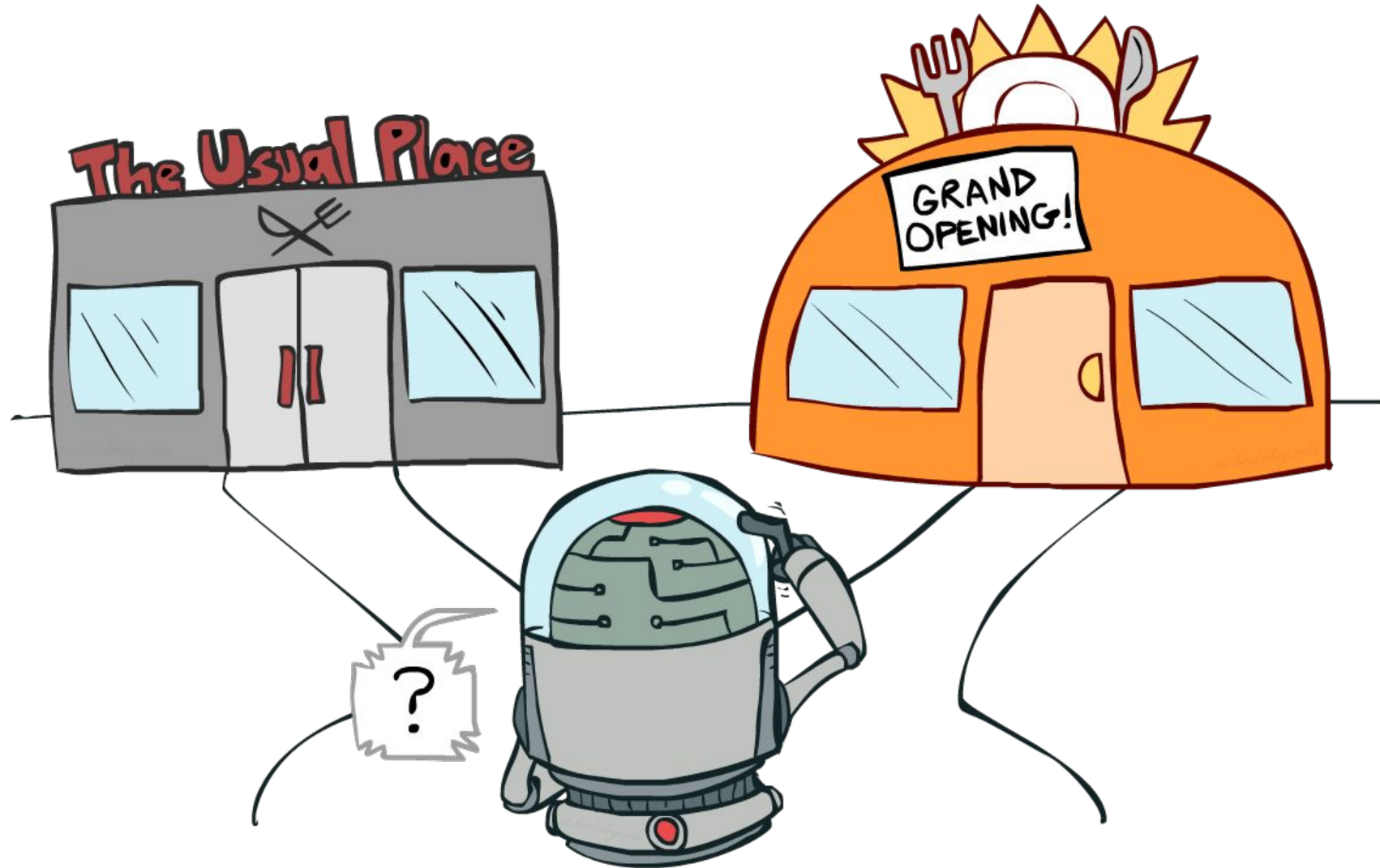○ but also explore…

# Model-Based Learning

act according to current optimal policy
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 ($\varepsilon$-greedy)
    o Every time step, flip a coin
    o With (small) probability $\varepsilon$, act randomly
    o With (large) probability $1-\varepsilon$, 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 $\varepsilon$ 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. $\quad f(u, n) = u + k/n$ is a predetermined constant

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

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

    N (s,a): number of times q-state (s,a) has been visited

o Note: this propagates the "bonus" back to states that lead to unknown states as well!

