

---

# CS 471/571 (Fall 2023): Introduction to Artificial Intelligence

## Lecture 15: Probability

---

Thanh H. Nguyen

Source: <http://ai.berkeley.edu/home.html>



# Reminder

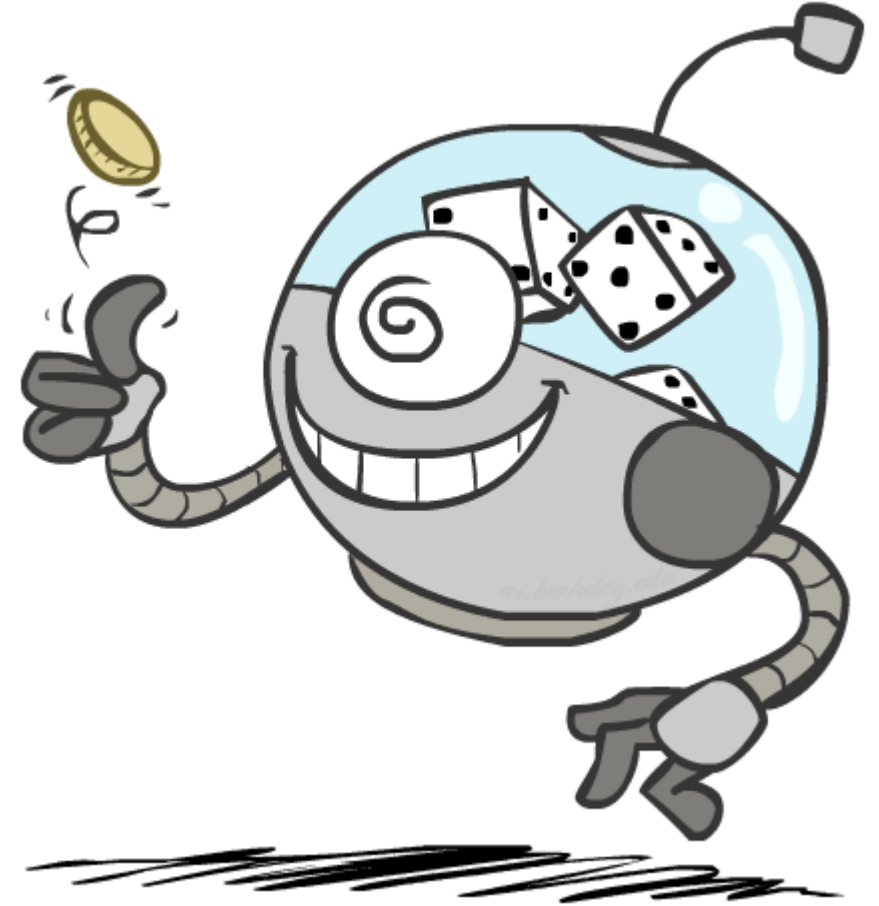
---

- Written assignment 3: MDPs and Reinforcement Learning
  - Deadline: Nov 8th, 2023

# Today

---

- Complete Approximate Q-Learning
- Probability
  - Random Variables
  - Joint and Marginal Distributions
  - Conditional Distribution
- You'll need all this stuff A LOT for the next few weeks, so make sure you go over it now!



