# Q-Learning

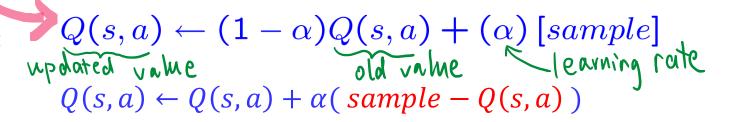- Q-Learning: sample-based Q-value iteration

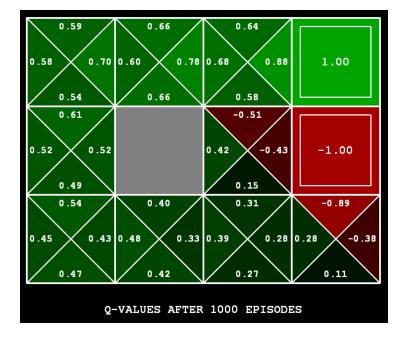$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

- 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:  ⌐next State

$$sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$

reward

Your **update** Function will implement this

- Incorporate the new estimate into a running average:

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

updated value        old value        learning rate

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



| 0.59 | 0.66 | 0.64 | |
| 0.58 | 0.70 0.60 | 0.78 0.68 | 0.88 | 1.00 |
| 0.54 | 0.66 | 0.58 | |
| 0.61 | | -0.51 | |
| 0.52 | 0.52 | 0.42 -0.43 | -1.00 |
| 0.49 | | 0.15 | |
| 0.54 | 0.40 | 0.31 | -0.89 |
| 0.45 | 0.43 0.48 | 0.33 0.39 | 0.28 0.28 | -0.38 |
| 0.47 | 0.42 | 0.27 | 0.11 |

Q-VALUES AFTER 1000 EPISODES

[Demo: Q-learning – gridworld (L10D2)]
[Demo: Q-learning – crawler (L10D3)]