# Video of Demo Q-learning – Exploration Function – Crawler
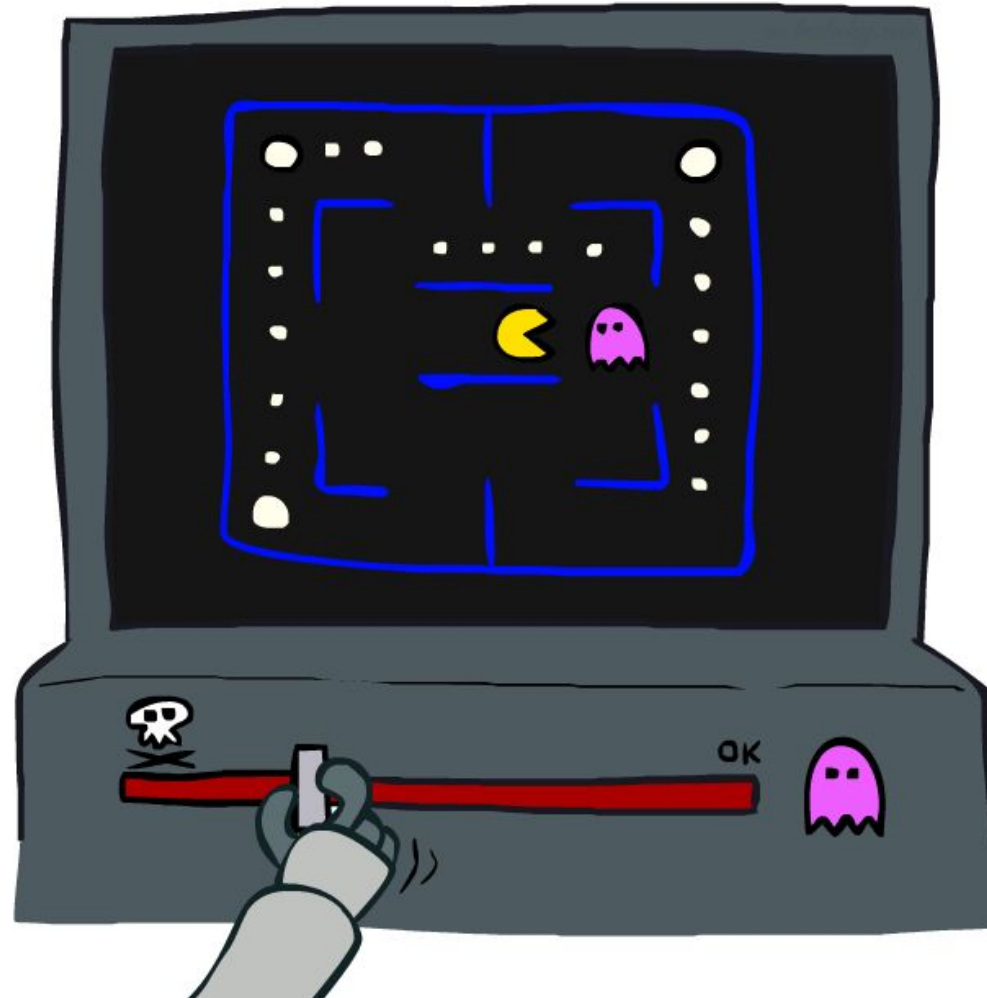
# Regret

- Even if you learn the optimal policy, you still make mistakes along the way!
- Regret is a measure of your total mistake cost: the difference between your (expected) rewards, including youthful suboptimality, and optimal (expected) rewards
- Minimizing regret goes beyond learning to be optimal – it requires optimally learning to be optimal
