



AIL 722: Reinforcement Learning

Lecture 29: Double Q-Learning

Raunak Bhattacharyya



ScAI

YARDI SCHOOL OF ARTIFICIAL INTELLIGENCE
INDIAN INSTITUTE OF TECHNOLOGY DELHI

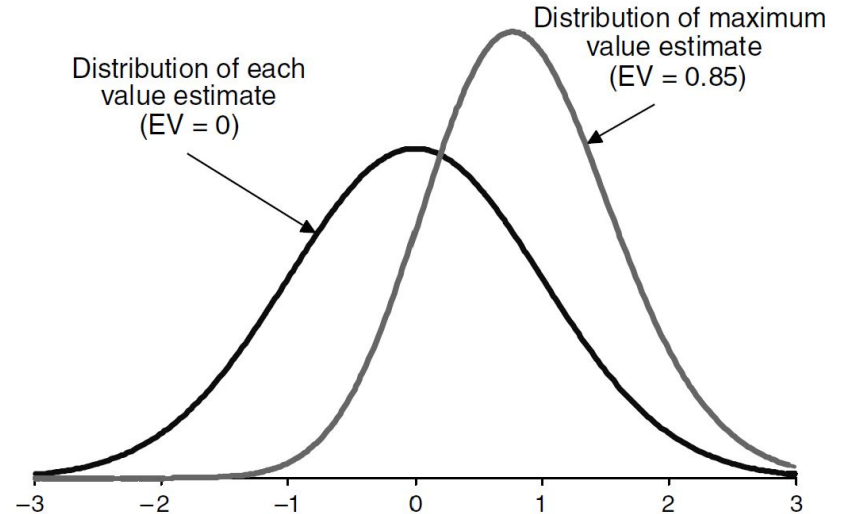
Recap & Today's Outline

- Overestimation Bias
- Roulette and ALE
- Deep Q Network case study
- Single and double estimator
- Double Q Learning
- Double DQN

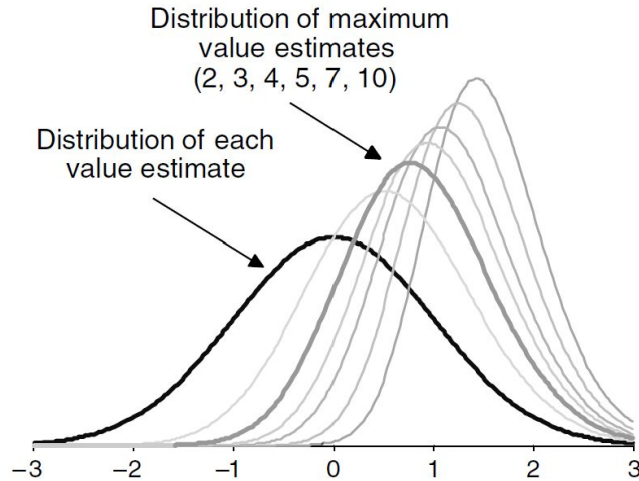
The Optimizer's Curse

- Decision science literature
- n alternatives
- Computing true values might cost many crores
- Computing estimates of the values might be lakhs of consulting effort
- Consulting firm tells us to pick the maximum expected estimated value

- 3 alternatives example



Increasing Number of Alternatives

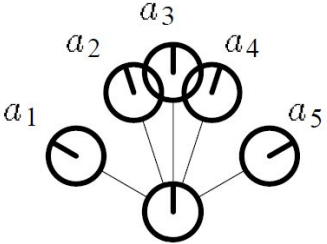
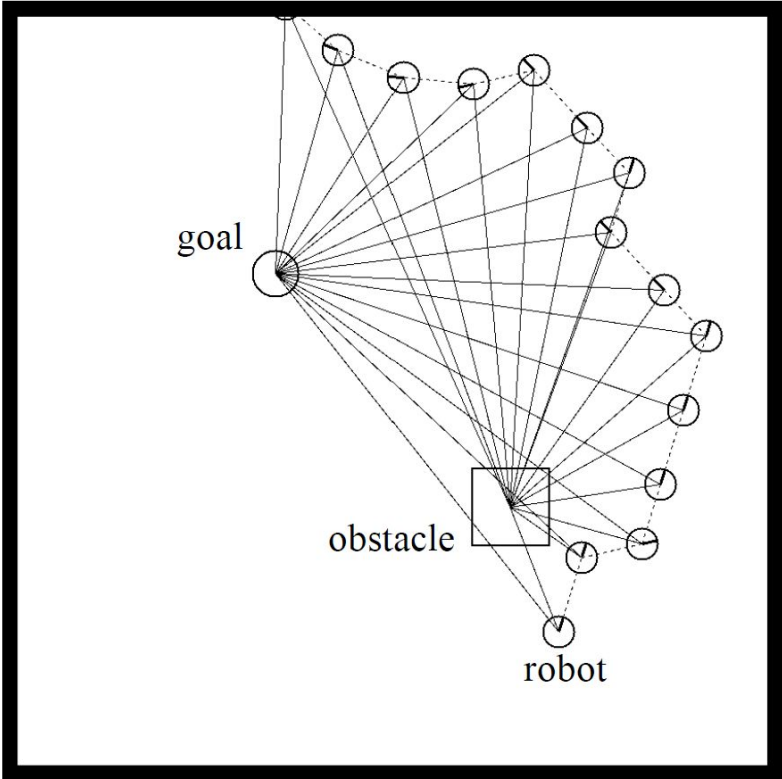


