

# Homework 1

10-403 Recitation 3

# Homework 1

## Clarifications

- Boltzman exploration:

- $\mathbb{P}(A_t = a) = \pi_t(a) = \frac{\exp(Q_t(a)/T)}{\sum_{a'} \exp(Q_t(a')/T)}$

- As  $T \rightarrow \infty$  , behaviour converges to uniformly at random policy

- As  $T \rightarrow 0$  , behaviour converges to pure greedy policy

- $\epsilon$ -greedy exploration:

- w.p.  $1 - \epsilon$  :  $A_t = \arg \max_a Q_t(a)$

- w.p.  $\epsilon$  :  $A_t \sim \text{Uniform}(\{a_1, \dots, a_N\})$

# Homework 1

## Clarifications

- Boltzman exploration:

- $\mathbb{P}(A_t = a) = \pi_t(a) = \frac{\exp(Q_t(a)/T)}{\sum_{a'} \exp(Q_t(a')/T)}$

- As  $T \rightarrow \infty$ , behaviour converges to uniformly at random policy

- As  $T \rightarrow 0$ , behaviour converges to pure greedy policy

- $\epsilon$ -greedy exploration:

- w.p.  $1 - \epsilon$  :  $A_t = \arg \max_a Q_t(a)$

- w.p.  $\epsilon$  :  $A_t \sim \text{Uniform}(\{a_1, \dots, a_N\} - \{\arg \max_a Q_t(a)\})$

Both forms of  $\epsilon$ -greedy are acceptable  
But previous slide is easier to implement

# Homework 1

## Clarifications

- “Fraction of population to keep at each iteration: 10%”
  - $\implies$  Elite Size = 10

# Homework 1

## Reminder

- Please respond to Team Information poll on Piazza!