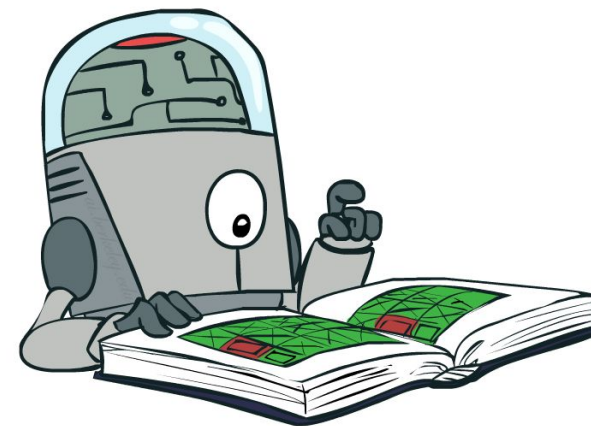- Example: random exploration and exploration functions both end up optimal, but random exploration has higher regret
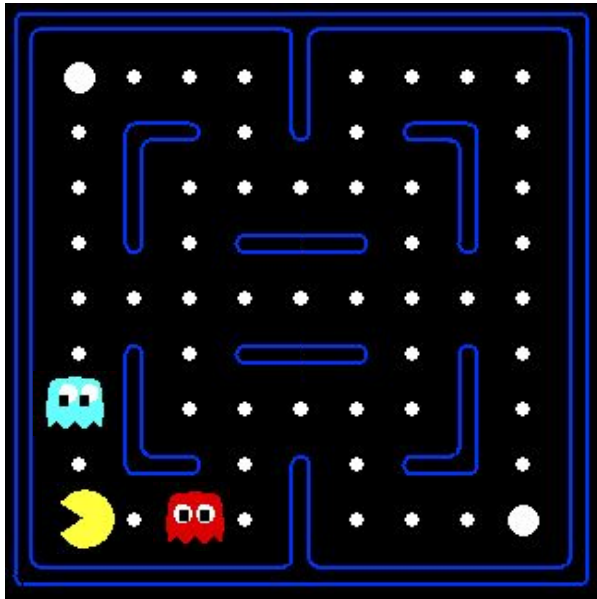
# Approximate Q-Learning

# 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
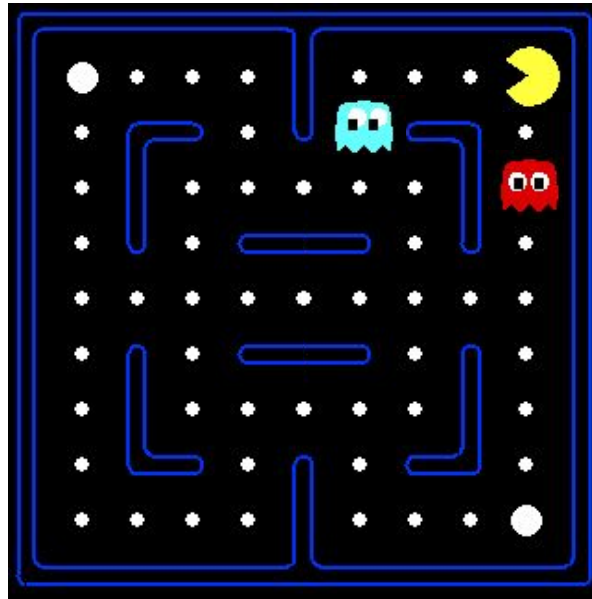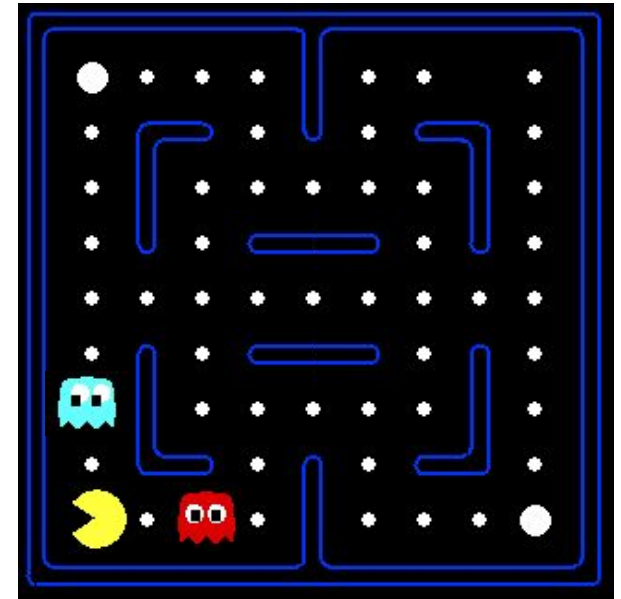
# Example: Pacman

