

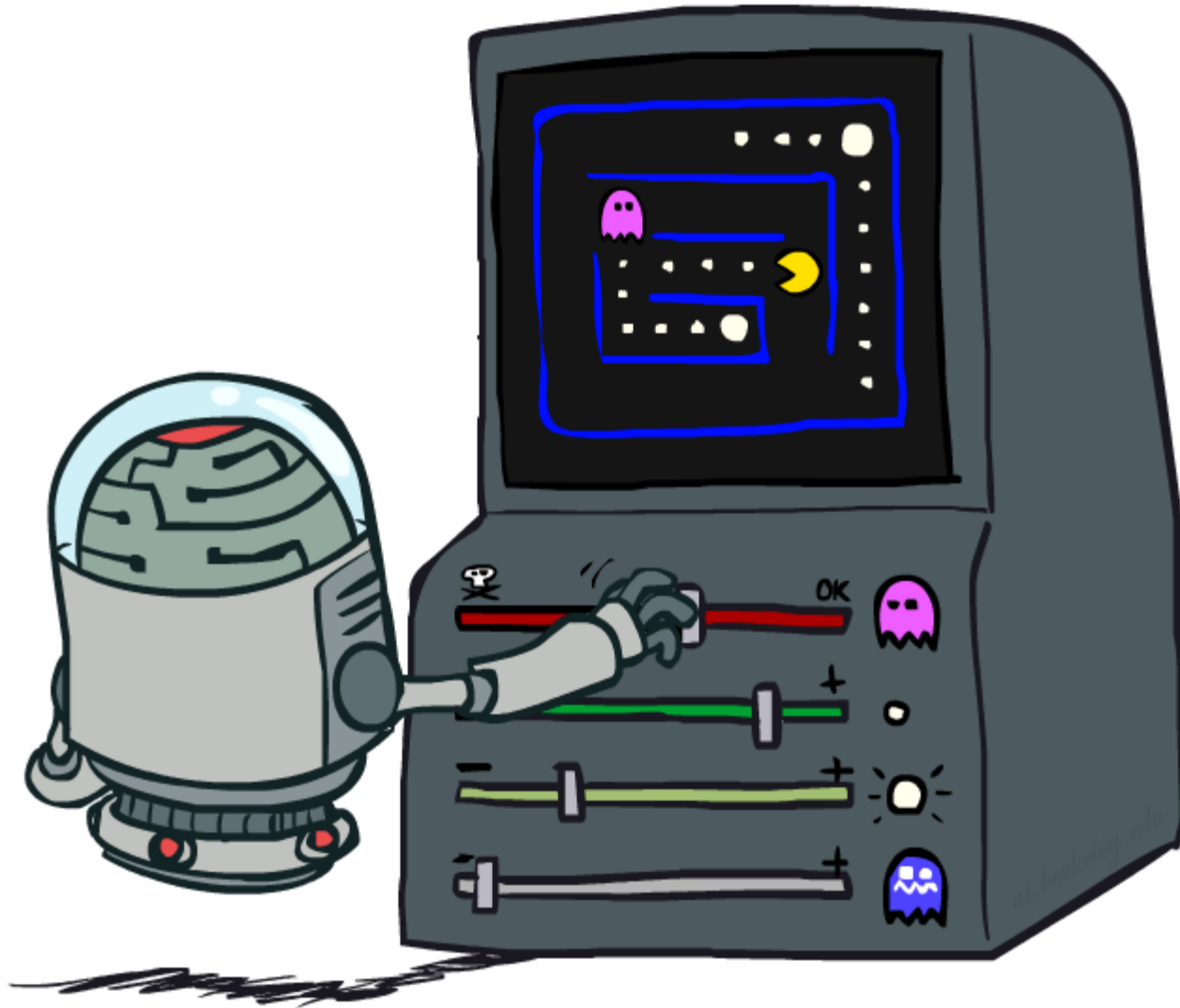
## Search with Uncertainty



# Reinforcement Learning

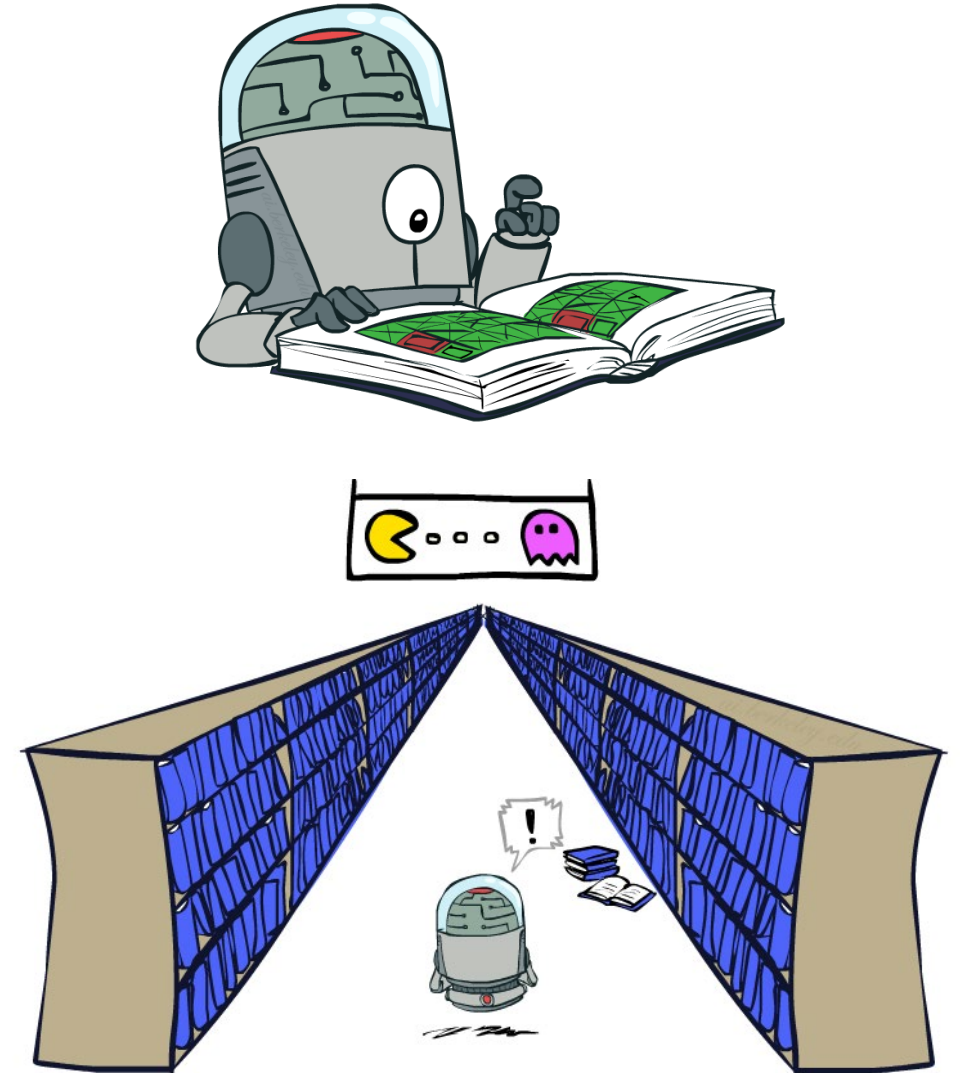
# Review- Stories so far

# Approximate Q-Learning



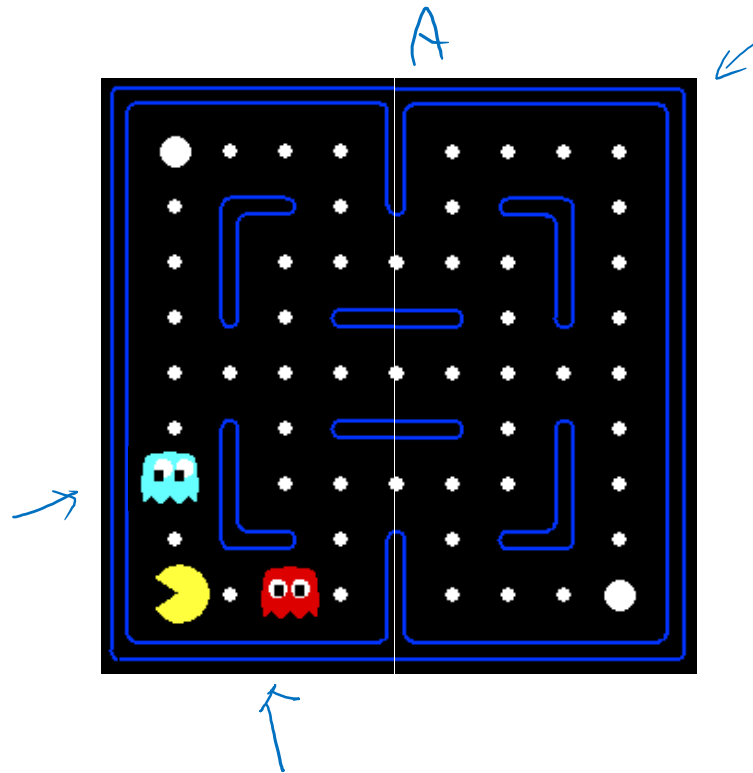
# Generalizing Across States

- 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

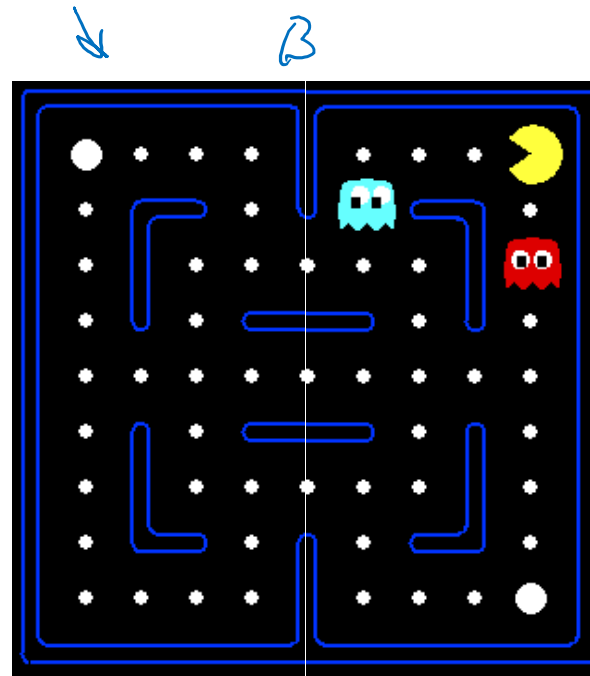


# Example: Pacman

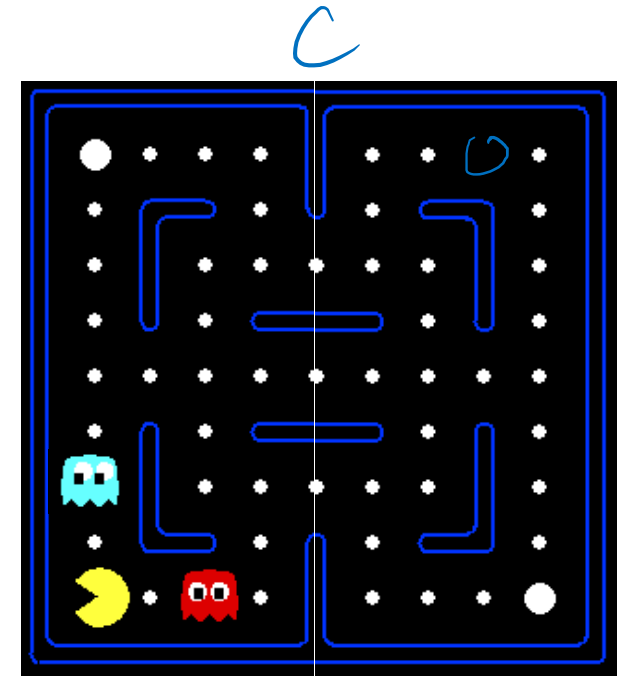
Let's say we discover through experience that this state is bad:



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



