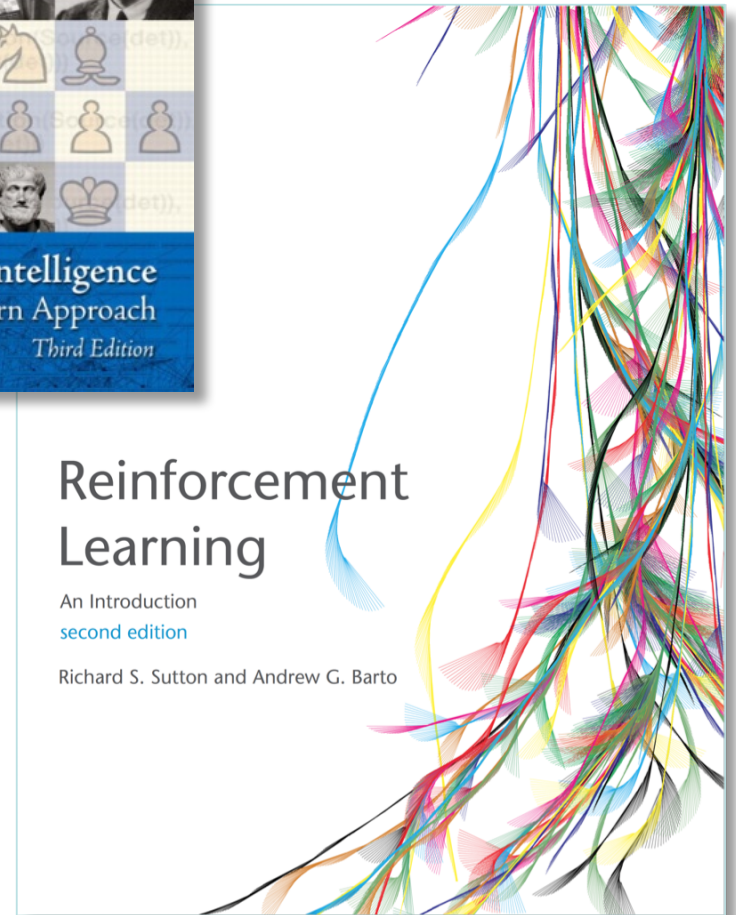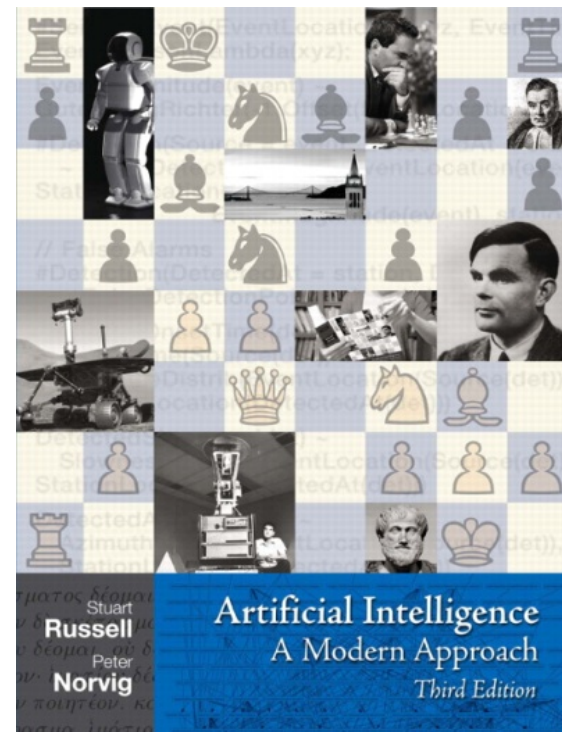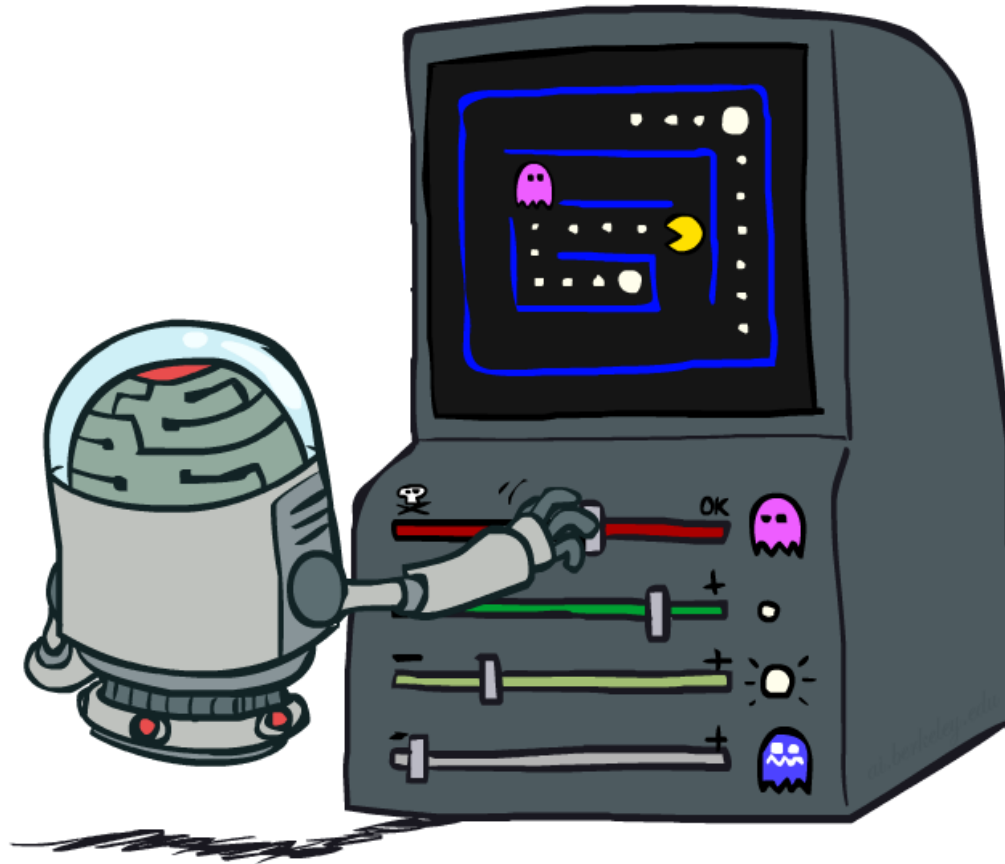# Reinforcement Learning 2

## Sutton and Barto, Chapter 6.1, 6.2, 6.5
## AIMA Chapter 21

Slides courtesy of Dan Klein and Pieter Abbeel
University of California, Berkeley