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

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

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

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

○ 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

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

○ Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

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

$$Q(s, a) \leftarrow Q(s, a) + \alpha \, [\text{difference}] \qquad \text{Exact Q's}$$

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

○ 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: online least squares

# Example: Q-Pacman

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



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

$s$

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

$s'$

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

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

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

$r + \gamma \max_{a'} Q(s', a') = -500 + 0$
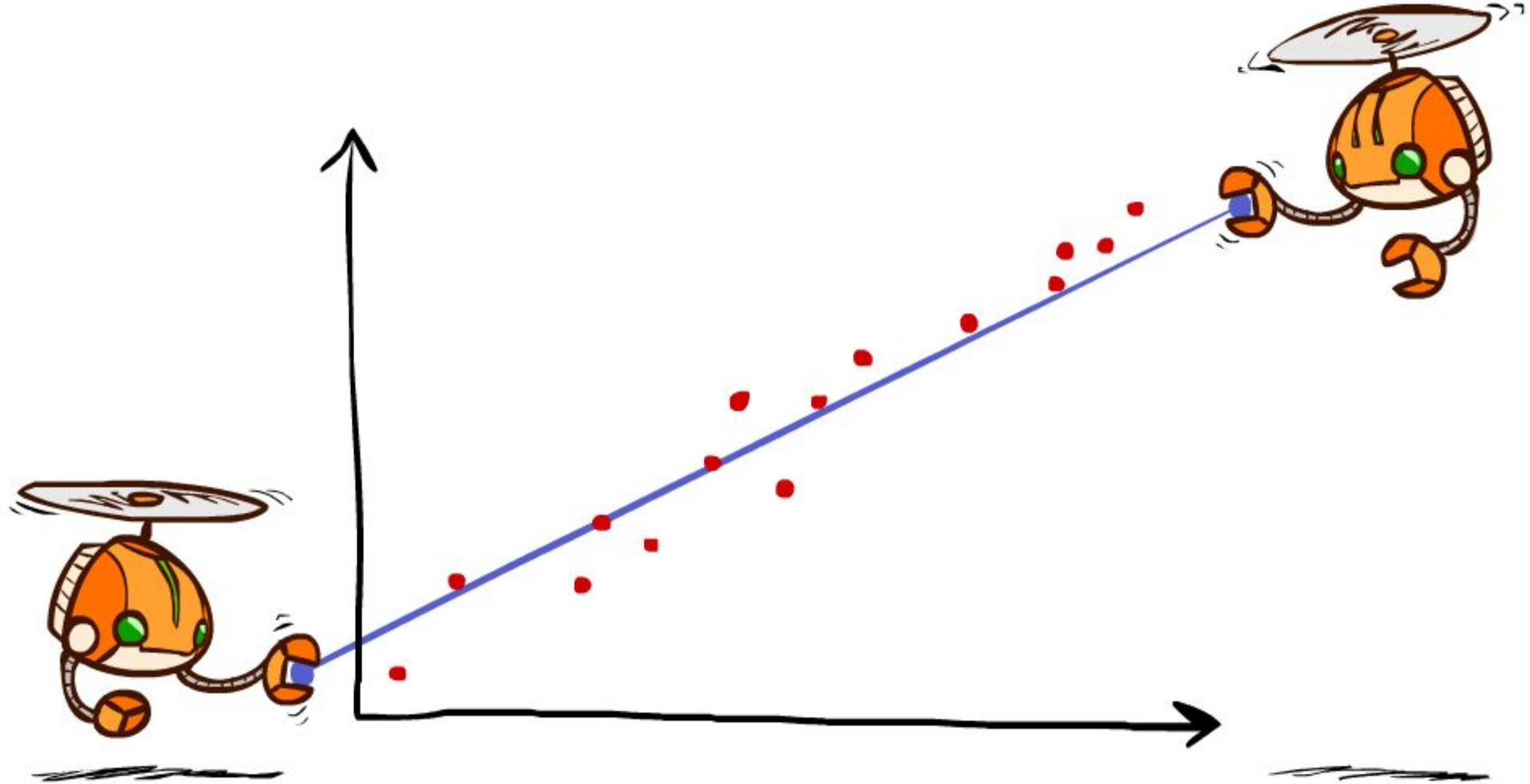
difference $= -501$