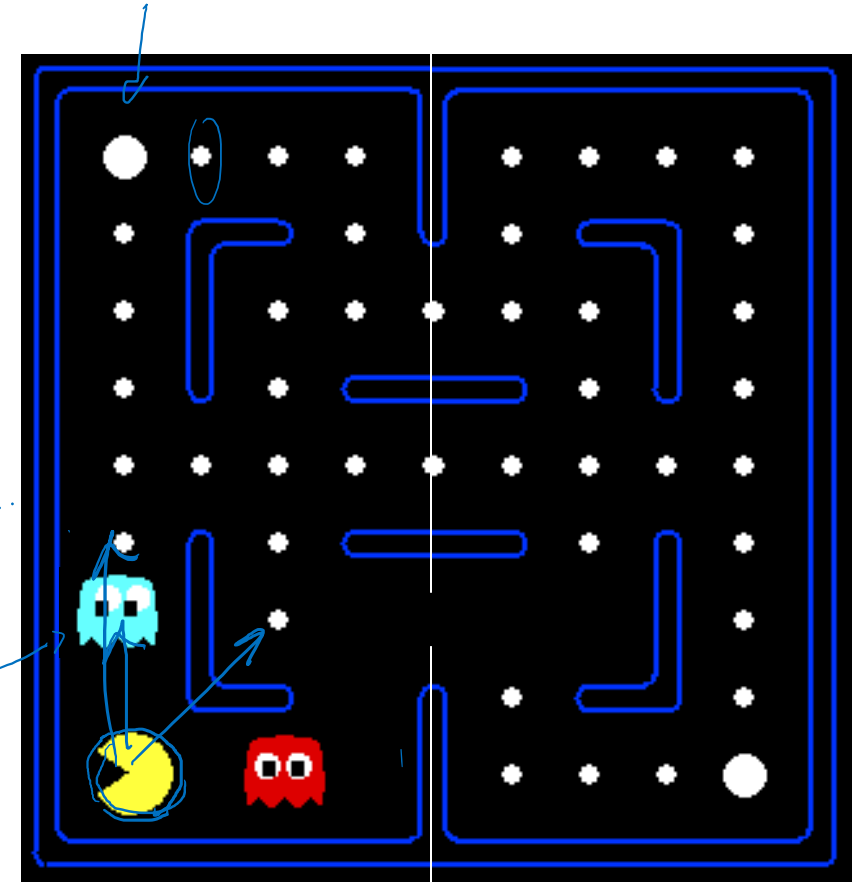
Or even this one!



# 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)

$$V(s) = F(x_1, x_2, \dots, x_n)$$



# Linear Value Functions

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

$$\begin{aligned} V(s) &= w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s) \\ Q(s, a) &= w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a) \end{aligned}$$

Weights = parameters

Linear regression

$s$

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

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