# Reinforcement Learning 2



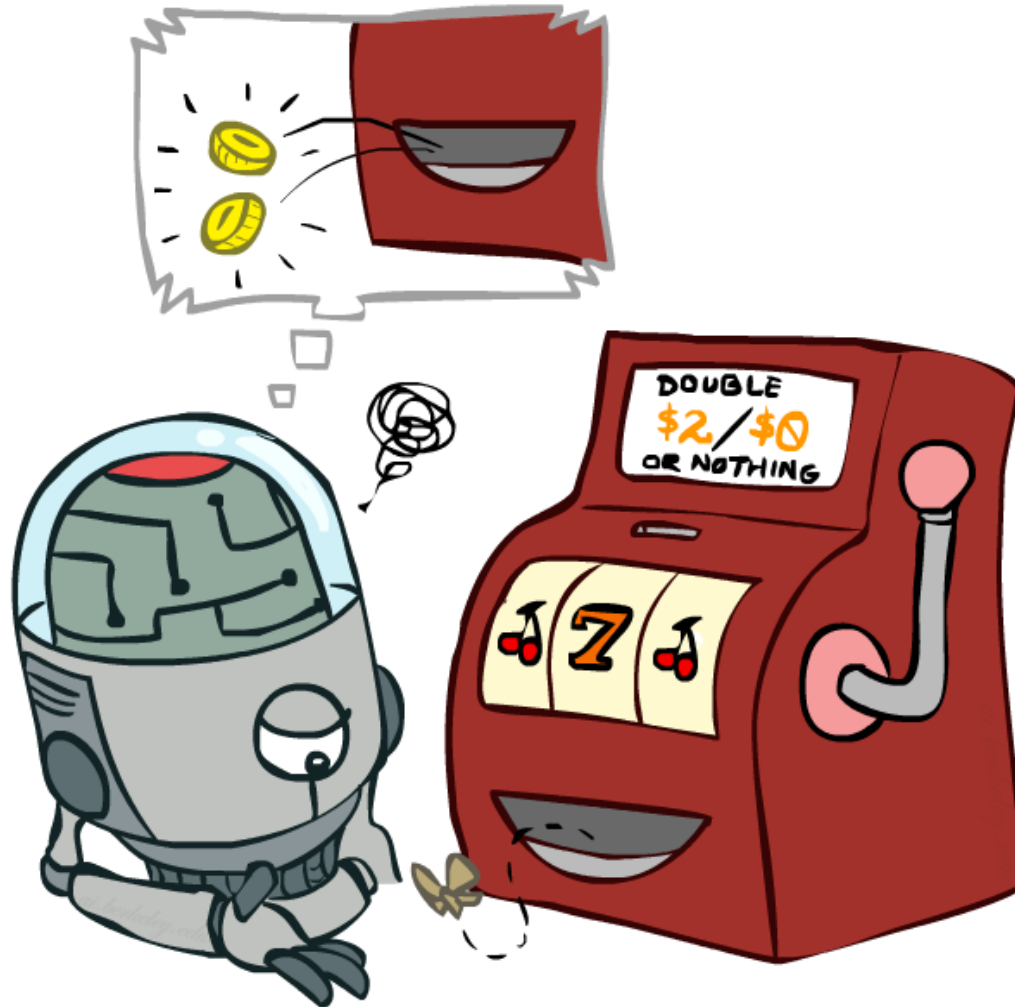Slides courtesy of Dan Klein and Pieter Abbeel – University of California, Berkeley
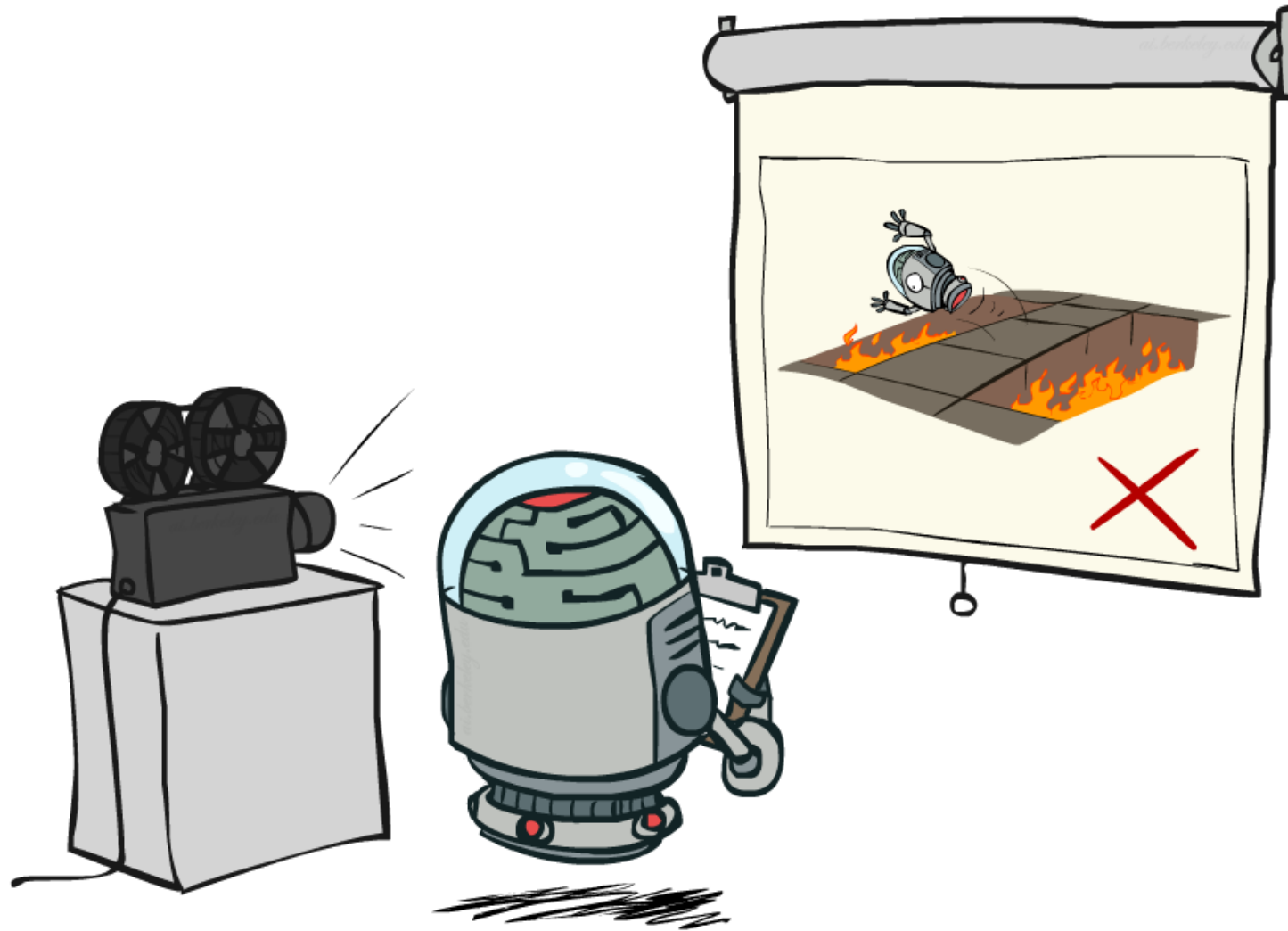
# Review: Reinforcement Learning

- We still assume an 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, so must try out actions

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

# Review: Passive Reinforcement Learning

- **Simplified task: policy evaluation**
  - Input: a fixed policy $\pi(s)$
  - You don't know the transitions $T(s,a,s')$
  - You don't know the rewards $R(s,a,s')$
  - Goal: learn the state values

- **In this case:**
  - Learner is "along for the ride"
  - No choice about what actions to take
  - Just execute the policy and learn from experience
  - This is NOT offline planning!  You actually take actions in the world.

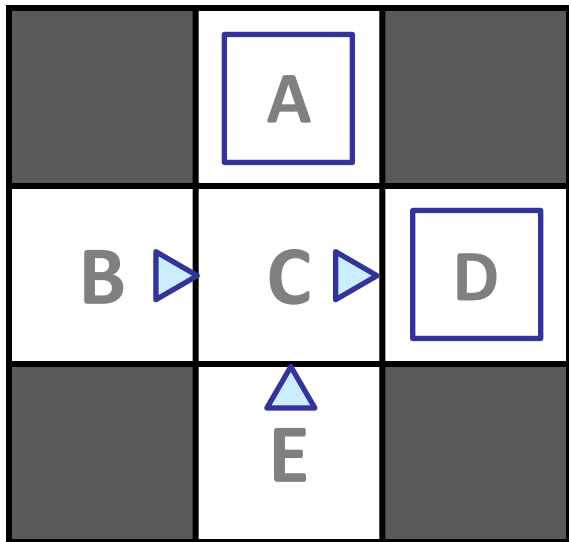# Review: Direct Evaluation

- Goal: Compute values for each state under $\pi$

- Idea: Average together observed sample values
  - Act according to $\pi$
  - Every time you visit a state, write down what the sum of discounted rewards turned out to be
  - Average those samples

- This is called direct evaluation

# Review: Example: Direct Evaluation

## Input Policy π



Assume: γ = 1

## Observed Episodes (Training)

### Episode 1

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 2

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 3

E, north, C, -1
C, east,    D, -1
D, exit,     x, +10

### Episode 4

E, north, C, -1
C, east,    A, -1
A, exit,     x, -10

## Output Values

# Review: Problems with Direct Evaluation

- ■ **What's good about direct evaluation?**
    - ■ It's easy to understand
    - ■ It doesn't require any knowledge of T, R
    - ■ It eventually computes the correct average values, using just sample transitions