# 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 =  $(s, a, r, s')$

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

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

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

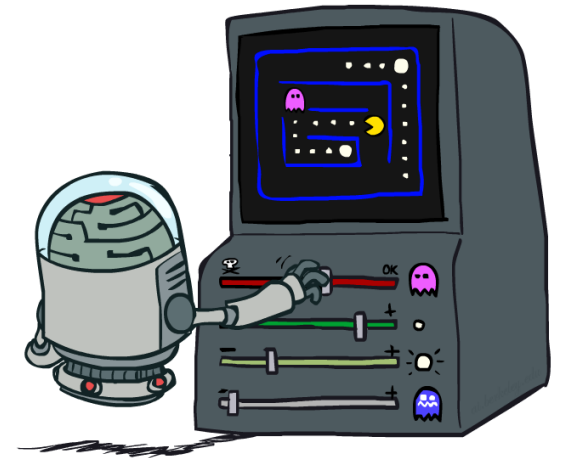
Exact Q's

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 with Linear Approximation

---

**Algorithm 4:** Q-learning with linear approximation.

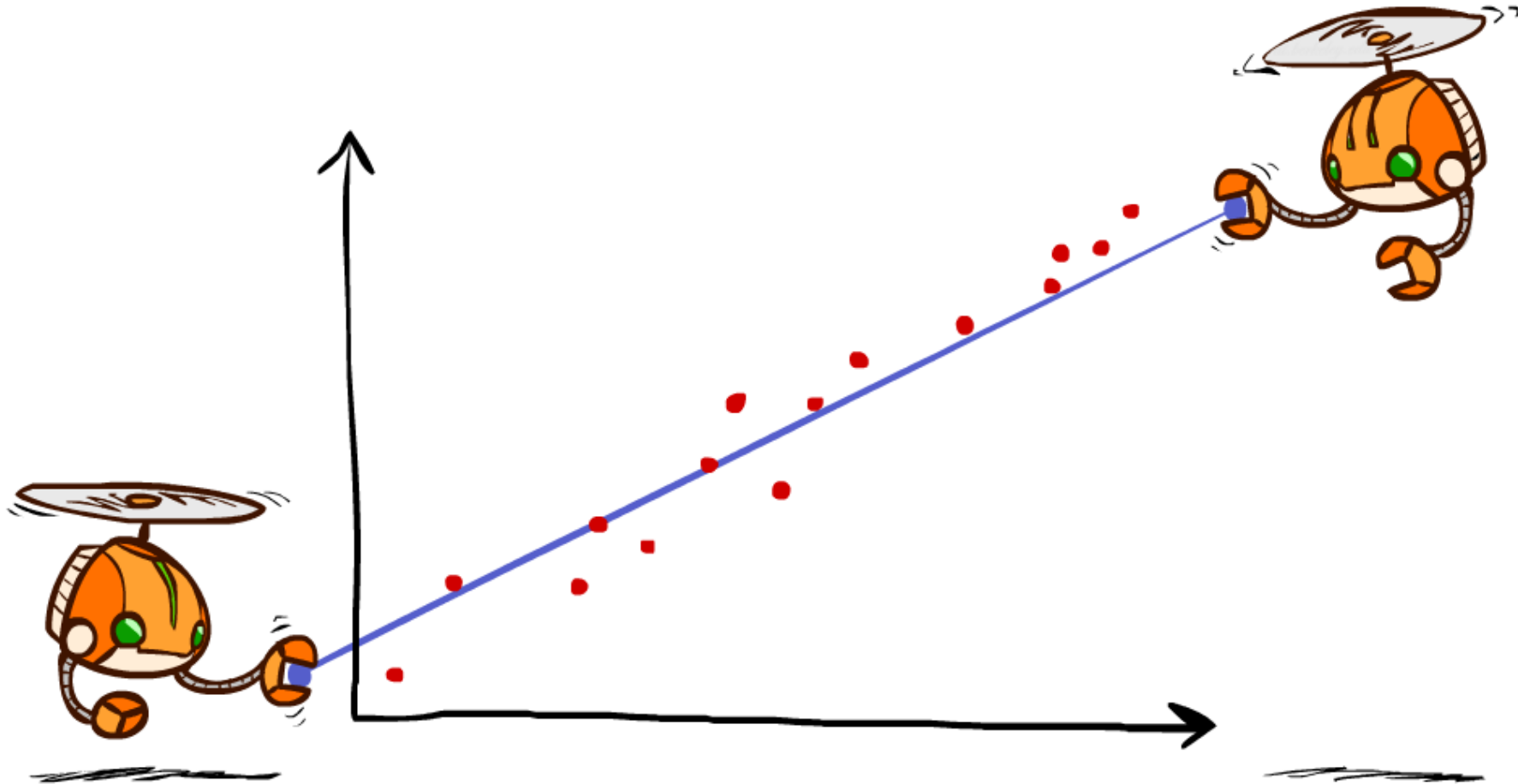
---

```
1 Initialize q-value function  $Q$  with random weights  $w$ :  $Q(s, a; w) = \sum_m w_m f_m(s, a)$ ;  
2 for  $episode = 1 \rightarrow M$  do  
3   Get initial state  $s_0$ ;  
4   for  $t = 1 \rightarrow T$  do  
5     With prob.  $\epsilon$ , select a random action  $a_t$ ;  
6     With prob.  $1 - \epsilon$ , select  $a_t \in \operatorname{argmax}_a Q(s_t, a; w)$ ;  
7     Execute selected action  $a_t$  and observe reward  $r_t$  and next state  $s_{t+1}$ ;  
8     Set target  $y_t = \begin{cases} r_t & \text{if episode terminates at step } t + 1; \\ r_t + \gamma \max_{a'} Q(s_{t+1}, a'; w) & \text{otherwise} \end{cases}$ ;  
9     Perform a gradient descent step to update  $w$ :  $w_m \leftarrow w_m + \alpha [y_t - Q(s_t, a_t; w)] f_m(s, a)$ ;
```