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

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

TD-control (off-policy)

Exact Q's

Approximate Q's

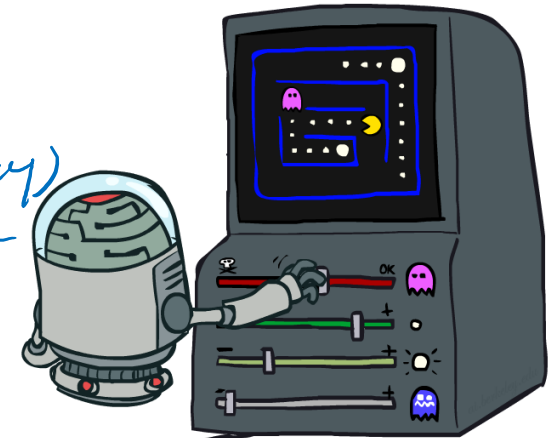
Weight update rule

- 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

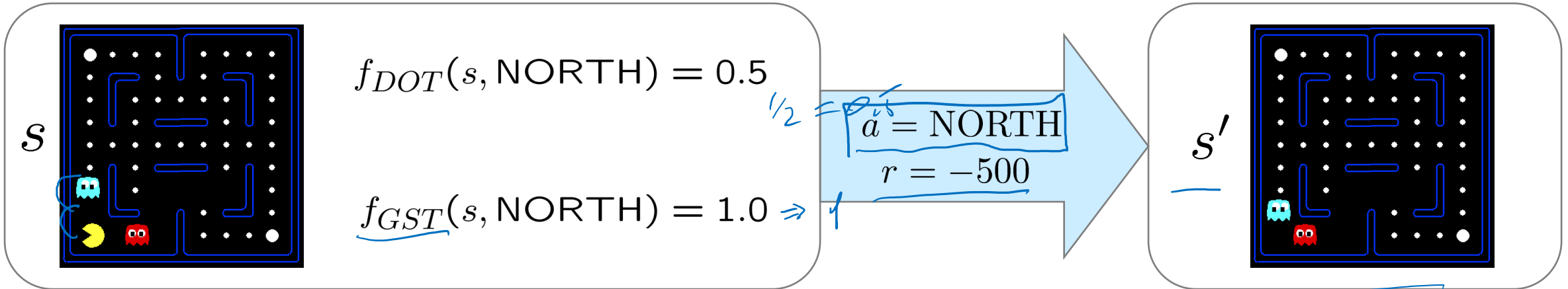
- Formal justification: online least squares

stepsize optimization



# Example: Q-Pacman

$$Q(s, a) = \overset{w_{\text{dot}}}{4.0} f_{\text{DOT}}(s, a) - \overset{w_g}{1.0} f_{\text{GST}}(s, a)$$