- ■ **What bad about it?**
    - ■ It wastes information about state connections
    - ■ Each state must be learned separately
    - ■ So, it takes a long time to learn

**Output Values**



*If B and E both go to C under this policy, how can their values be different?*
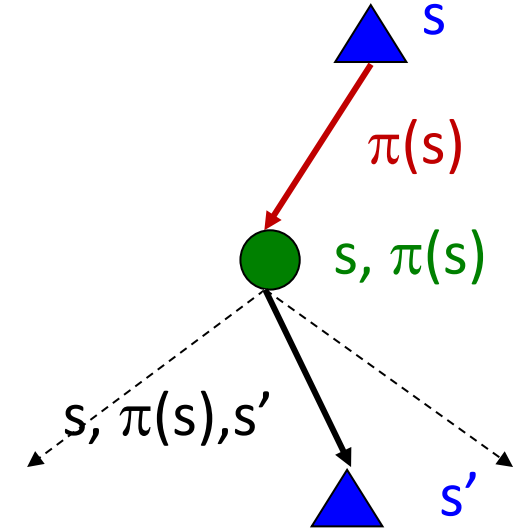
# Review: Why Not Use Policy Evaluation Instead?

- ■ Simplified Bellman updates calculate V for a fixed policy:
  - ■ Each round, replace V with a one-step-look-ahead layer over V

$$V_0^\pi(s) = 0$$

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

s

$\pi(s)$

s, $\pi(s)$

s, $\pi(s)$,s'

s'

  - ■ This approach fully exploited the connections between the states
  - ■ Unfortunately, we need T and R to do it!

- ■ Key question: how can we do this update to V without knowing T and R?
  - ■ In other words, how to we take a weighted average without knowing the weights?
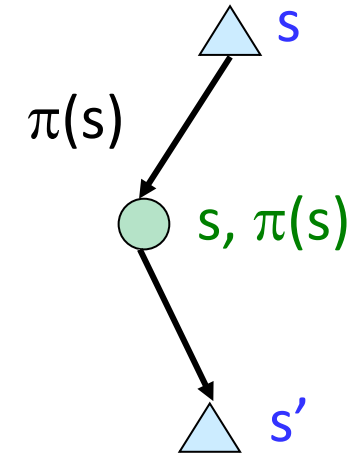
# Review: Temporal Difference Learning

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

- **Temporal difference learning of values**
  - Policy still fixed, still doing evaluation!
  - 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):     $V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$

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

# Review: Exponential Moving Average

- **Exponential moving average**
  - The running interpolation update: $\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$

  - Makes recent samples more important:

  $$\bar{x}_n = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \ldots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \ldots}$$

  - Forgets about the past (distant past values were wrong anyway)

- **Decreasing learning rate (alpha) can give converging averages**

# Example: Temporal Difference Learning

## States



A

B C D

E

*Assume: $\gamma = 1$, $\alpha = 1/2$*

## Observed Transitions

B, east, C, -2          C, east, D, -2



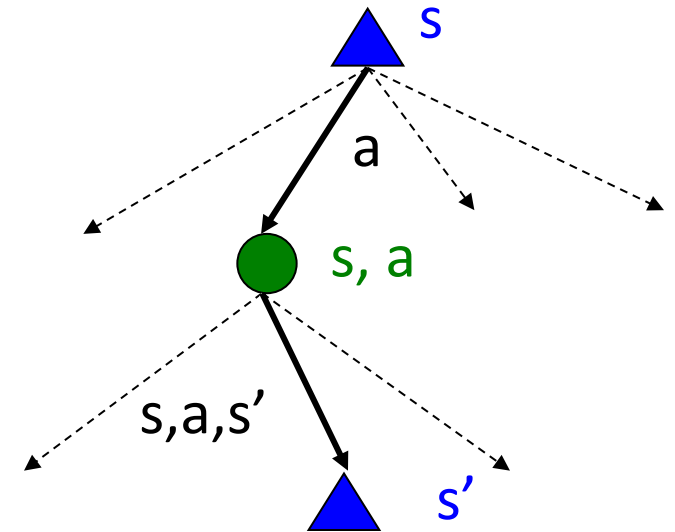$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha\left[R(s, \pi(s), s') + \gamma V^\pi(s')\right]$$

# Review: Problems with TD Value Learning

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

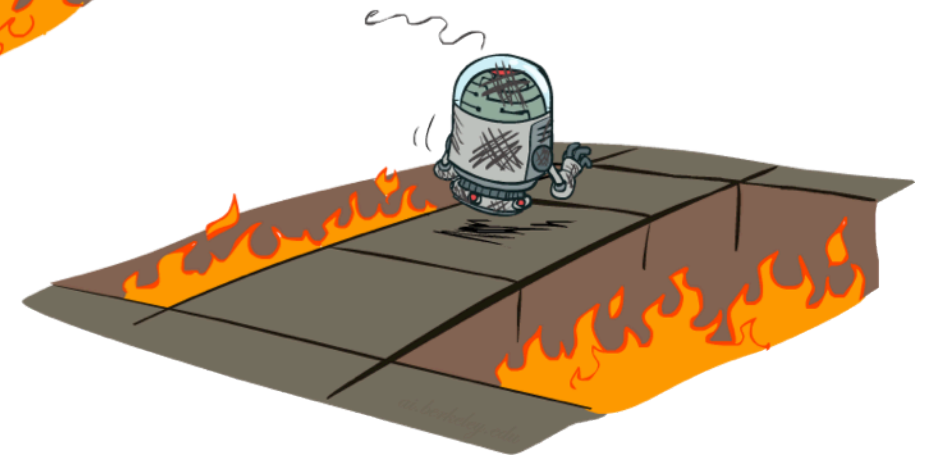- However, if we want to turn values into a (new) policy, we're sunk:
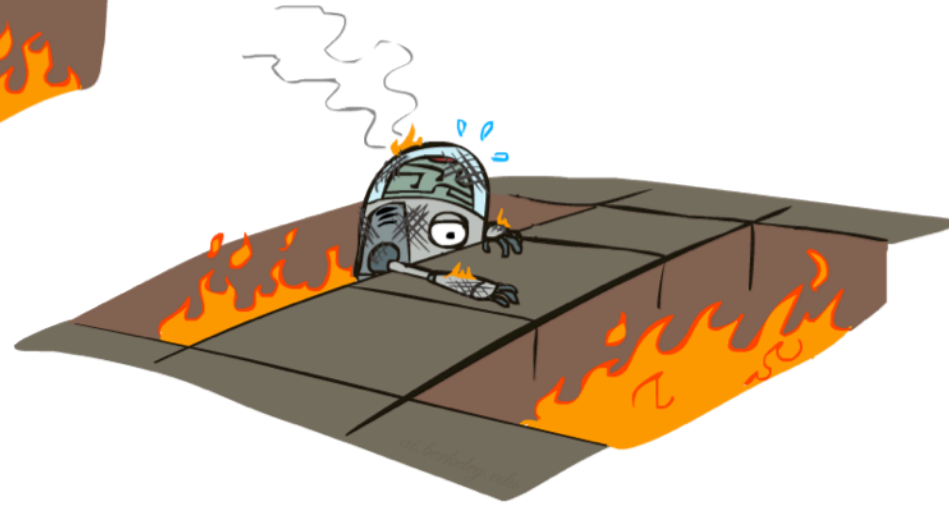
$$\pi(s) = \arg\max_a Q(s, a)$$

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

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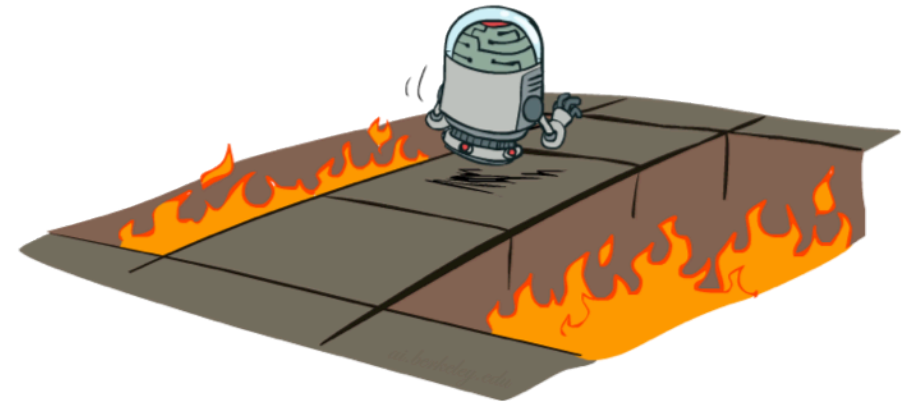- Makes action selection model-free too!