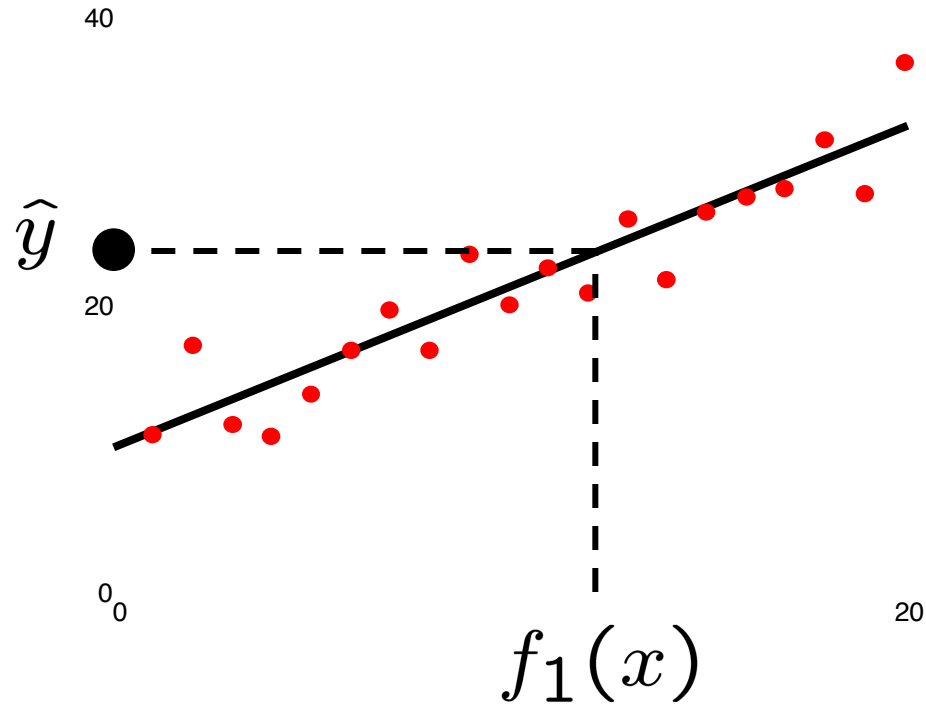
---

# Q-Learning and Least Squares

---

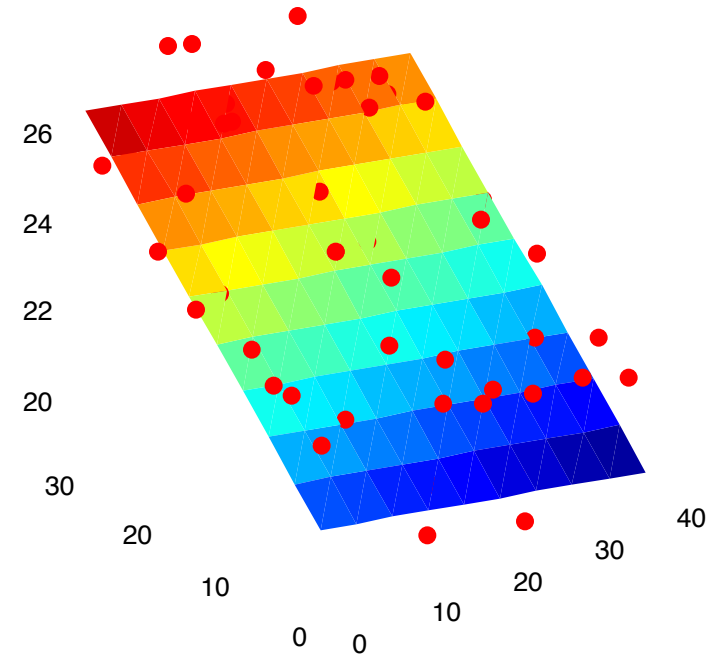


# Linear Approximation: Regression\*



Prediction:

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

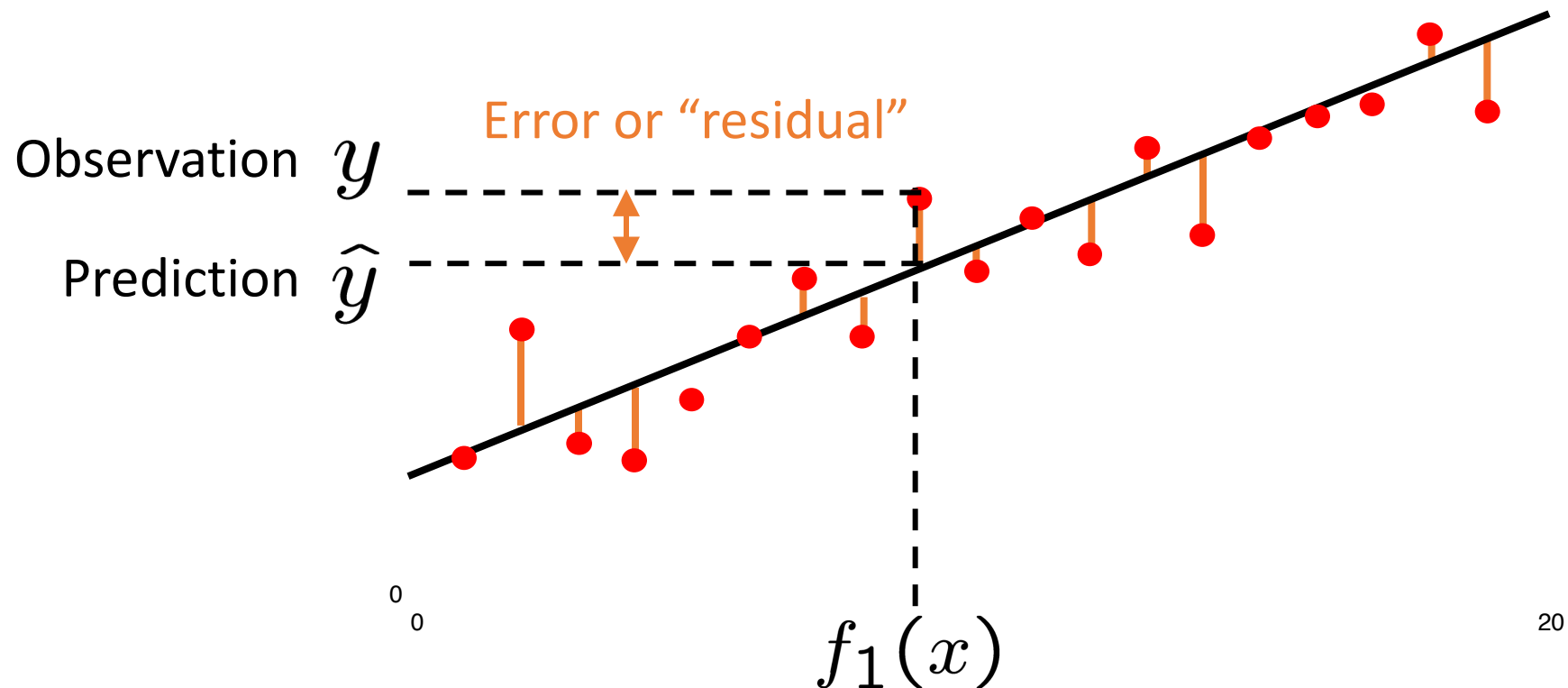


Prediction:

$$\hat{y}_i = w_0 + w_1 f_1(x) + w_2 f_2(x)$$

# Optimization: Least Squares\*

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

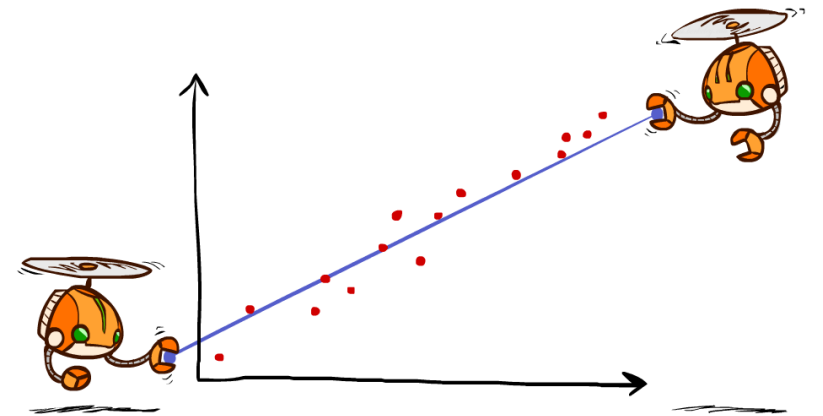




# Minimizing Error\*

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

$$\begin{aligned}\text{error}(w) &= \frac{1}{2} \left( y - \sum_k w_k f_k(x) \right)^2 \\ \frac{\partial \text{error}(w)}{\partial w_m} &= - \left( y - \sum_k w_k f_k(x) \right) f_m(x) \\ w_m &\leftarrow w_m + \alpha \left( y - \sum_k w_k f_k(x) \right) f_m(x)\end{aligned}$$



