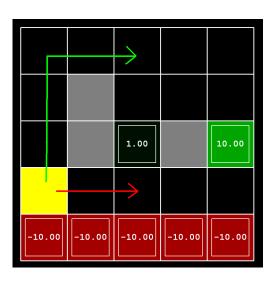
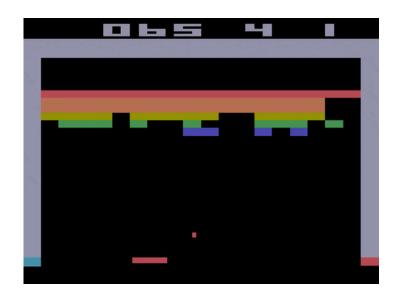
Can tabular methods scale?

Discrete environments





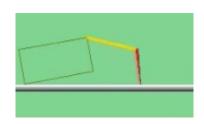
Gridworld 10^1

Tetris 10^60

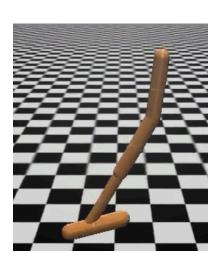
Atari 10^308 (ram) 10^16992 (pixels)

Can tabular methods scale?

Continuous environments (by crude discretization)



Crawler 10^2



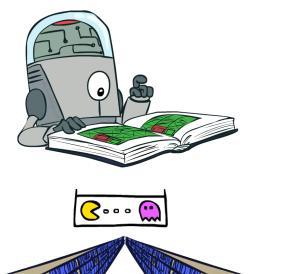
Hopper 10¹0



Humanoid 10^100

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





Approximate Q-Learning

- ullet Instead of a table, we have a parametrized Q function: $Q_{ heta}(s,a)$
 - Can be a linear function in features:

$$Q_{\theta}(s,a) = \theta_0 f_0(s,a) + \theta_1 f_1(s,a) + \dots + \theta_n f_n(s,a)$$

- Or a complicated neural net
- Learning rule:
 - Remember: $target(s') = R(s, a, s') + \gamma \max_{a'} Q_{\theta_k}(s', a')$
 - Update:

$$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta} \left[\frac{1}{2} (Q_{\theta}(s, a) - \text{target}(s'))^2 \right] \Big|_{\theta = \theta_k}$$

Connection to Tabular Q-Learning

• Suppose $\theta \in \mathbb{R}^{|S| \times |A|}, \quad Q_{\theta}(s, a) \equiv \theta_{sa}$

$$\nabla_{\theta_{sa}} \left[\frac{1}{2} (Q_{\theta}(s, a) - \text{target}(s'))^{2} \right]$$

$$= \nabla_{\theta_{sa}} \left[\frac{1}{2} (\theta_{sa} - \text{target}(s'))^{2} \right]$$

$$= \theta_{sa} - \text{target}(s')$$

Plug into update: $\theta_{sa} \leftarrow \theta_{sa} - \alpha(\theta_{sa} - \mathrm{target}(s'))$ $= (1 - \alpha)\theta_{sa} + \alpha[\mathrm{target}(s')]$

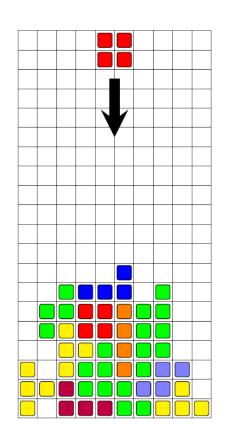
Compare with Tabular Q-Learning update:

$$Q_{k+1}(s,a) \leftarrow (1-\alpha)Q_k(s,a) + \alpha \left[\operatorname{target}(s') \right]$$

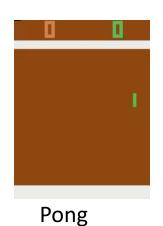
Engineered Approximation Example: Tetris

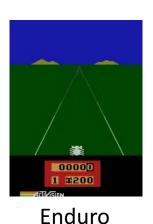
- state: naïve board configuration + shape of the falling piece ~10⁶⁰ states!
- action: rotation and translation applied to the falling piece
- ullet 22 features aka basis functions $\,\phi_i$
 - Ten basis functions, $0, \ldots, 9$, mapping the state to the height h[k] of each column.
 - Nine basis functions, $10, \ldots, 18$, each mapping the state to the absolute difference between heights of successive columns: |h[k+1] h[k]|, $k = 1, \ldots, 9$.
 - One basis function, 19, that maps state to the maximum column height: $\max_k h[k]$
 - One basis function, 20, that maps state to the number of 'holes' in the board.
 - One basis function, 21, that is equal to 1 in every state.

$$\hat{V}_{\theta}(s) = \sum_{i=0}^{21} \theta_i \phi_i(s) = \theta^{\top} \phi(s)$$

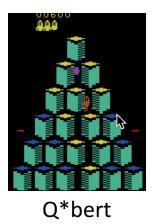


Deep Reinforcement Learning









From pixels to actions

- Same algorithm (with effective tricks)
- CNN function approximator, w/ 3M free parameters

Recap: Approximate Q-Learning

```
Algorithm:
      Start with Q_0(s,a) for all s, a.
       Get initial state s
       For k = 1, 2, ... till convergence
              Sample action a, get next state s'
                                                         Chasing a nonstationary target!
              If s' is terminal:
                   target = R(s, a, s')
                   Sample new initial state s'
              else:
                   target = R(s, a, s') + \gamma \max_{s'} Q_k(s', a')
             \theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta} \mathbb{E}_{s' \sim P(s'|s,a)} \left[ (Q_{\theta}(s,a) - \text{target}(s'))^2 \right] \Big|_{\theta = \theta}
              s \leftarrow s'
                                       Updates are correlated within a trajectory!
```