# Active Reinforcement Learning

# Active Reinforcement Learning

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

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

# Q-Learning

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

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$
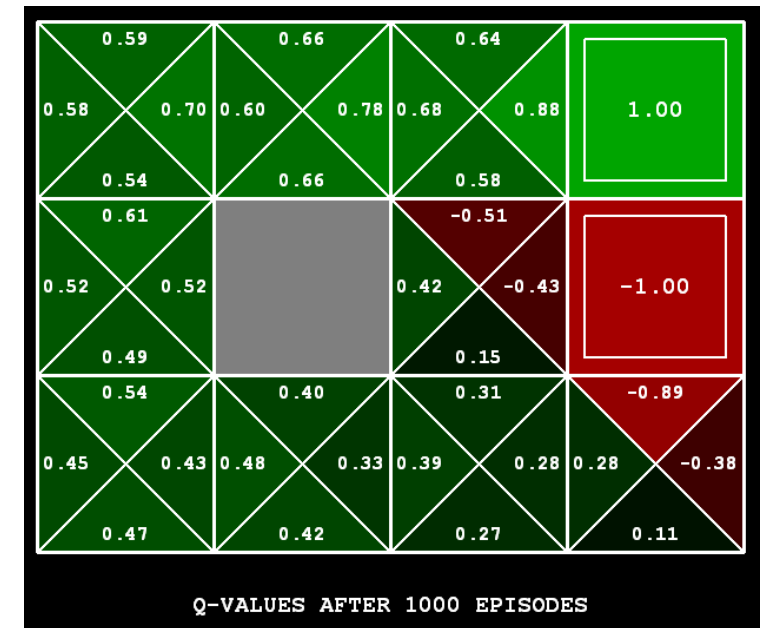
- **Learn Q(s,a) values as you go**
  - Receive a sample (s,a,s',r)
  - Consider your old estimate: $Q(s,a)$
  - Consider your new sample estimate:
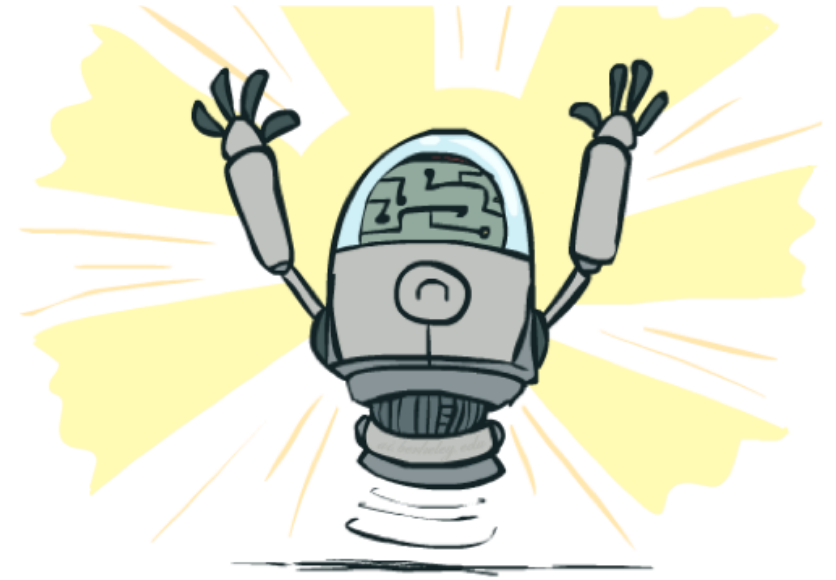
  $$sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$

  - Incorporate the new estimate into a running average:

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



Q-VALUES AFTER 1000 EPISODES

# Q-Learning Properties

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

- 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 (!)

# Review: Reinforcement Learning

- We still assume an 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, so must try out actions

- 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

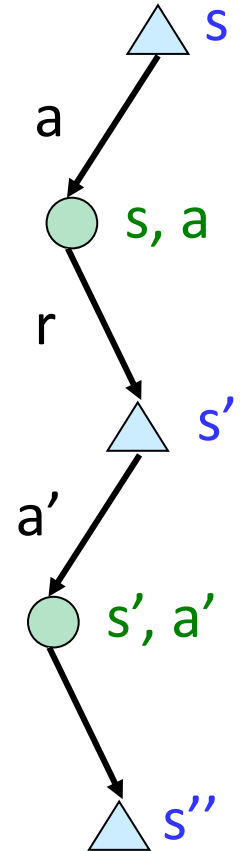| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | Q-learning |
| Evaluate a fixed policy $\pi$ | Value Learning |

# Model-Free Learning

- ## Model-free (temporal difference) learning
  - Experience world through episodes

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

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

  - Over time, updates will mimic Bellman updates

# Q-Learning

- We'd like to do Q-value updates to each Q-state:

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

  - But can't compute this update without knowing T, R

- Instead, compute average as we go

  - Receive a sample transition (s,a,r,s')
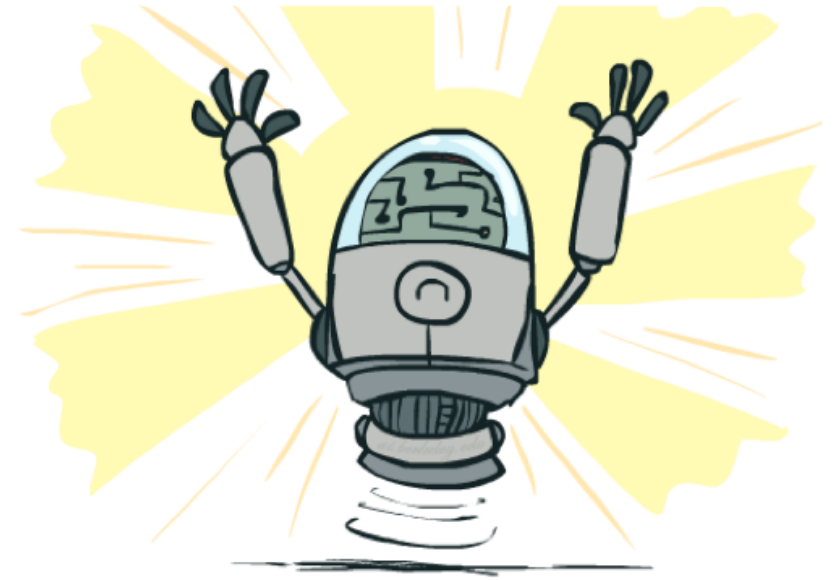
  - This sample suggests