Number of alternatives	Expected disappointment
1	0.00
2	0.56
3	0.85
4	1.03
5	1.16
6	1.27
7	1.35
8	1.43
9	1.48
10	1.54

PROPOSITION 1. Let V_1, \dots, V_n be estimates of μ_1, \dots, μ_n that are conditionally unbiased in that $E[V_i | \mu_1, \dots, \mu_n] = \mu_i$ for all i . Let i^* denote the alternative with the maximal estimated value $V_{i^*} = \max\{V_1, \dots, V_n\}$. Then,

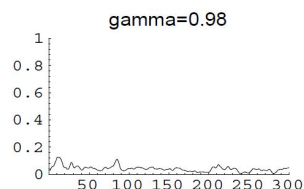
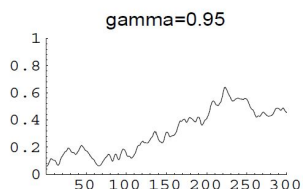
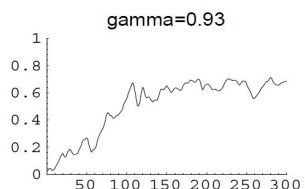
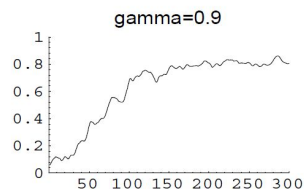
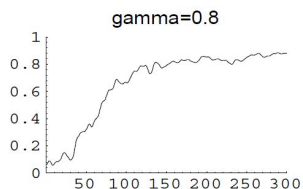
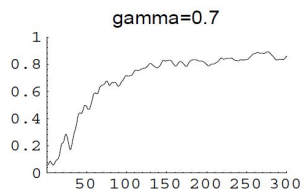
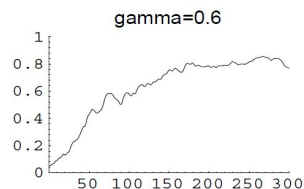
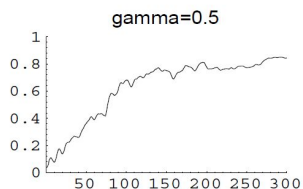
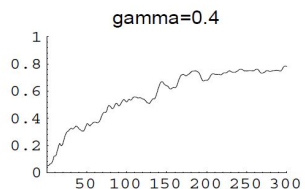
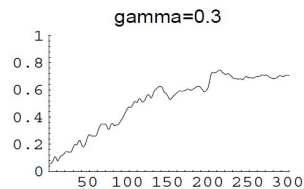
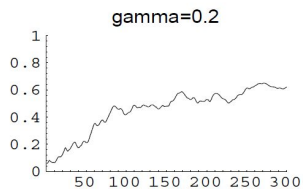
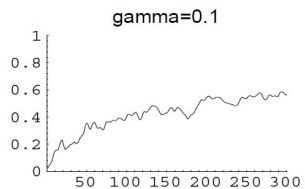
$$E[\mu_{i^*} - V_{i^*}] \leq 0. \quad (1)$$

Example from Robotics



Derives an upper bound on the overestimation

Overestimation



$$0 \leq E[Z_s] \leq \gamma c \text{ with } c = \varepsilon \frac{n-1}{n+1}$$

Towards Double Q-Learning

M random variables $X = \{X_1, \dots, X_M\}$

$$\max_i E\{X_i\}$$

Notation alert

$$S = \bigcup_{i=1}^M S_i$$

S_i contains samples of X_i

$$\mu_i(S) \stackrel{\text{def}}{=} \frac{1}{|S_i|} \sum_{s \in S_i} s$$

$$\max_i E\{X_i\} = \max_i E\{\mu_i\} \approx \max_i \mu_i(S)$$

Single estimator

Double Estimator

$$S^A \cap S^B = \emptyset$$

$$\mu_i^A(S) = \frac{1}{|S_i^A|} \sum_{s \in S_i^A} s \quad \mu_i^B(S) = \frac{1}{|S_i^B|} \sum_{s \in S_i^B} s$$

$$\mu^A = \{\mu_1^A, \dots, \mu_M^A\} \text{ and } \mu^B = \{\mu_1^B, \dots, \mu_M^B\}$$

μ_i^A and μ_i^B are unbiased

What does this mean?

Double Estimator

$$\text{Max}^A(S) \stackrel{\text{def}}{=} \{j \mid \mu_j^A(S) = \max_i \mu_i^A(S)\}$$

$$E\{\mu_j^B\} = E\{X_j\} \text{ for all } j, \text{ including all } j \in \text{Max}^A$$

$$\mu_{a^*}^A(S) \stackrel{\text{def}}{=} \max_i \mu_i^A(S)$$

$$\max_i E\{X_i\} = \max_i E\{\mu_i^B\} \approx \mu_{a^*}^B$$