Approximate q update explained:

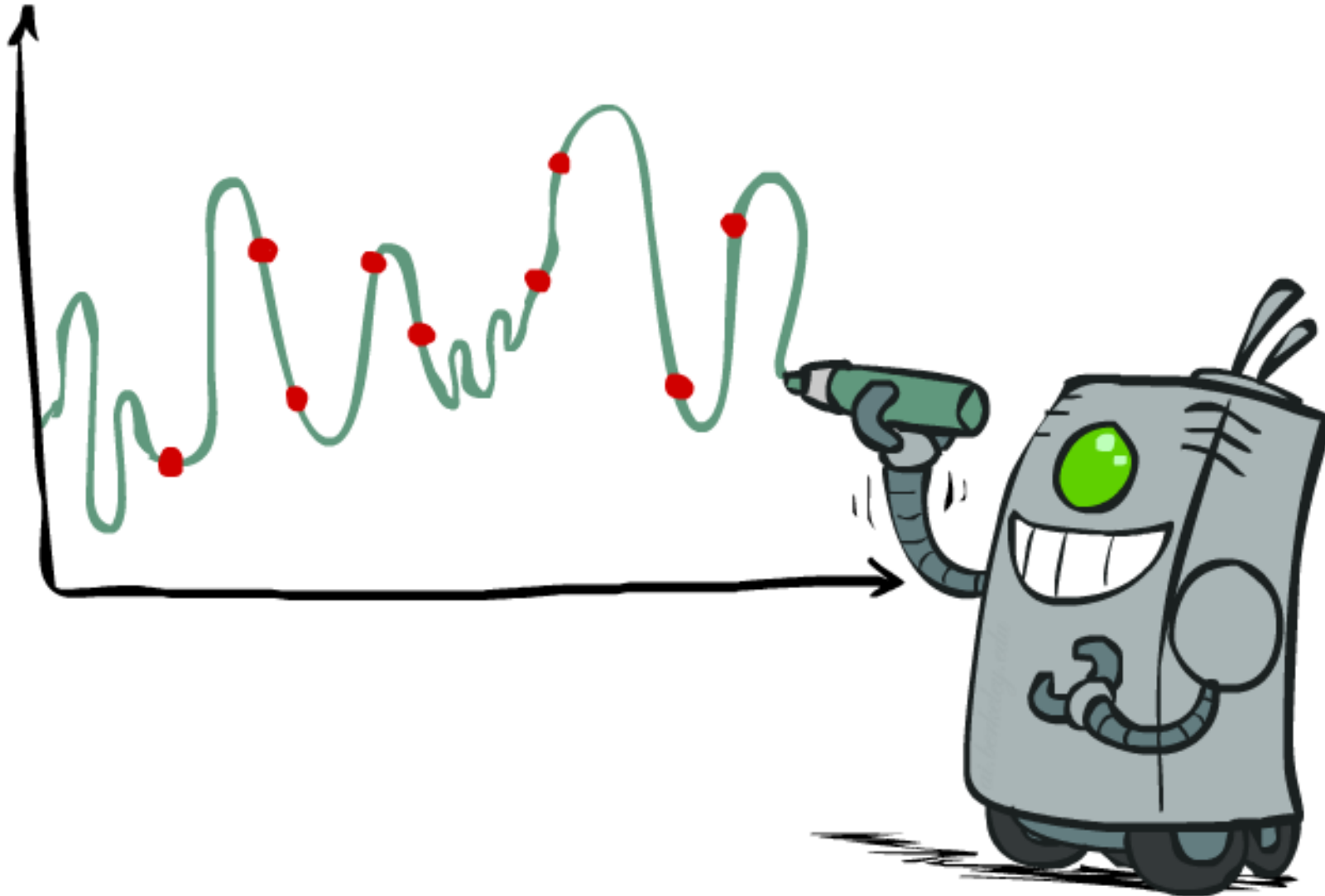
$$w_m \leftarrow w_m + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right] f_m(s, a)$$

“target”