    $$Q(s, a) \approx r + \gamma \max_{a'} Q(s', a')$$

  - But we want to average over results from (s,a)  (Why?)
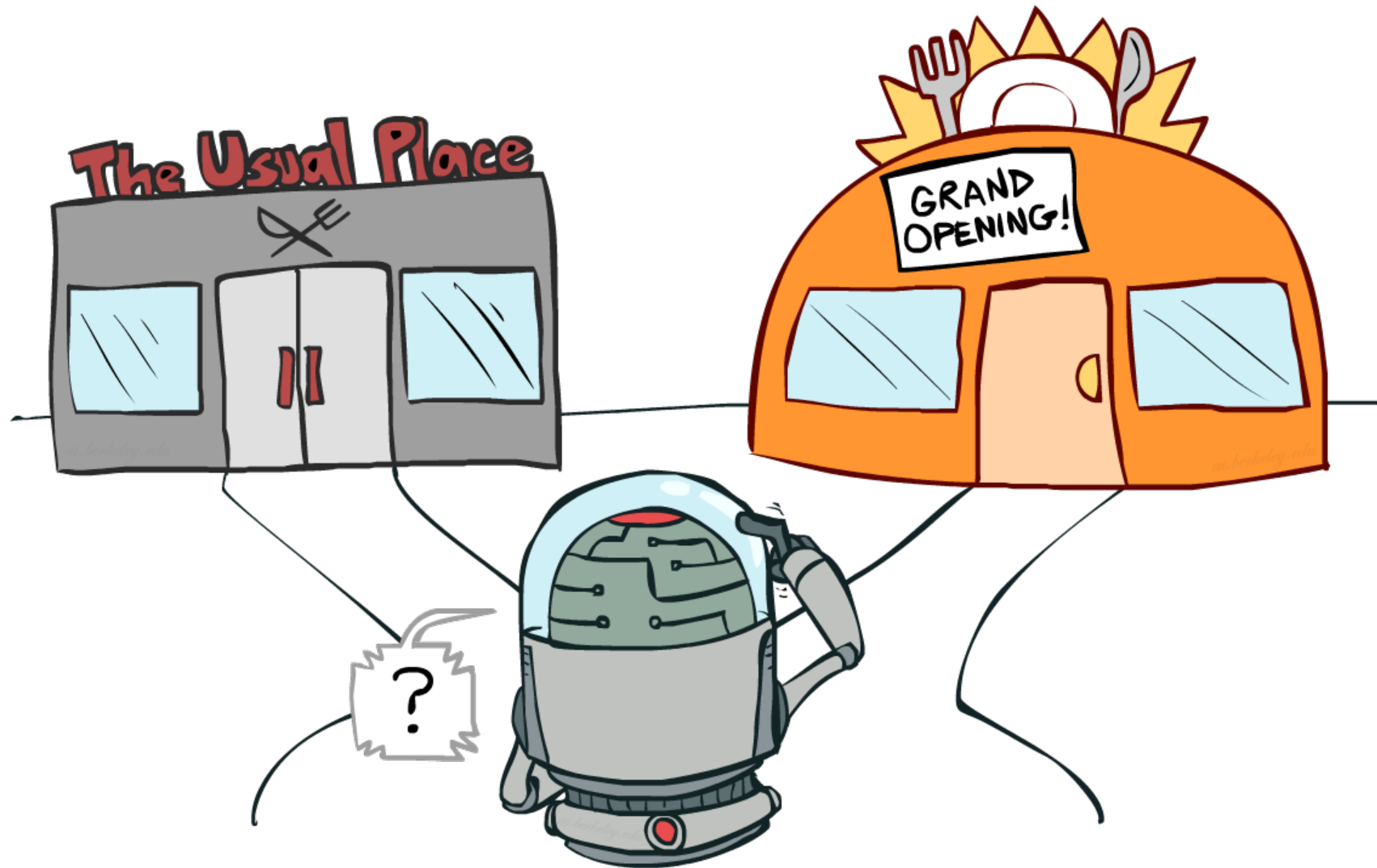
  - So keep a running average

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

# Q-Learning Properties

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

- 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 (!)

[Demo: Q-learning – auto – cliff grid (L11D1)]

# Exploration vs. Exploitation

# How to Explore?

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

# How to Explore?

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

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

# Exploration Functions

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

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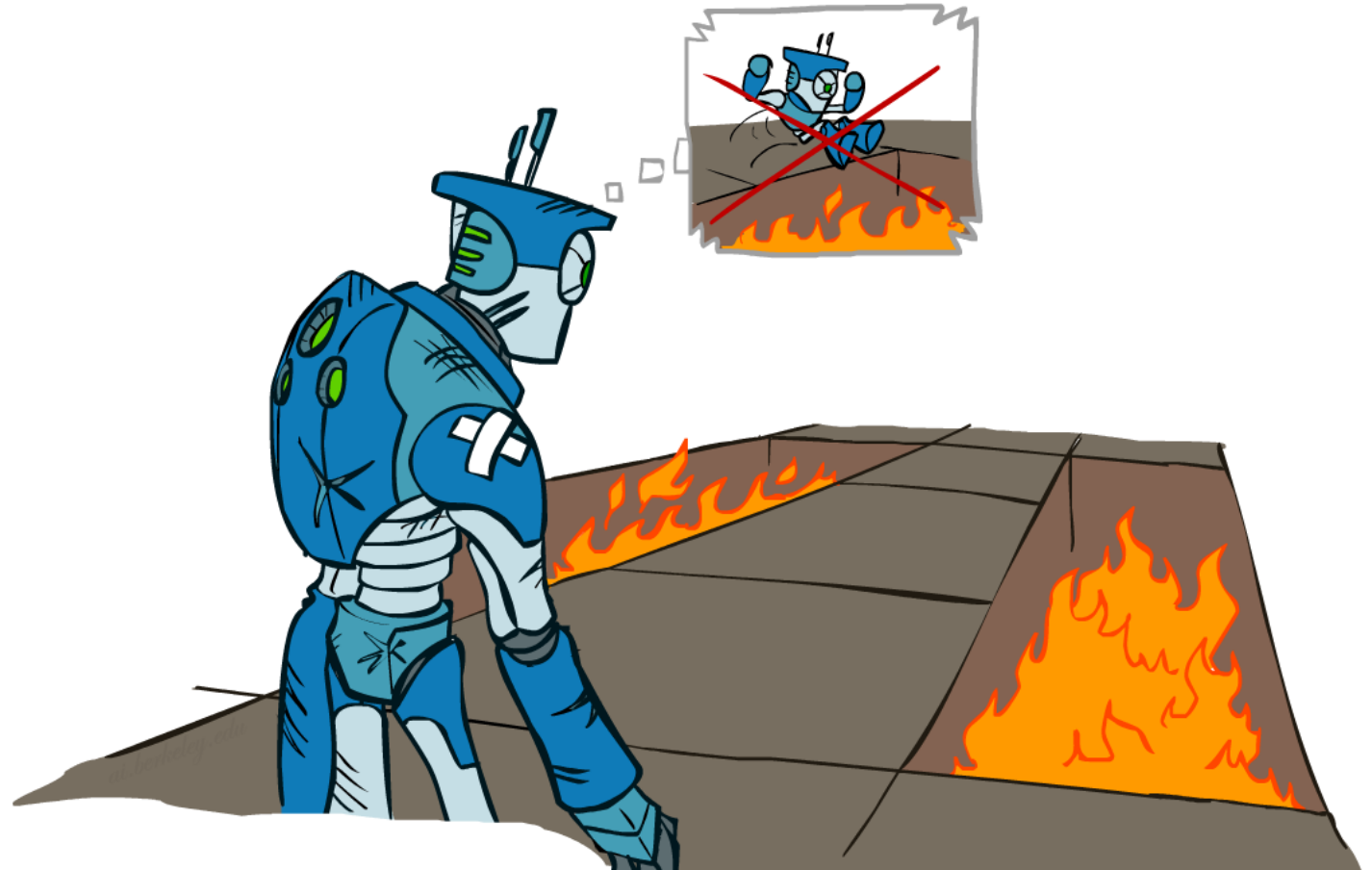
    Modified Q-Update: $Q(s, a) \leftarrow_\alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$

