# Maximization Bias
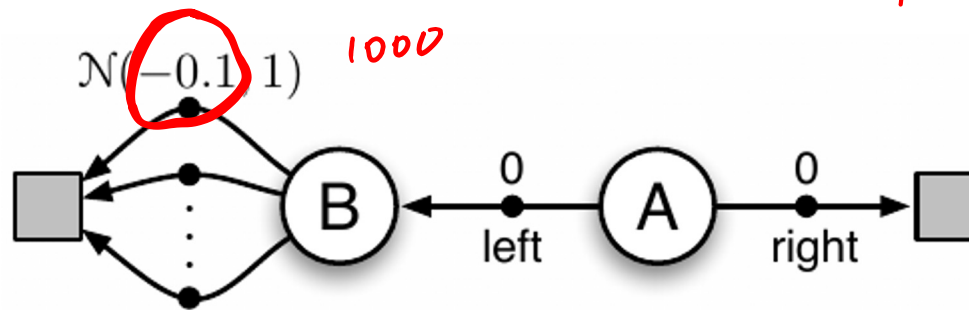
B. Ravindran

# Maximization Bias
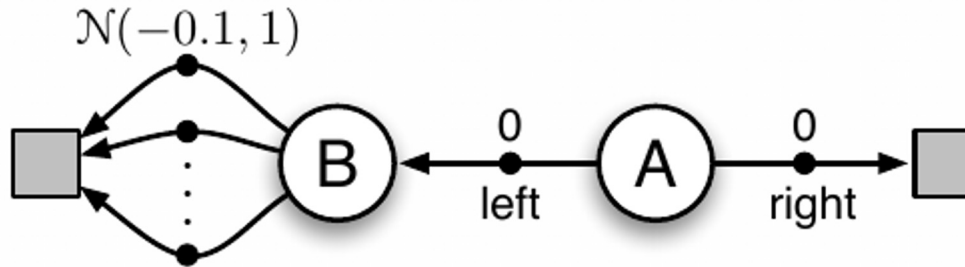
Consider the simple example below:



1. A is the starting state.
2. T(A, left, B) = 1
3. R(A, left) = 0, R(A, right) = 0
4. From B, there are |N| actions available, each of which results in a terminal state. And these |N| actions are normally distributed with mean = -0.1 and std = 1

# Maximization Bias

Consider the simple example below:



1. A is the starting state.
2. T(A, left, B) = 1
3. R(A, left) = 0, R(A, right) = 0
4. From B, there are |N| actions available, each of which results in a terminal state. And these |N| actions are normally distributed with mean = -0.1 and std = 1
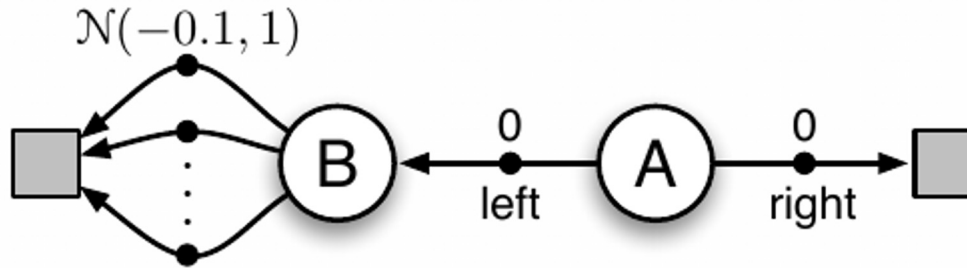
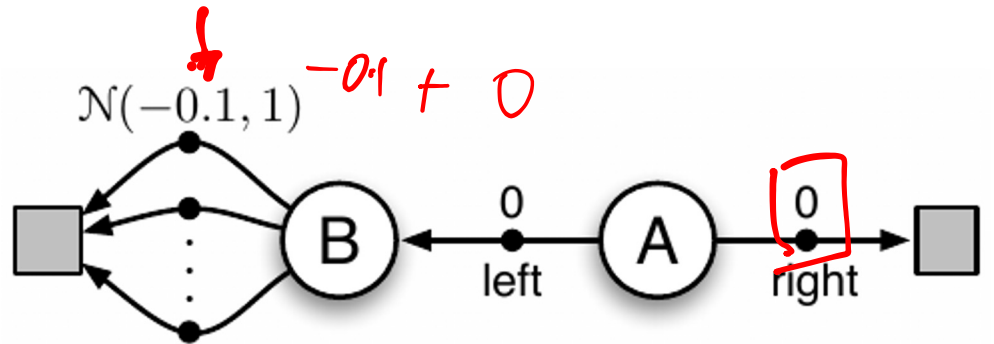# Maximization Bias

Which direction to move from A?