“prediction”

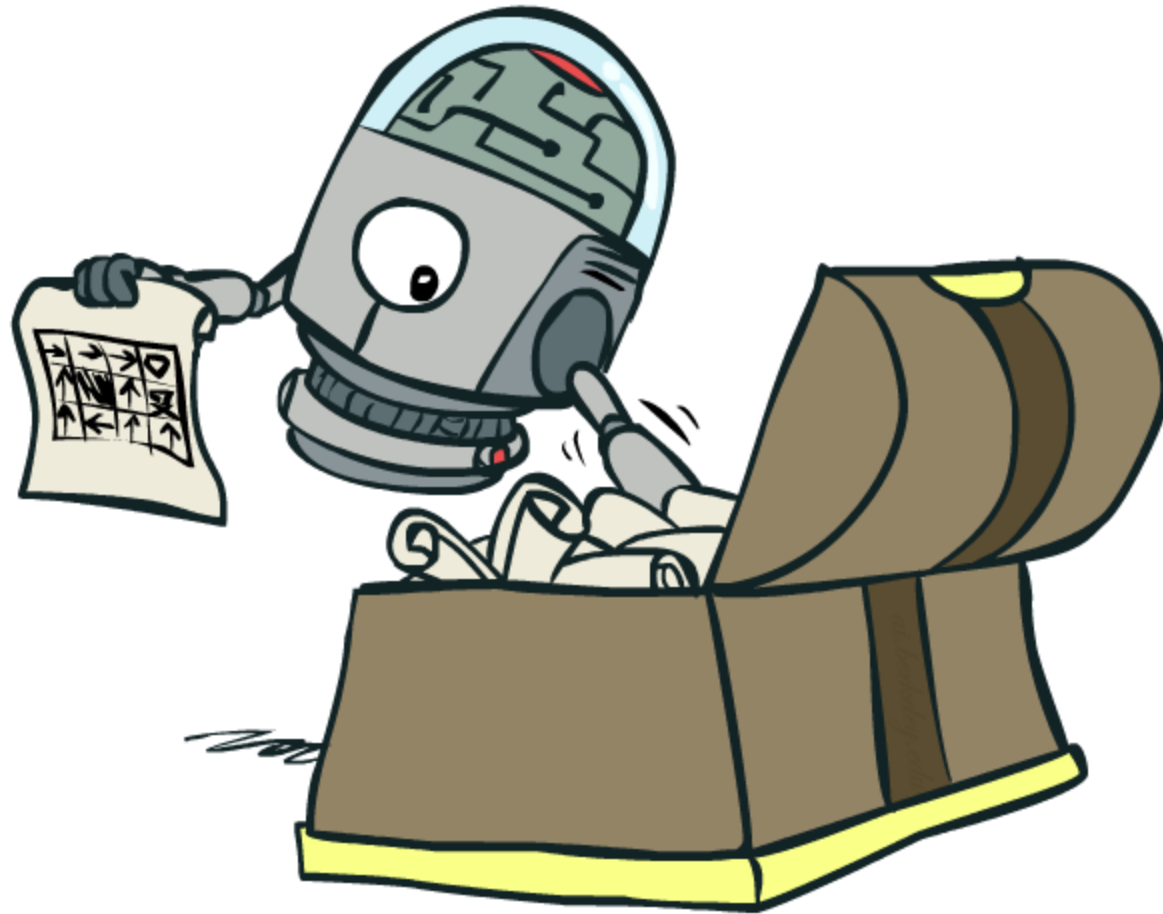


# Overfitting: Why Limiting Capacity Can Help\*



# Policy Search

---



# Policy Search

---

- 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 lookahead structure, sample wisely, change multiple parameters...



# Conclusion

- We're done with Part I: Search and Planning!
- We've seen how AI methods can solve problems in:
  - Search
  - Constraint Satisfaction Problems
  - Games
  - Markov Decision Problems
  - Reinforcement Learning
- Next up: Part II: Uncertainty and Learning!