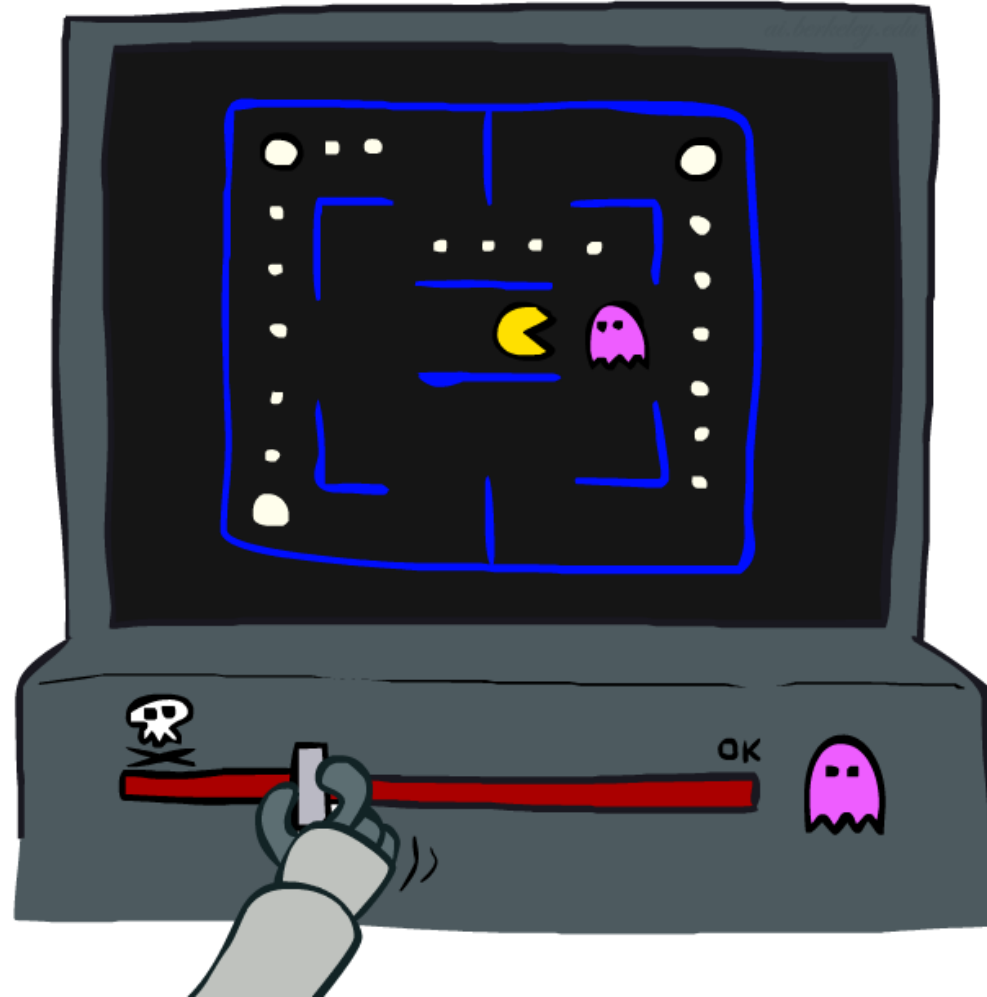  - Note: this propagates the "bonus" back to states that lead to unknown states as well!

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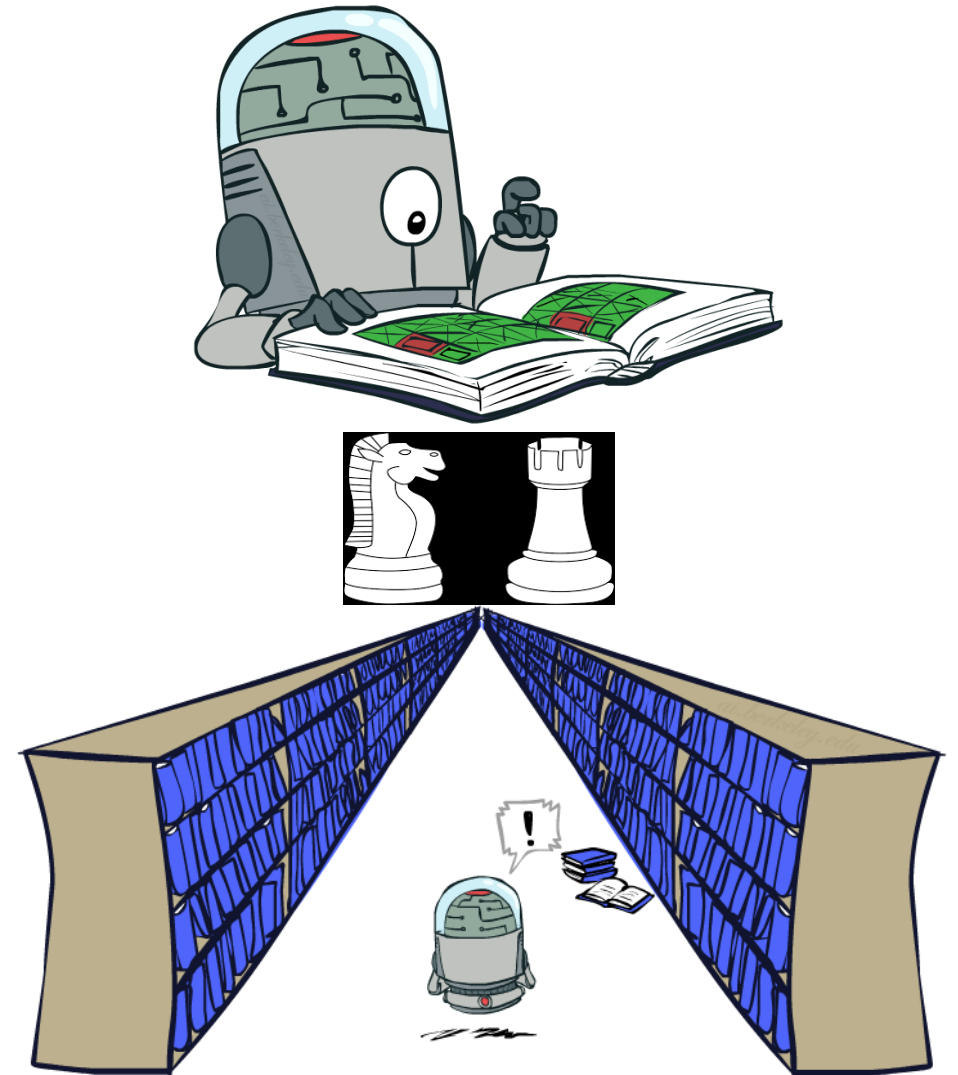- Example: random exploration and exploration functions both end up optimal, but random exploration has higher regret

# Approximate Q-Learning

# Generalizing Across States

- 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
  - This is a fundamental idea in machine learning, and we'll see it over and over again