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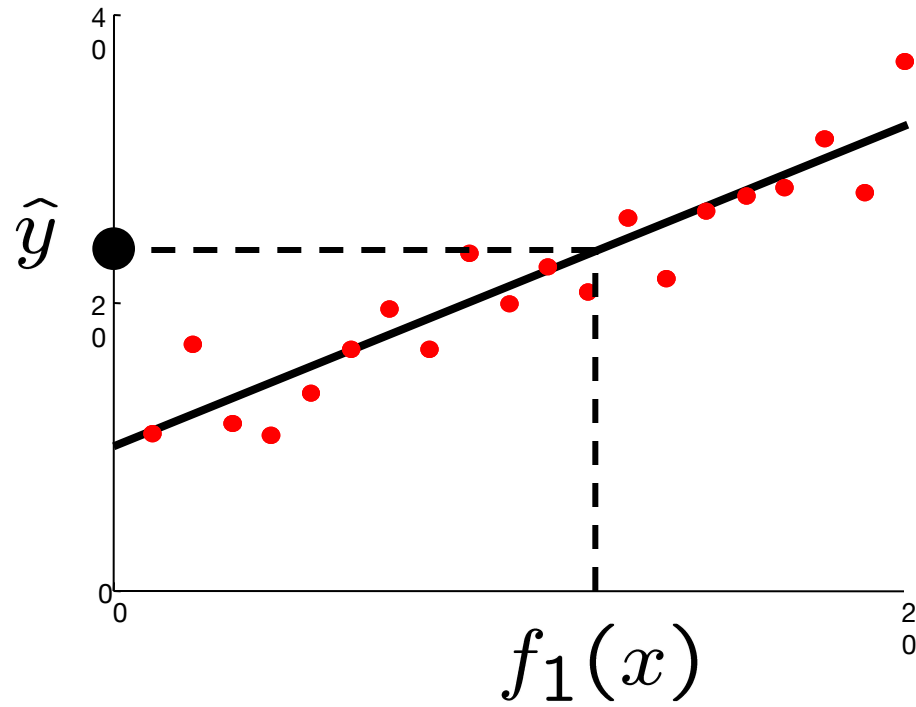
$$Q(s,a) = 3.0 f_{DOT}(s,a) - 3.0 f_{GST}(s,a)$$