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

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

$Q_{\text{new}} - Q_{\text{old}} = -501 \quad 0 < \alpha < 1$

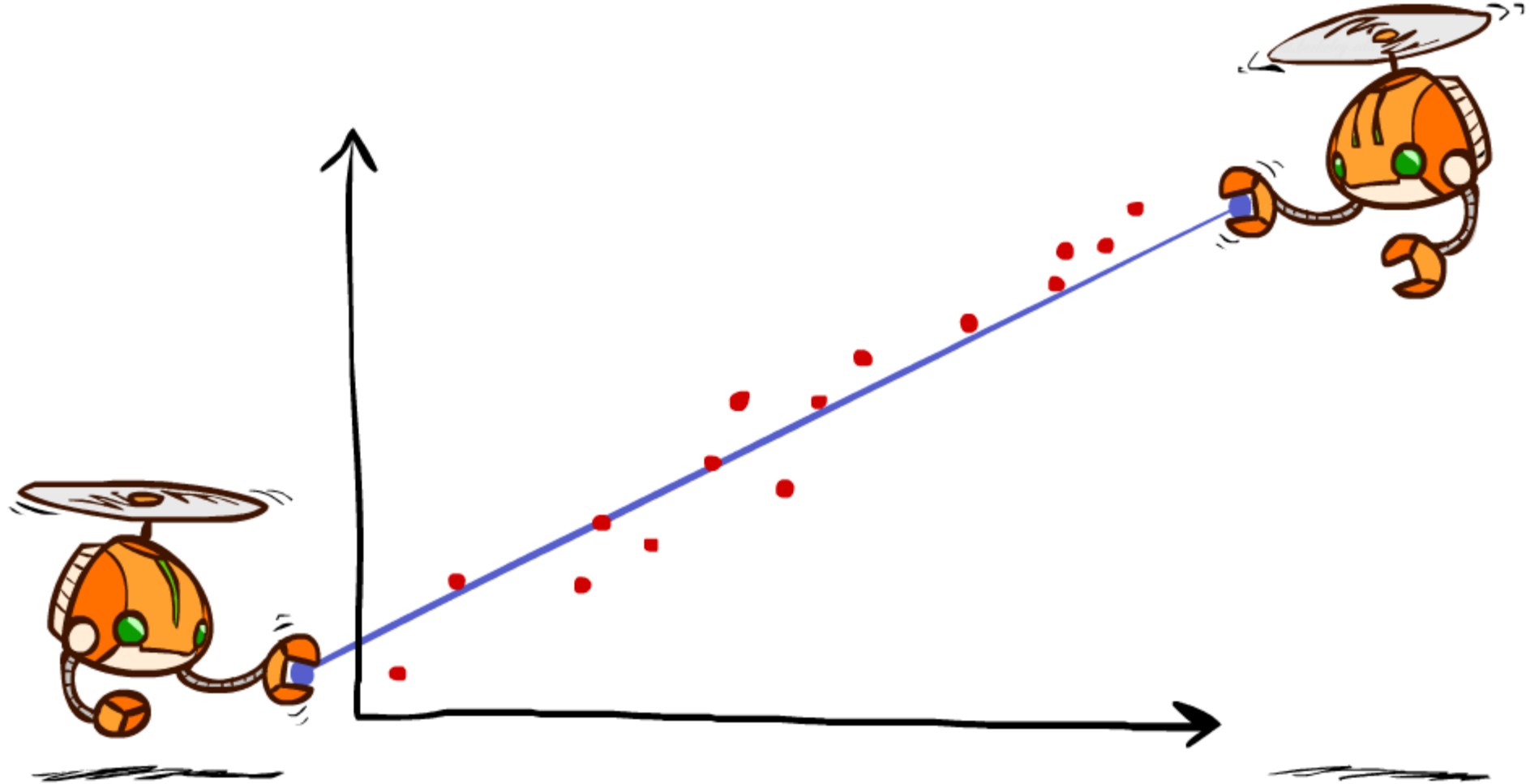
difference = -501

$w_{\text{DOT}} \leftarrow 4.0 + \alpha [-501] 0.5$

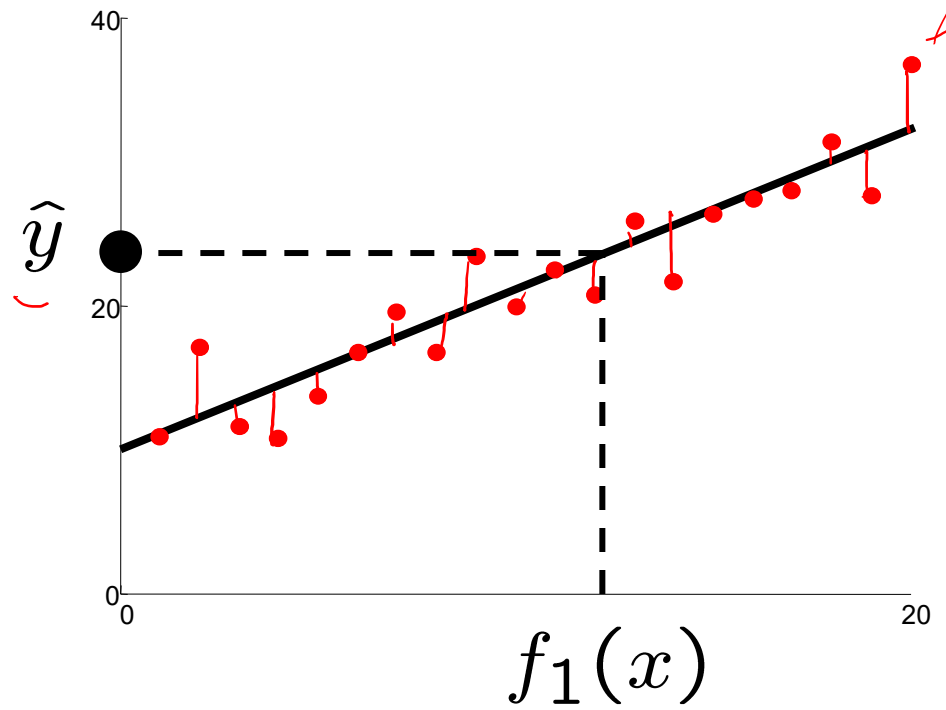
