

# Maximisation bias

Van Hasselt (2010, 2015)

$$Q(s, a) = Q(s, a) + \alpha \left( r + \sum \max_{a' \in A} Q(s', a') - Q(s, a) \right)$$

→ tends to overshoot (by a lot)

Abstract setting:

$X_1, \dots, X_n$

random variables

$\delta_1, \dots, \delta_n$

$\max_i E[X_i]$

single estimator:

for each  $i$ ,

Sample  $X_i$  a number of times

→ choose  $\arg \max_i \hat{E}[X_i] = *$

this is the candidate best machine

→ evaluate using  $\hat{E}[X_*]$

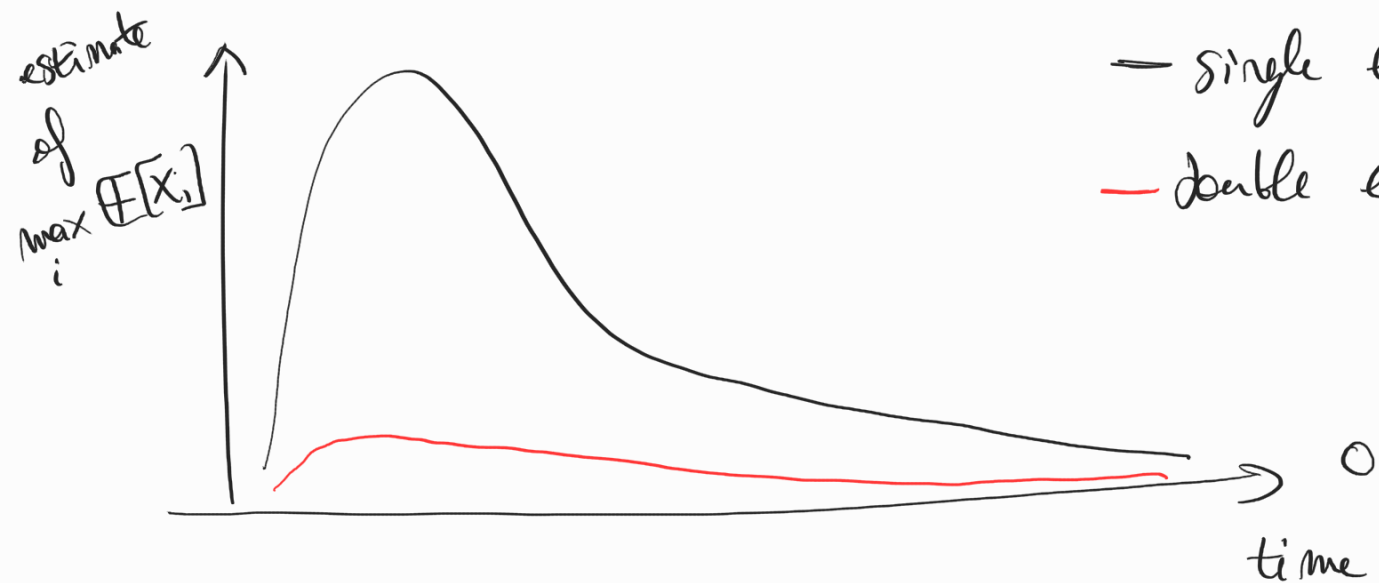
double estimator:

for each  $i$ ,

Sample  $X_i$  a number of times, twice

→ choose  $\arg \max_i \hat{E}[X_i] = *$  using the first set of samples  
this is the candidate best machine

→ evaluate using  $\hat{E}[X_*]$  using the second set of samples



Double Q-learning

$$a' = \arg \max_{a'} Q_2(s', a')$$

$$Q(s, a) = Q_1(s, a) + \alpha(r + \gamma - Q_1(s, a))$$

$$Q(s', a') = Q_2(s, a)$$

Symmetrically for updating  $Q_2$ .

Exercise: take yesterday's code and  
implement double Q-learning

REPEAT:

$S \leftarrow S_0$

while (trajectory not over):

choose action  $a$  using  $\epsilon$ -greedy strategy  
from  $q$

sample  $(S, a) : R, S'$

with probability  $1/2$ :

$$a' = \operatorname{argmax}_{a'} Q_2(s', a')$$

$$Q_1(s, a) = Q_1(s, a) + \alpha (R + \gamma$$

$$Q_1(s', a') - Q_1(s, a)$$

$Q_1(s, a)$

$$a' = \arg \max_{a'} Q_1(s', a')$$

$$Q_2(s, a) = Q_2(s, a) + \alpha (r + \gamma$$

$$Q_2(s', a') - Q_2(s, a)$$