# Video of Demo Approximate Q-Learning -- Pacman

# Q-Learning and Least Squares

# Linear Approximation: Regression



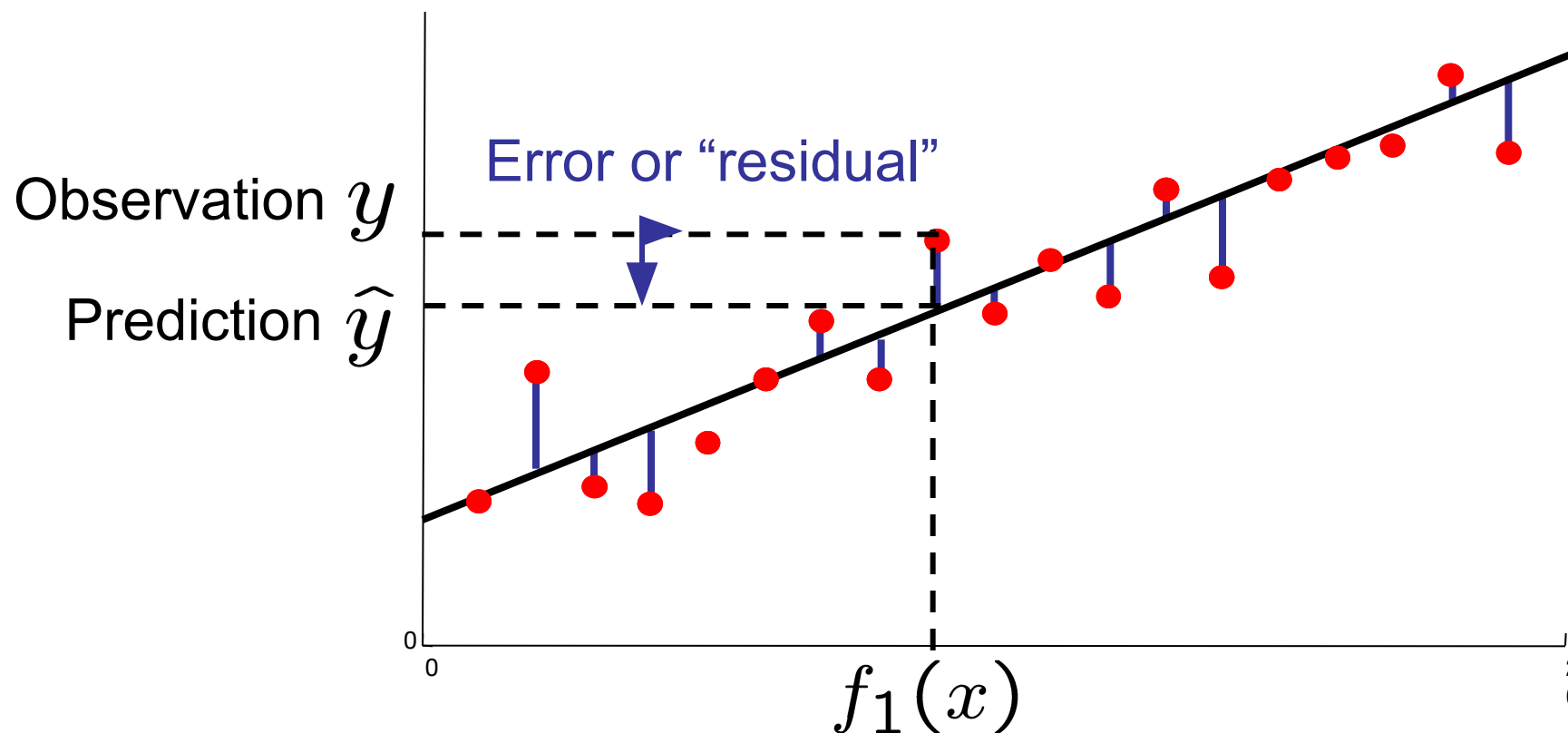Prediction:

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

Prediction:

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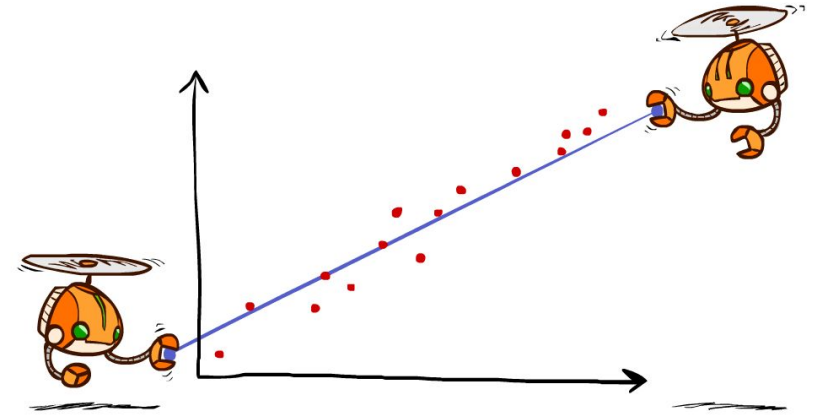
# Optimization: Least Squares