$w_{\text{GST}} \leftarrow -1.0 + \alpha [-501] 1.0$

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

# Q-Learning and Least Squares



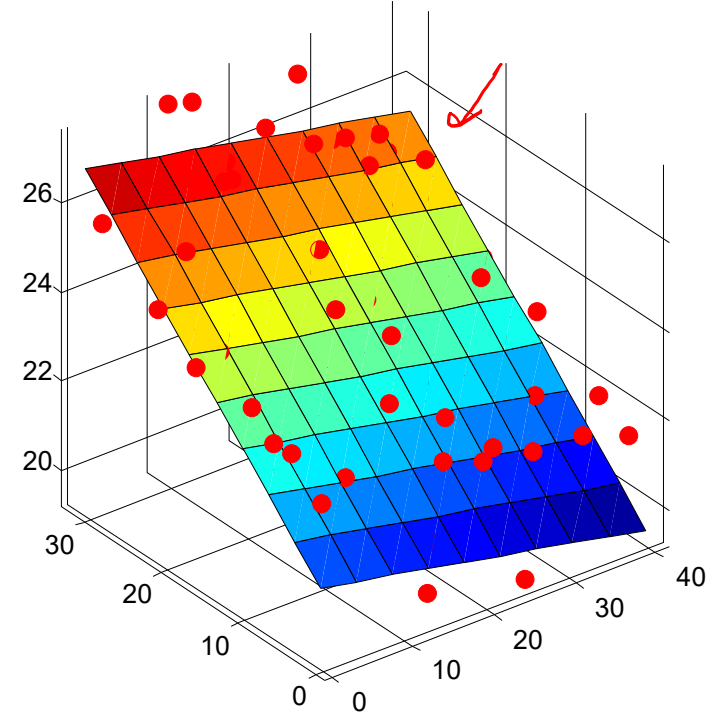
# Linear Approximation: Regression\*



Prediction:

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

$\underbrace{\quad\quad}_b$



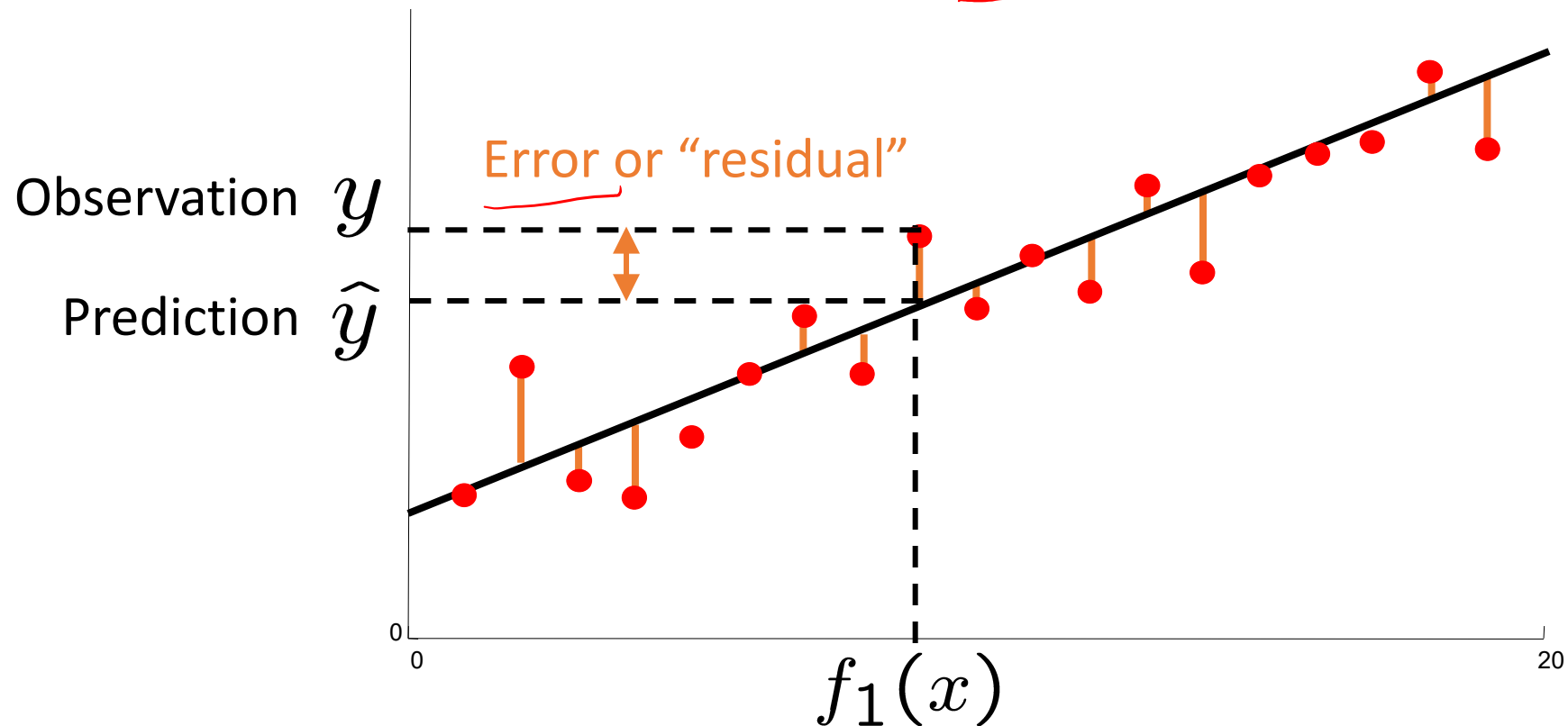
Prediction:

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

# Optimization: Least Squares\*

↓

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



# Weight update rule- Gradient Descent\*

## Least Squares Method (Linear Regression)

Loss fun  $L(W, X) =$  RSS or MSE  
(J)  $RSS = \sum_{\hat{x}=1}^N (y_i - \hat{y}_i)^2$

$$MSE = \frac{1}{N} \cdot RSS$$

$$\rightarrow w_j^{\text{new}} \leftarrow w_j^{\text{old}} - \alpha \frac{\partial L}{\partial w_j}$$

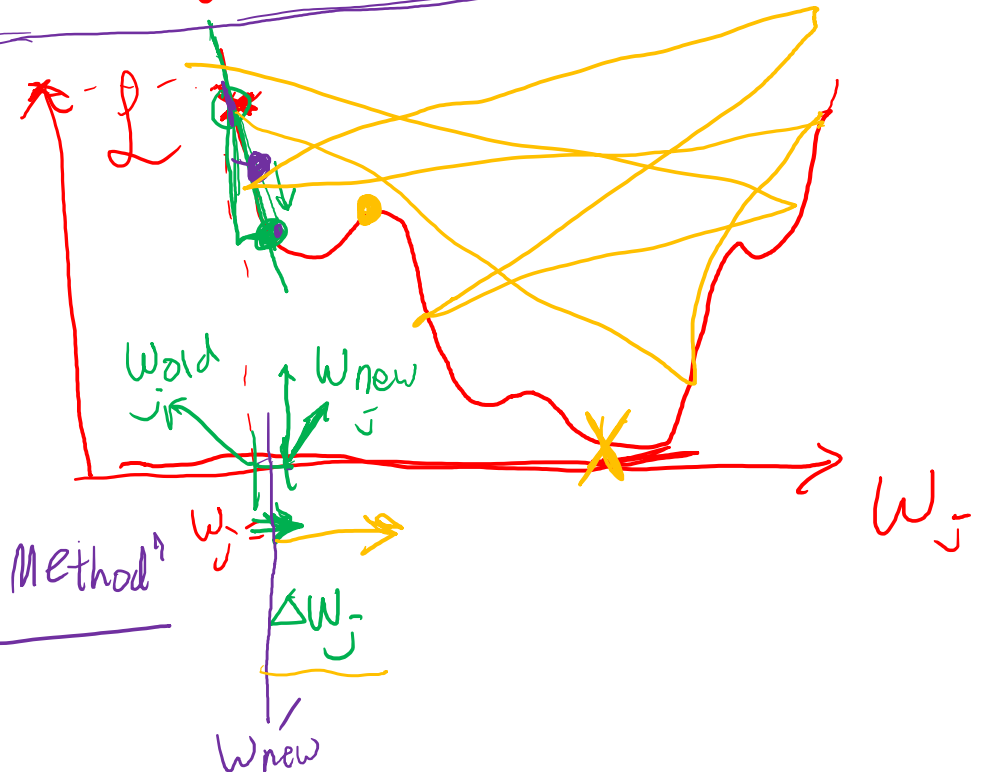
$$W_j^{\text{new}} \leftarrow W_j^{\text{old}} + \alpha (y - \hat{y}) \cdot X_j$$

## "Gradient Descent Method"

→ ML

$$\underline{f(w, x) \rightarrow \hat{y} \in \mathbb{R}}$$

$$\frac{\partial \mathcal{L}}{\partial w_j} = \sum_i 2 \cdot (y_i - \sum_j w_j x_{ji}) \cdot (-x_{ji})$$



# Minimizing Error\*

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

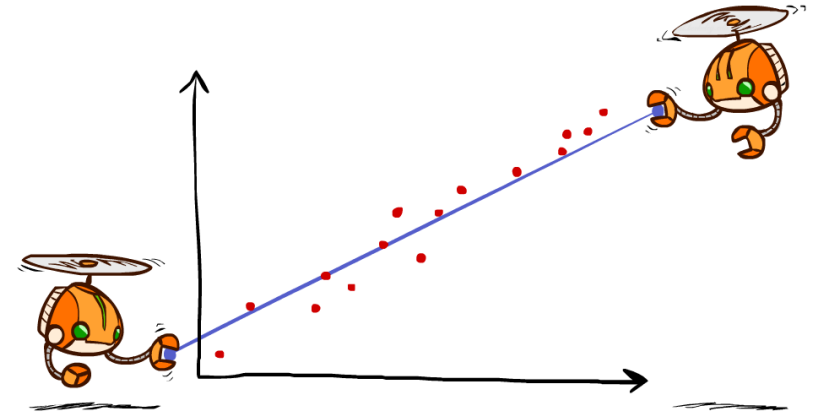
$$\text{error}(w) = \frac{1}{2} \left( y - \sum_k w_k f_k(x) \right)^2$$