# Uncertainty

- General situation:
  - **Observed variables (evidence):** Agent knows certain things about the state of the world (e.g., sensor readings or symptoms)
  - **Unobserved variables:** Agent needs to reason about other aspects (e.g. where an object is or what disease is present)
  - **Model:** Agent knows something about how the known variables relate to the unknown variables
- Probabilistic reasoning gives us a framework for managing our beliefs and knowledge

0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

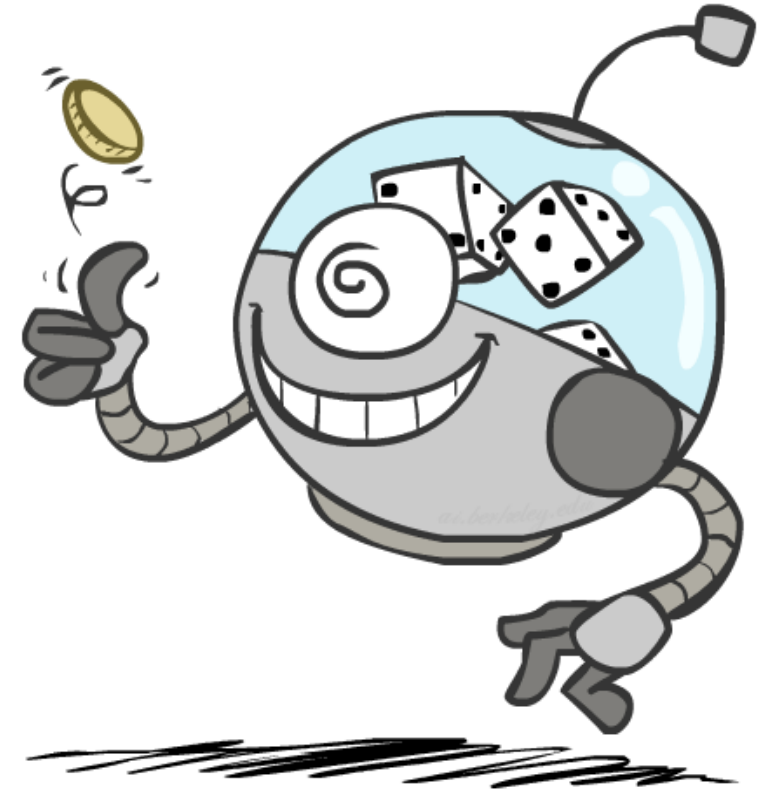
0.17	0.10	0.10
0.09	0.17	0.10
<0.01	0.09	0.17

<0.01	<0.01	0.03
<0.01	0.05	0.05
<0.01	0.05	0.81



# Random Variables

- A random variable is some aspect of the world about which we (may) have uncertainty
  - $R$  = Is it raining?
  - $T$  = Is it hot or cold?
  - $D$  = How long will it take to drive to work?
  - $L$  = Where is the ghost?
- We denote random variables with capital letters
- Like variables in a CSP, random variables have domains
  - $R$  in  $\{\text{true}, \text{false}\}$  (often write as  $\{+r, -r\}$ )
  - $T$  in  $\{\text{hot}, \text{cold}\}$
  - $D$  in  $[0, \infty)$
  - $L$  in possible locations, maybe  $\{(0,0), (0,1), \dots\}$

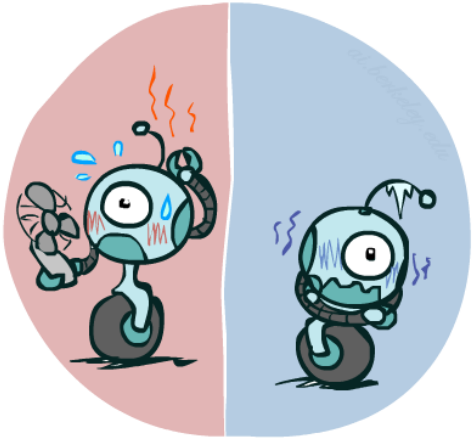




# Probability Distributions

- Associate a probability with each value

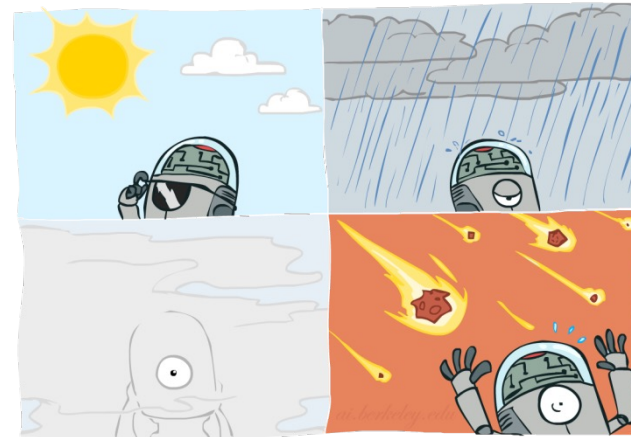
- Temperature:



$P(T)$

T	P
hot	0.5
cold	0.5

- Weather:



$P(W)$

W	P
sun	0.6
rain	0.1
fog	0.3
meteor	0.0



# Probability Distributions

- Unobserved random variables have distributions

$P(T)$		$P(W)$	
T	P	W	P
hot	0.5	sun	0.6
cold	0.5	rain	0.1
		fog	0.3
		meteor	0.0

Shorthand notation:

$$P(hot) = P(T = hot),$$

$$P(cold) = P(T = cold),$$

$$P(rain) = P(W = rain),$$

...