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

# 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"

# Summary: 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

*use features to generalize*

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

## Unknown MDP: Model-Free

*use features to generalize*

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

# RL and dopamine



o Dopamine signal generated by parts of the striatum

o Encodes predictive error in value function (as in TD learning)
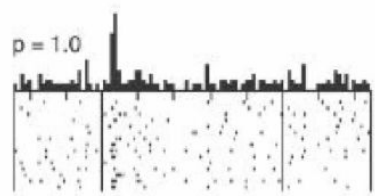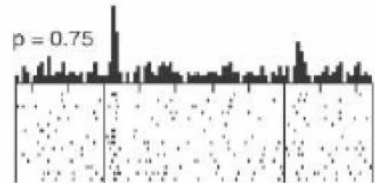
reward following 0% predictive cue

reward following 50% predictive cue

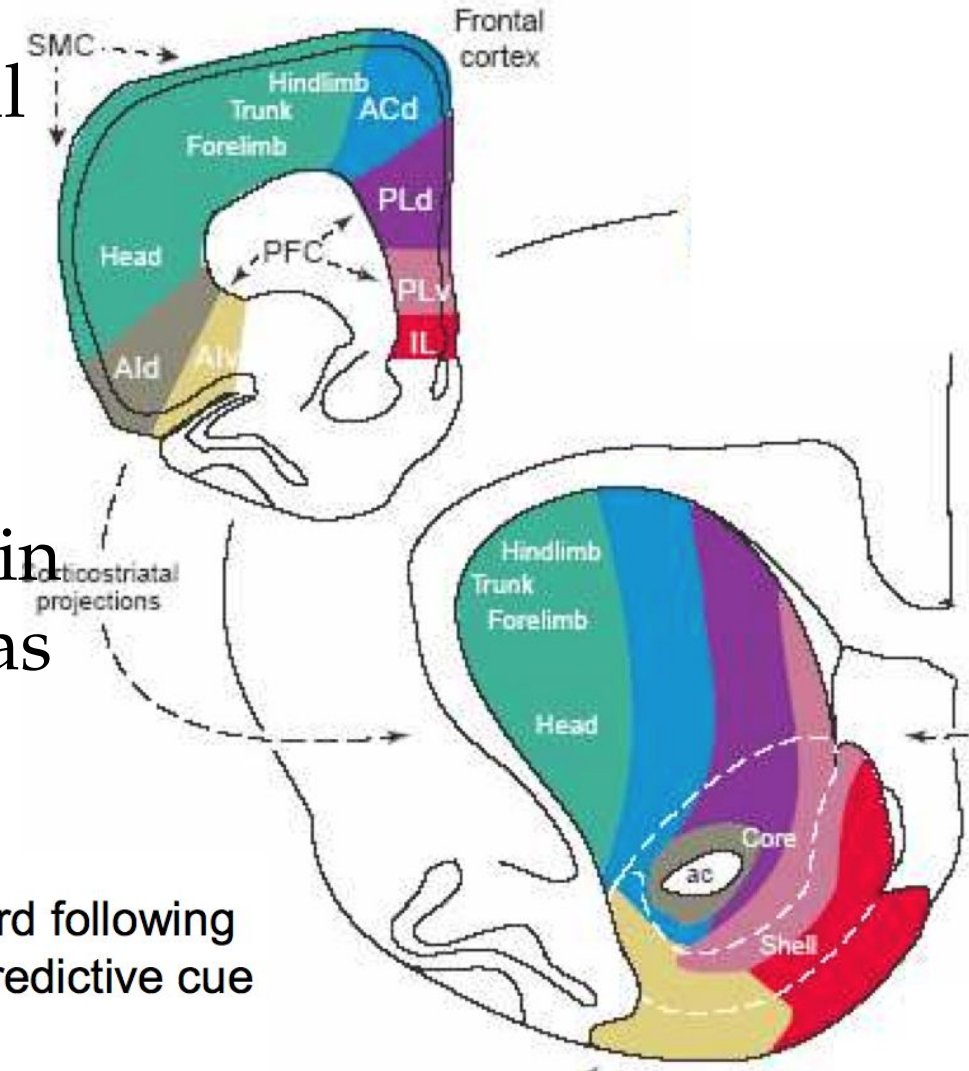reward following 100% predictive cue

no reward following 100% predictive cue

(Fiorillo et al 2003)

Voorn et al 2004

# Next Section: Advanced Topics

o Advanced topic I: Adversarial machine learning

o Advanced topic II: Fairness in machine learning

o Advanced topic III: CLIP

o Final lecture: AI safety (Stuart)