

- **bold**
- underline
- *italic*

A collection of notes from the textbook, *Reinforcement Learning* by Richard Sutton and Andrew Barto. Available at <http://incompleteideas.net/book/the-book.html>

1 Action selection

Most RL methods require some form of policy or action-value based action selection algorithm.

- **Greedy Selection:** Choosing the best action.

$$A = \operatorname{argmax}_a Q(a)$$

- **ϵ -greedy Selection:** Simple exploration with ϵ -probability.

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \epsilon \text{ (breaking ties randomly)} \\ \text{a random action} & \text{with probability } \epsilon \end{cases}$$

- **Upper Confidence Bound (UCB):** Takes into account the proximity of the estimate to being maximal and the uncertainty in the estimates. Does not perform well on large state spaces.

$$A_t = \operatorname{argmax}_a \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Where:

- $c > 0$ is the degree of exploration
- $N_t(a)$ is the number of times that action a has been selected prior to time t .
If $N_t(a) = 0$, then a is considered to be a maximizing action.

2 Performance Measures