DQN

- High-level idea make Q-learning look like supervised learning.
- Two main ideas for stabilizing Q-learning.
- Apply Q-updates on batches of past experience instead of online:
 - Experience replay (Lin, 1993).
 - Previously used for better data efficiency.
 - Makes the data distribution more stationary.
- Use an older set of weights to compute the targets (target network):
 - Keeps the target function from changing too quickly.

$$L_i(\theta_i) = \mathbb{E}_{s,a,s',r \sim D} \left(\underbrace{r + \gamma \ \max_{a'} Q(s', a'; \boldsymbol{\theta}_i^-)}_{\text{target}} - Q(s, a; \theta_i) \right)^2$$

Target Network Intuition

- Changing the value of one action will change the value of other actions and similar states.
- The network can end up chasing its own tail because of bootstrapping.
- Somewhat surprising fact bigger networks are less prone to this because they alias less.

$$L_i(\theta_i) = \mathbb{E}_{s,a,s',r \sim D} \left(\underbrace{r + \gamma \, \max_{a'} Q(s', a'; \boldsymbol{\theta_i^-})}_{\text{target}} - Q(s, a; \theta_i) \right)^2$$







$$y' = \sum \omega_1 x_1$$

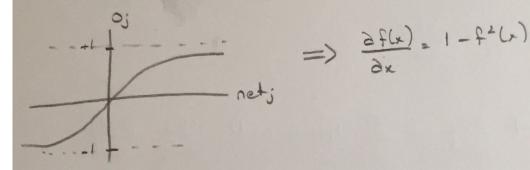
 $E = \frac{1}{2} (\hat{y} - y')^2$

$$\Delta \omega_{i} = -\alpha \cdot \frac{\partial E}{\partial \omega_{i}} = -\alpha \cdot (\hat{g} - \hat{g}') \cdot (-1) \cdot \frac{\partial g'}{\partial \omega_{i}}$$

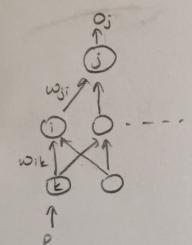
$$= \alpha \cdot (\hat{g} - \hat{g}') \cdot \times_{i}$$

$$net_j = \sum_i w_j \cdot x_i$$
 $g = \begin{cases} 1 & \text{if } net_j \ge 0 \\ 0 & \text{otherwise} \end{cases}$

For normalized or precise output, we are going to use squashing function f(x) = tanh(x) => 0; = tanh(net;)



$$\Rightarrow \frac{9x}{9t(r)} = 1 - t_{5}(r)$$



$$t_{Pj} = desired output for node j$$
 $t_{Pj} = desired output for node j$
 $t_{Pj} = desired output for node j

 $t_{Pj} = desired output for node j

Oj = tor$$

net
$$j = \sum_{k} \omega_{jk} \cdot x_{i}$$
 net $i = \sum_{k} \omega_{ik} \cdot x_{k}$
 $O_{j} = tanh (net_{j})$ $O_{i} = tanh (net_{i})$

At we want to calculate and, the rate of change of error with respect to the given connective weight, so we can minimize it.

$$\frac{\partial E}{\partial \omega_{i}} = \frac{\partial E}{\partial O_{i}} \cdot \frac{\partial O_{j}}{\partial \omega_{j}} \cdot \frac{\partial net_{j}}{\partial \omega_{j}}$$

$$= (t_{p_{i}} - O_{j}) \cdot (-1) \cdot (1 - O_{j}^{2}) \cdot x_{i}$$

$$= (t_{p_{i}} - O_{j}) \cdot (-1) \cdot (1 - O_{j}^{2}) \cdot x_{i}$$

$$= + \alpha \cdot (t_{p_{j}} - O_{j}) \cdot (t_{j}) \cdot (1 - O_{j}^{2}) \cdot x_{i}$$

$$\begin{cases} \delta_{p_{j}} = -\frac{\partial E}{\partial \omega_{j}} \\ \partial \omega_{j} = \frac{\partial E}{\partial \omega_{j}} \end{cases}$$

$$\Delta \omega_{ik} = -\alpha \cdot \frac{\partial E}{\partial \omega_{ik}} = -\alpha \cdot \frac{\partial E}{\partial o_i} \cdot \frac{\partial O_i}{\partial net_i} \cdot \frac{\partial net_i}{\partial \omega_{ik}}$$

$$-\delta \rho_i$$

$$\delta p_i = \sum_{j} (+p_j - o_j).(-1). \frac{\partial O_j}{\partial net_i} = \sum_{j} (+p_j - o_j).(-1). \frac{\partial O_j}{\partial net_j}. \frac{\partial net_j}{\partial net_i}$$

$$\frac{\sum_{j} (+\rho_{3}-o_{j}).(-1). \frac{\partial o_{j}}{\partial nel_{j}}. \frac{\partial nel_{j}}{\partial o_{i}}. \frac{\partial o_{j}}{\partial nel_{j}}}{\frac{\partial o_{j}}{\partial nel_{j}}} = \left(\frac{\sum_{j} (+\rho_{3}-o_{j}).(-1).(1-o_{j}^{2}). \omega_{j}}{(1-o_{j}^{2}). \omega_{j}}\right).(1-o_{j}^{2})}$$

$$\frac{-\delta_{\rho_{j}}}{\delta_{\rho_{i}}} = (1-o_{i}^{2}). \sum_{j} -\delta_{\rho_{j}}.\omega_{j}i$$