OK if all domain entries are unique

- A distribution is a TABLE of probabilities of values
- A probability (lower case value) is a single number

$$P(W = rain) = 0.1$$

- Must have:  $\forall x \ P(X = x) \geq 0$  and  $\sum_x P(X = x) = 1$



# Joint Distributions

- A *joint distribution* over a set of random variables:  $X_1, X_2, \dots, X_n$  specifies a real number for each assignment (or *outcome*):

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

$$P(x_1, x_2, \dots, x_n)$$

- Must obey:  $P(x_1, x_2, \dots, x_n) \geq 0$

$$\sum_{(x_1, x_2, \dots, x_n)} P(x_1, x_2, \dots, x_n) = 1$$

- Size of distribution if  $n$  variables with domain sizes  $d$ ?
  - For all but the smallest distributions, impractical to write out!

$$P(T, W)$$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

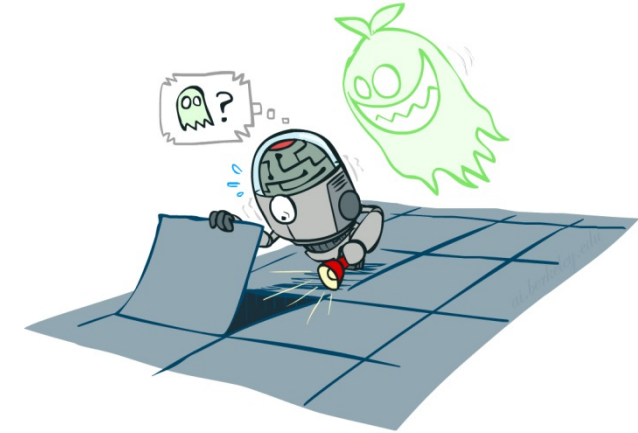


# Probabilistic Models

- A probabilistic model is a joint distribution over a set of random variables
- Probabilistic models:
  - (Random) variables with domains
  - Assignments are called *outcomes*
  - Joint distributions: say whether assignments (outcomes) are likely
  - *Normalized*: sum to 1.0
  - Ideally: only certain variables directly interact
- Constraint satisfaction problems:
  - Variables with domains
  - Constraints: state whether assignments are possible
  - Ideally: only certain variables directly interact

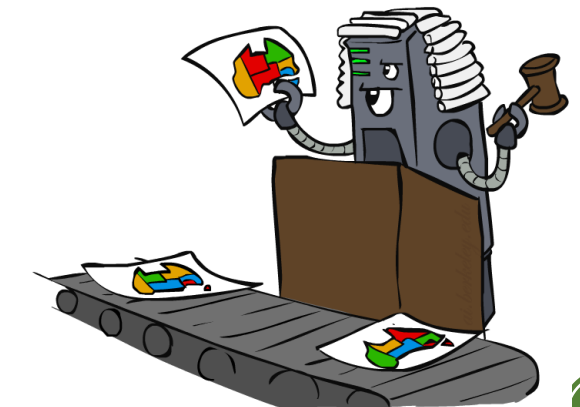
Distribution over T,W

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3



Constraint over T,W

T	W	P
hot	sun	T
hot	rain	F
cold	sun	F
cold	rain	T



# Events

- An *event* is a set  $E$  of outcomes

$$P(E) = \sum_{(x_1 \dots x_n) \in E} P(x_1 \dots x_n)$$

- From a joint distribution, we can calculate the probability of any event
  - Probability that it's hot AND sunny?
  - Probability that it's hot?
  - Probability that it's hot OR sunny?
- Typically, the events we care about are *partial assignments*, like  $P(T=\text{hot})$

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3



# Quiz: Events

■  $P(+x, +y)$  ?

■  $P(+x)$  ?

■  $P(-y \text{ OR } +x)$  ?

$P(X, Y)$

X	Y	P
+x	+y	0.2
+x	-y	0.3
-x	+y	0.4
-x	-y	0.1