# Maximization Bias

Which direction to move from A?
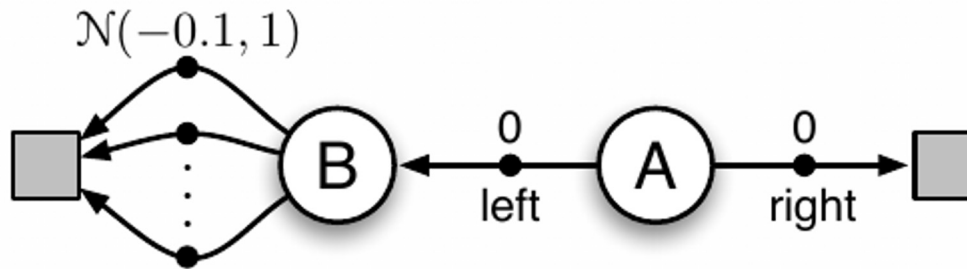


$\mathcal{N}(-0.1, 1)$ -0.1 + 0

$E[\ G_t\ |s_0 = A, a_0 = \text{left}\ ] = -0.1$

$E[\ G_t\ |\ s_0 = A, a_0 = \text{right}\ ] = 0$

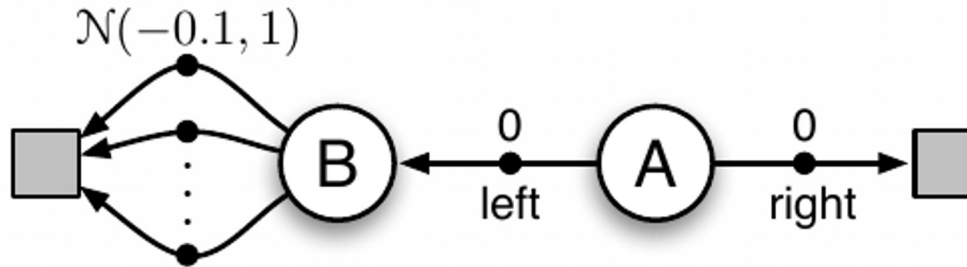# Maximization Bias

Which direction to move from A?



$$\mathcal{N}(-0.1, 1)$$

$$E[\ G_t\ |\ s_0 = A,\ a_0 = \text{left}\ ] = -0.1$$

$$E[\ G_t\ |\ s_0 = A,\ a_0 = \text{right}\ ] = 0$$

# Maximization Bias

$Q(B, a) \doteq 0$
$\forall a$

What happens when we learn a policy using Q-learning?
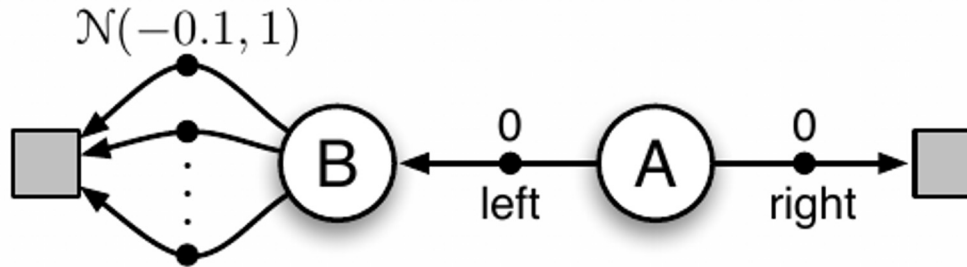


$Q(A, \text{Left})$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

$$Q(A, \text{left}) + \alpha \left[ 0 + \gamma \max_a Q(B, a) - Q(A, \text{left}) \right]$$

$$i = 1, \ldots, 1000 \quad Q(B, a_i) \leftarrow Q(B, a_i) + \alpha \left[ \gamma + \cdots - Q(B, a_i) \right]$$

# Maximization Bias

What happens when we learn a policy using Q-learning?



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

Using maximum over estimate
as an estimate of the maximum!

Leads to a positive bias, called **maximization bias**.

# Maximization Bias



% left actions from A

$\mathcal{N}(-0.1, 1)$

Q-learning

optimal

Episodes

ε = 0.1 for above example
Therefore, 10% of the actions are random.
Optimal => 5% can be right (random) and 95% should be left.