# 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:

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

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

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

$$w_i \leftarrow w_i + \alpha \, [\text{difference}] \, f_i(s, a) \qquad \text{Approximate Q's}$$
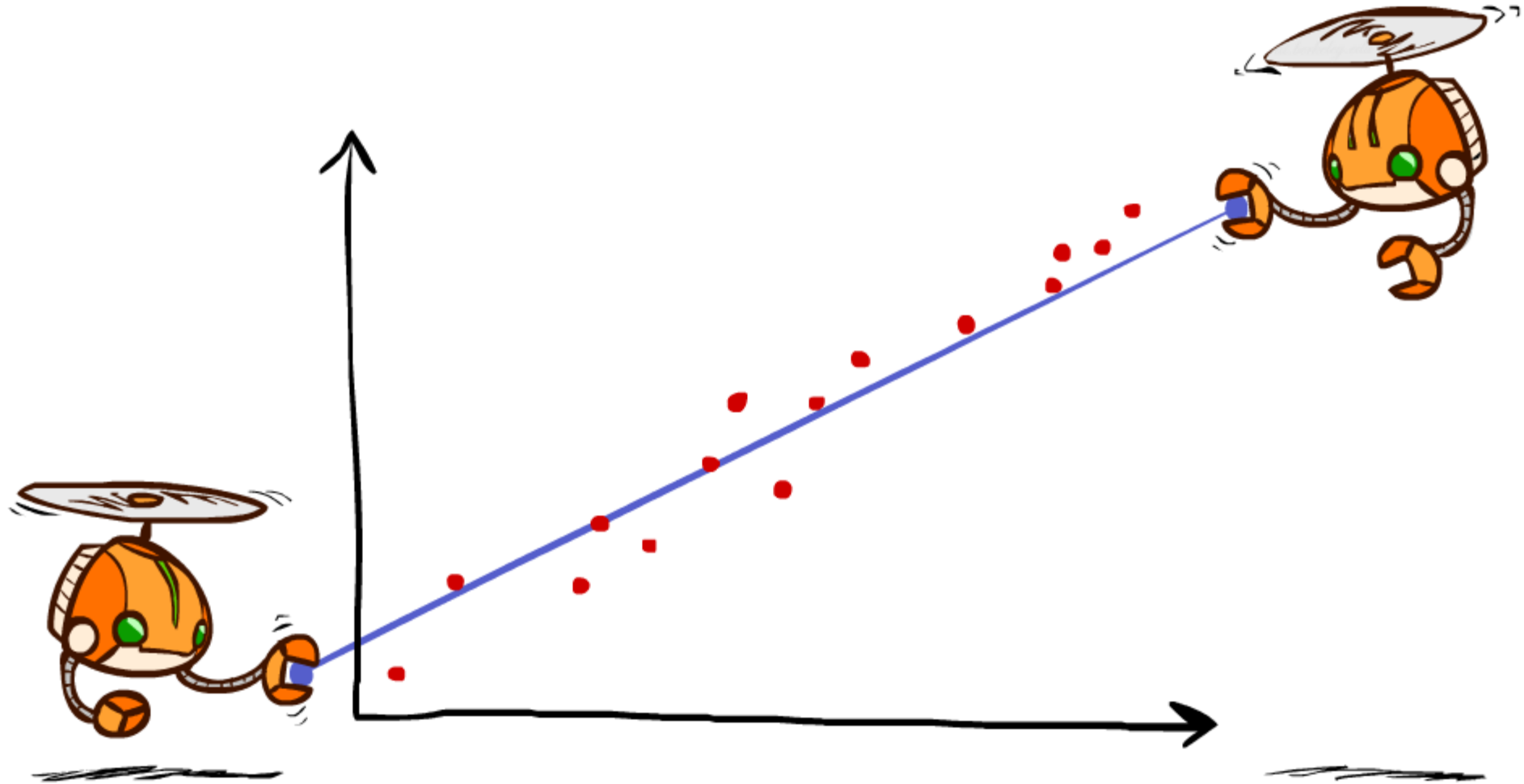
- Intuitive interpretation:
  - Adjust weights of active features
  - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features
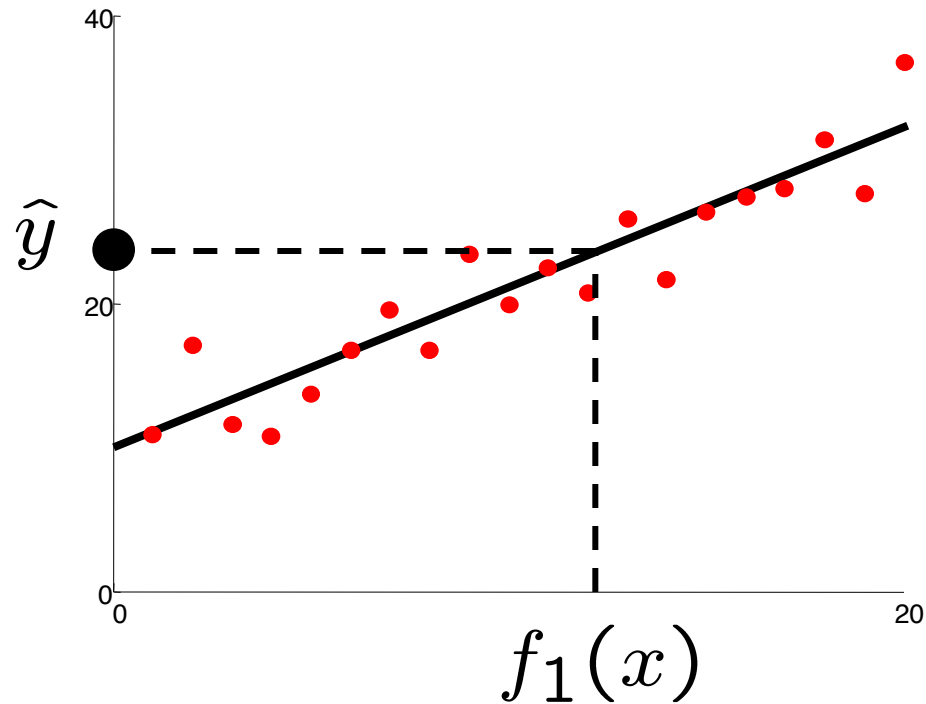
- Formal justification: online least squares
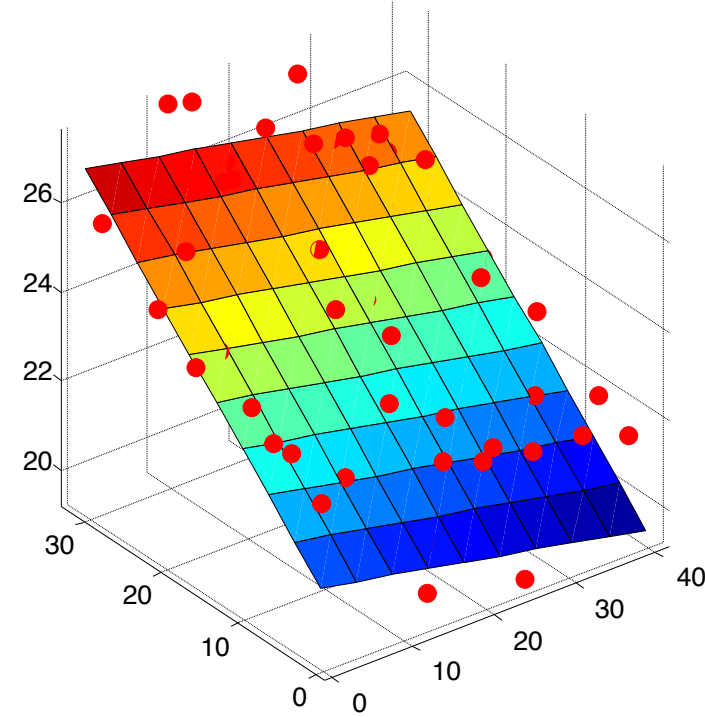
# Q-Learning and Least Squares