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

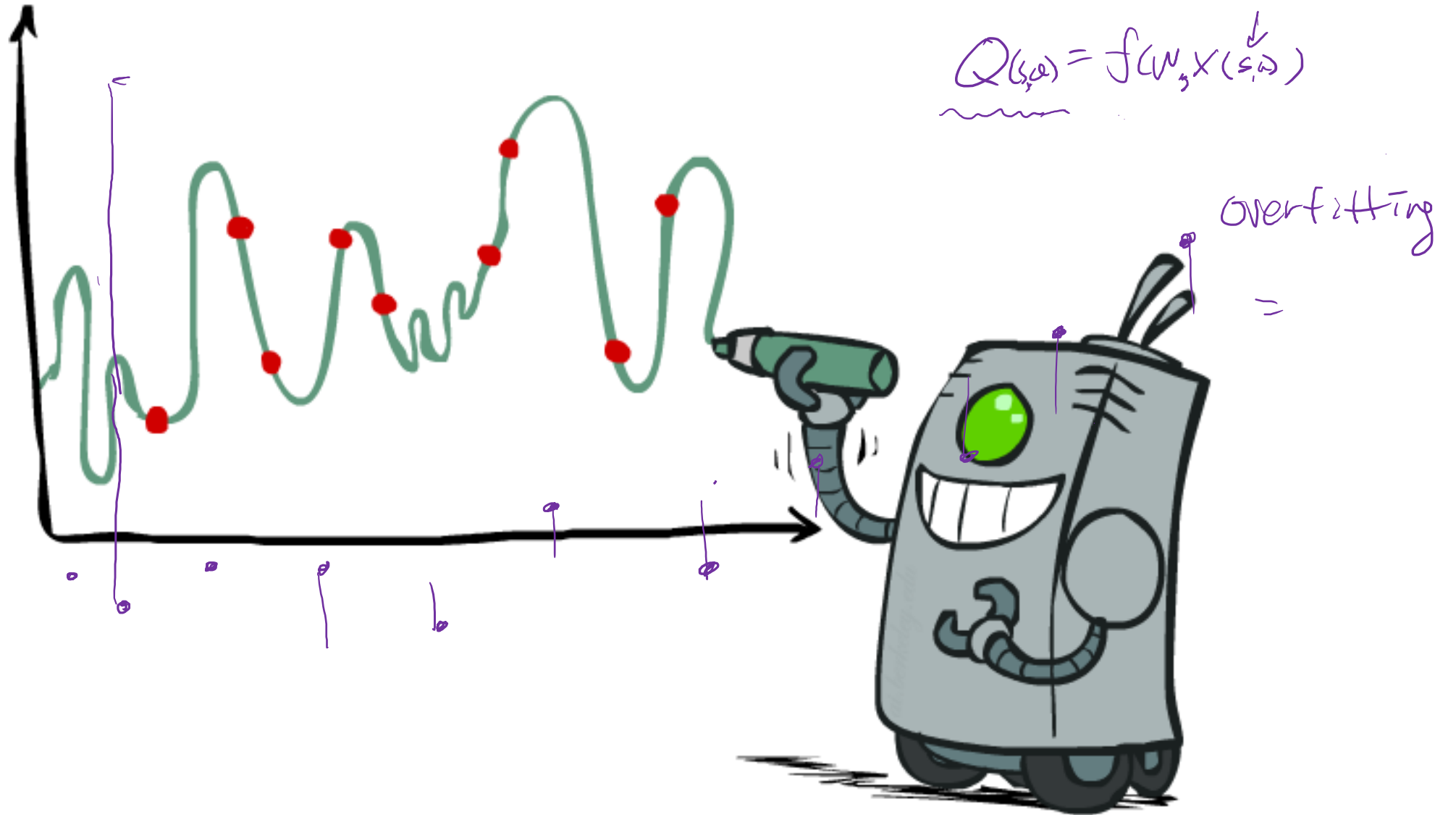
$$\underline{w_m} \leftarrow w_m + \alpha \left( y - \sum_k w_k f_k(x) \right) f_m(x)$$

Approximate q update explained:

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

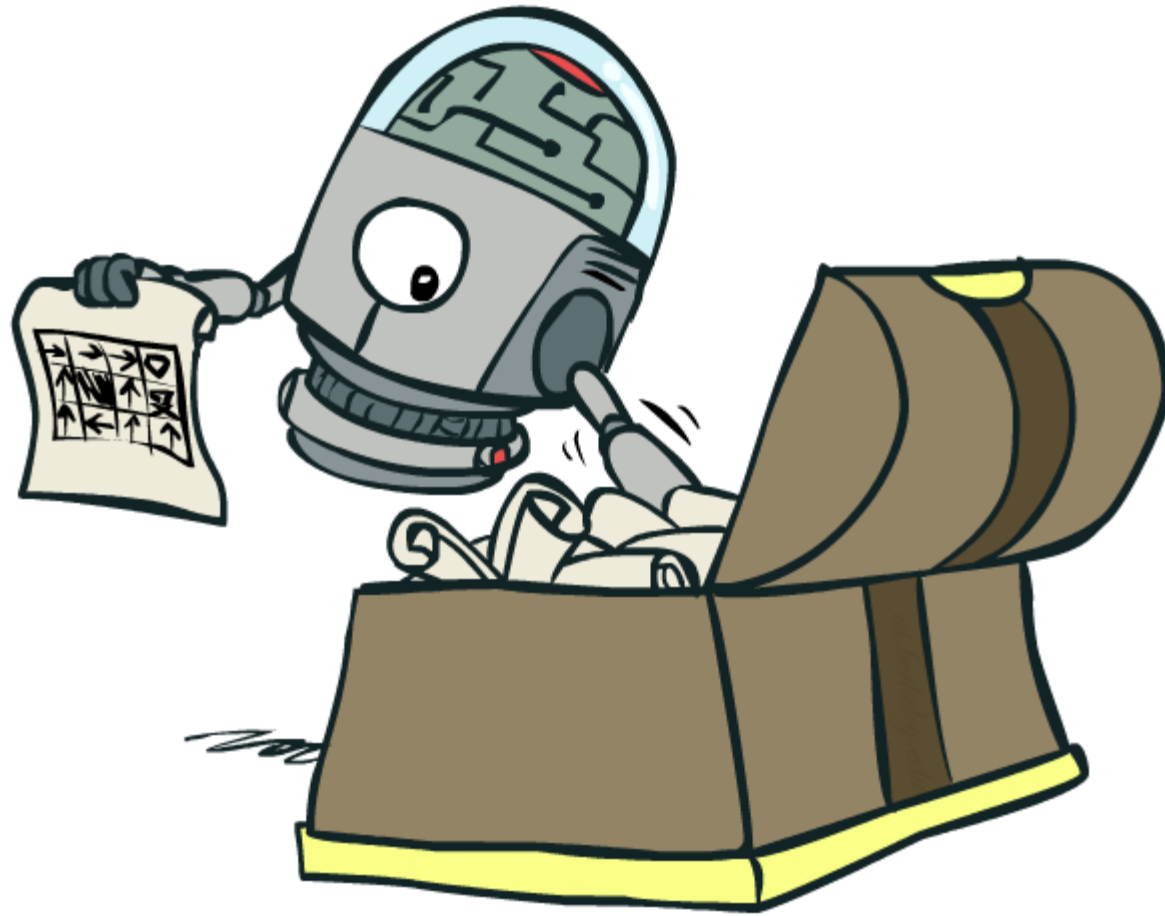


# Overfitting: Why Limiting Capacity Can Help\*



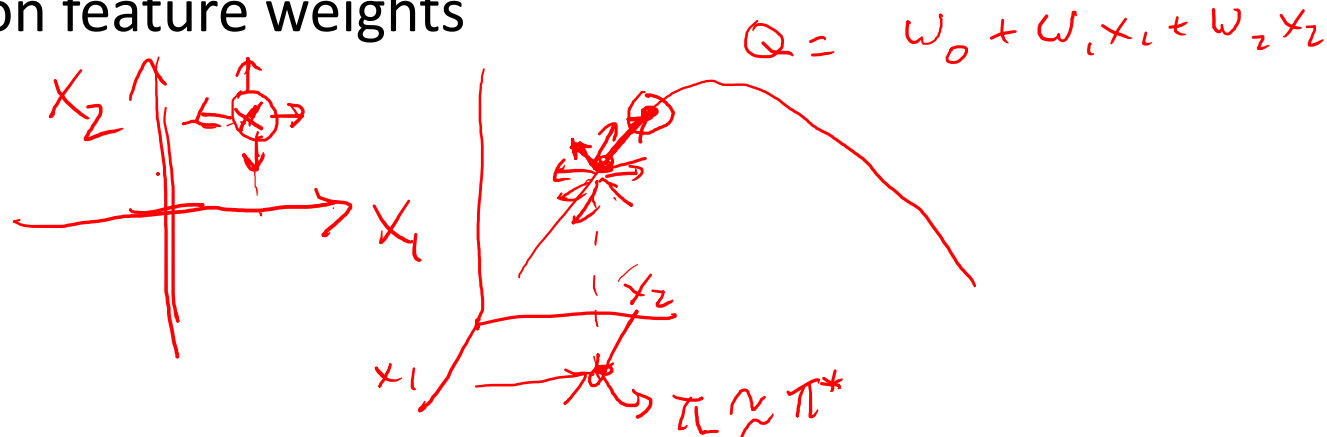


# Policy Search



# Policy Search

- Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate  $V / Q$  best
  - Q-learning's priority: get Q-values close (modeling)
  - Action selection priority: get ordering of Q-values right (prediction)
- 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

$N$  time  $a_j$

$$Q = \sum_j w_j x_j$$

$\pi \rightarrow ?$

$$\pi = P(a|s) \quad \begin{matrix} x = f(s) \\ \uparrow \end{matrix}$$

$$\pi_w(s, a) = \text{Softmax} = \frac{e^{\hat{Q}_w(s, a)}}{\sum_{a'} e^{\hat{Q}_w(s, a')}} \quad \leftarrow$$

$$a \leftarrow \{a'\}$$

$$\frac{ds}{f} = d(\log f)$$

hill climbing  $\rightarrow \frac{1}{N} \sum_j (\nabla_w \log \pi_w) R$

policy gradient

$P \leftarrow$  policy value  $\neq Q$

REINFORCE

$$\nabla_w P(w) = \nabla_w \left( \sum_a \pi_w(a|s) R(a) \right)$$

$$= \sum_a \frac{\nabla_w(\pi_w) R}{\pi_w} \cdot \pi_w = \sum_a \nabla_w \log \pi_w \cdot R$$

# Policy Search



# Conclusion

- e