9

# Maximization Bias
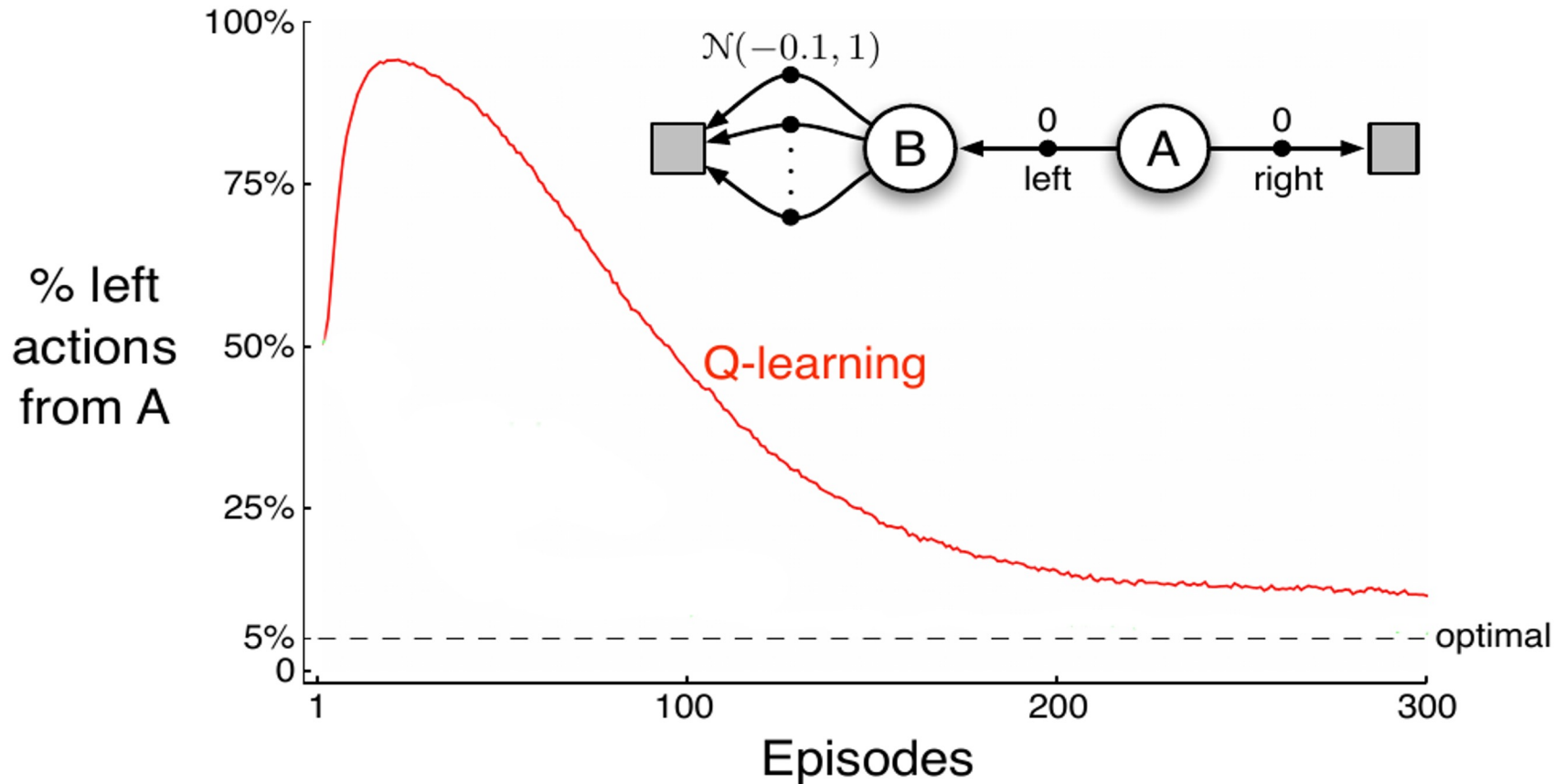


ε = 0.1 for above example
Therefore, 10% of the actions are random.
Optimal => 5% can be right (random) and 95% should be left.

# Double Q-learning

The problem can also be viewed as:

using the same samples both to determine the maximizing action and to estimate its value

# Double Q-learning

The problem can also be viewed as:

using the same samples both to determine the maximizing action and to estimate its value

Solution: Use different estimates for maximizing the action and estimating its value

$$\hat{a} = \operatorname*{argmax}_{a} Q_1(s, a)$$

$$Q_2(s, \hat{a})$$

$$Q_1$$

$$r_1 \quad r_2 \quad r_3 \quad r_4 \quad \text{----} \quad r_{100}$$

$$Q_2$$

$$Q_1 \text{ ave} ( r_1 + r_3 + \cdots )$$

$$Q_2 \text{ ave} ( r_2 + r_4 + \cdots )$$

12

# Double Q-learning

**Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, such that $Q(terminal, \cdot) = 0$

Loop for each episode:

    Initialize $S$

    Loop for each step of episode:

        Choose $A$ from $S$ using the policy $\varepsilon$-greedy in $Q_1 + Q_2$

        Take action $A$, observe $R, S'$

        With 0.5 probabillity:

            $Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \Big( R + \gamma Q_2 \big( S', \arg\max_a Q_1(S', a) \big) - Q_1(S, A) \Big)$

        else:

            $Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \Big( R + \gamma Q_1 \big( S', \arg\max_a Q_2(S', a) \big) - Q_2(S, A) \Big)$

        $S \leftarrow S'$

    until $S$ is terminal

*(handwritten annotation: — TD target)*

# Double Q-learning

Double Q-learning vs Q-learning (ε = 0.1)