# Linear Approximation: Regression*
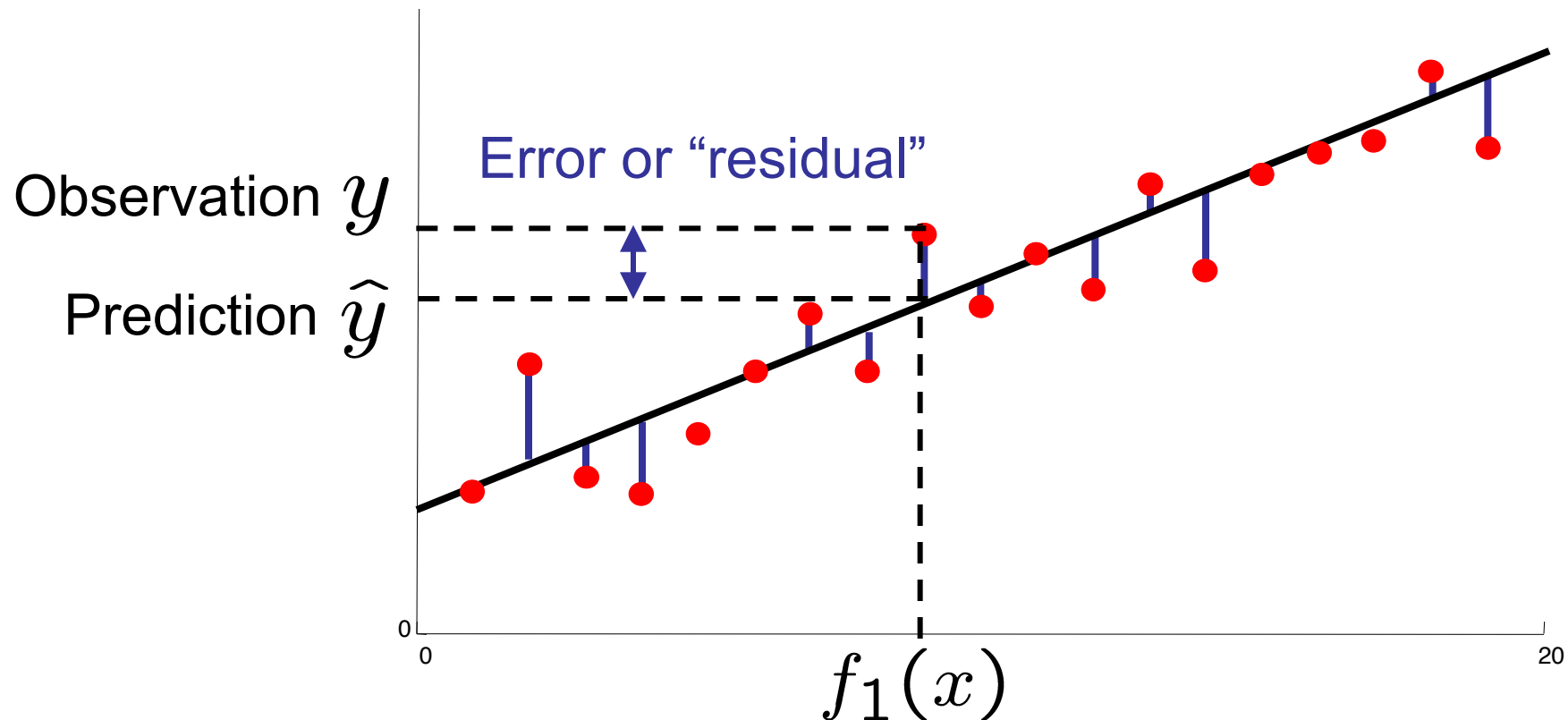


Prediction:
$$\hat{y} = w_0 + w_1 f_1(x)$$

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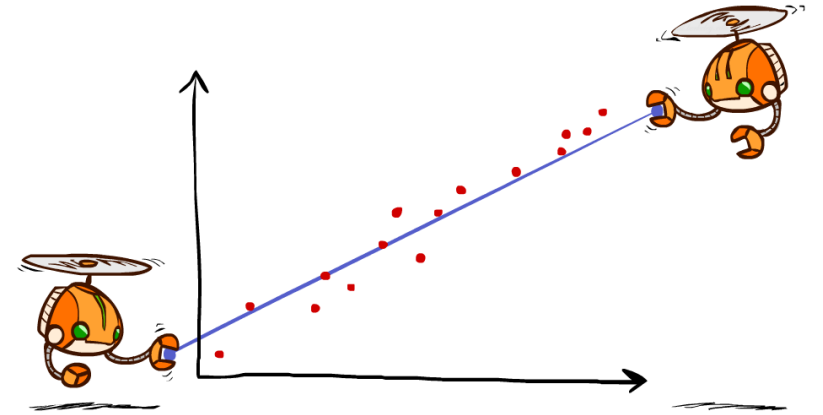
# Optimization: Least Squares*

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

# Minimizing Error*

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

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

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

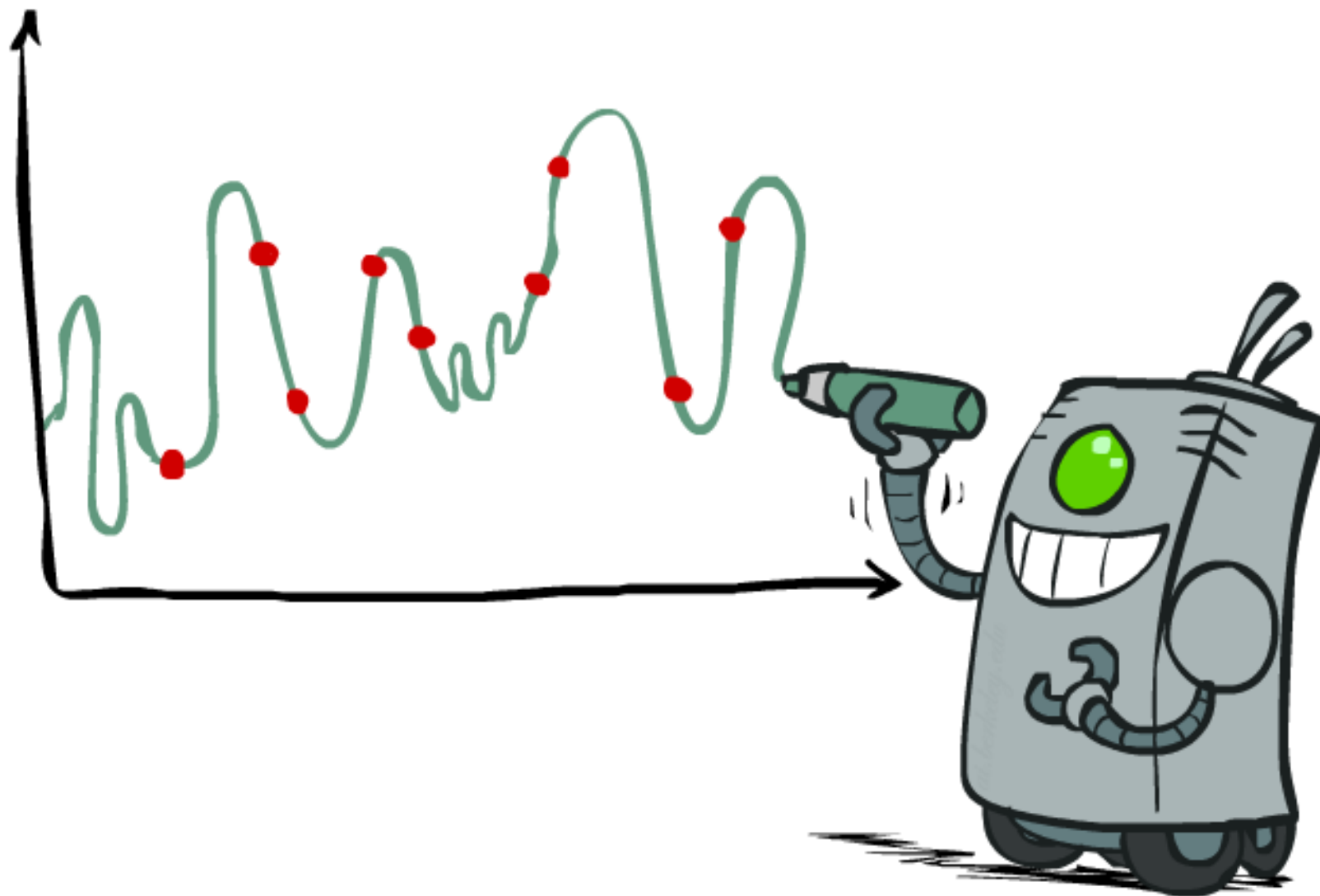$$w_m \leftarrow w_m + \alpha\left(y - \sum_k w_k f_k(x)\right) f_m(x)$$

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"

# Policy Search

# Policy Search

- Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate V / Q best
  - E.g. your value functions from project 2 were probably horrible estimates of future rewards, but they still produced good decisions
  - Q-learning's priority: get Q-values close (modeling)
  - Action selection priority: get ordering of Q-values right (prediction)
  - We'll see this distinction between modeling and prediction again later in the course

- Solution: learn policies that maximize rewards, not the values that predict them

- Policy search: start with an ok solution (e.g. Q-learning) then fine-tune by hill climbing on feature weights

# Policy Search

- Simplest policy search:
  - Start with an initial linear value function or Q-function
  - Nudge each feature weight up and down and see if your policy is better than before

- Problems:
  - How do we tell the policy got better?
  - Need to run many sample episodes!
  - If there are a lot of features, this can be impractical

- Better methods exploit look ahead structure, sample wisely, change multiple parameters…

# Term Project

- Pick a topic from the class, and write a program to demonstrate how it can apply to your robots.  Short write-up or video to describe what you did.

- Teams of up to 5 people.

- Examples:
  - 8 Puzzle to sort grid of robots into ascending order by ID
  - Adversarial search for a game between R2D2 and R2Q5
  - Natural language processing for robots (translate commands into code)

- Alternate: do one of last year's homework assignments (on your own, not with a team).

# In class exercise

- Robot Excercise 2: Robot Navigation with A*
- http://artificial-intelligence-class.org