# Reinforcement Learning II

## ★ How to Explore?

⟹ Several Schemes for forcing exploration

# Simplest : random action (==E-greedy==)

> → Every time step, flip a coin
> → With (small) probability $\varepsilon$, act randomly
> → With (large) probability $1-\varepsilon$, act on current policy

⟹ Problem with random actions?

> → You do eventually explore the space, but keep thrashing around once learning is done.
> → One Solution : lower $\varepsilon$ over time
> → Another Solution : Exploration functions

## ★ Exploration Function

▪ Random actions : Explore a fixed amount
▪ Better idea : Explore areas whose badness is not (Yet) established, eventually stop exploring

⟹ Takes a value estimate $u$ and a visit count $n$ & return an optimistic utility :

$$f(u,n) = u + \frac{k}{n}$$

Modified Q-update :
$$Q(s,a) \xleftarrow{\alpha} R(s,a,s') + \gamma \max_{a'} f\left(Q(s',a'), N(s',a')\right)$$

# * Regret

$\Rightarrow$ Measure of your (total mistake cost)

Difference between you rewards

        &              optimal rewards

$$regret = u(\text{Optimal action}) - u(\text{action taken})$$

# * Approximate Q-Learning

⇒ Basic Q-Learning keeps a table of all q-values.

⇒ In realistic situation's we cannot possibly learn about every single state.

       ↳ Too many states to visit them all in training.

       ↳ Too many states to hold the q-tables in memory.

→ Instead, we want to generalize:

       ↳ Learn about some small number of training states from experience.

       ↳ Generalize that experience to new, similar situations.

       ↳ This is a fundamental idea in machine learning, and we will see it over & over again.

## (#) Feature-Based Representations

"Describe state using a vector of features"

⇒ Features are functions from states to real numbers (often 0/1) that capture important properties of the state.

     → Examples:

       ↳ Distance to closest ghost

       ↳ Distance to closest dot

       ↳ Number of ghosts

       ↳ $1/(\text{dist to dot})^2$

       ↳ Is pacman in a tunnel? (0/1)

# Linear Value Function

⇒ Using feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = \omega_1 f_1(s) + \omega_2 f_2(s) + \cdots + \omega_m f_m(s)$$

$$Q(s) = \omega_1 f_1(s,a) + \omega_2 f_2(s,a) + \cdots + \omega_m f_m(s,a)$$

Advantage: Our experience is summed up in a few powerful numbers.

Disadvantage: States may share features but actually be very different in value.

⇒ Q-learning with linear Q functions:

transition $= (s, a, \mathcal{r}, s')$

difference $= \left[\mathcal{r} + \gamma \max\limits_{a'} Q(s', a)\right] - Q(s,a)$

Exact Q: $Q(s,a) \leftarrow Q(s,a) + \alpha \, [\text{difference}]$

Approxite Q: $\omega_i \leftarrow \omega_i + \alpha \, [\text{difference}] \, f_i(s,a)$

# ⊕ Q-Learning & Least Square

Let $y = \mathcal{r} + \gamma \max\limits_{a'} Q(s', a))$

error $= y - Q(s,a)$

$\searrow \omega_1 f_1 + \omega_2 f_2 \cdots \omega_m f_m$

$$\nabla_\omega \tfrac{1}{2} e^2 = \nabla_\omega (y - \omega^T f)^2 = -\tfrac{2}{2}(y - \omega^T f) f$$

$$\omega \leftarrow \omega + 2\alpha \, [\text{diff}] \, f(s,a)$$

# * Policy Search

⇒ Learn policies that maximize rewards, not the value that predict them.

⇒ Simplest policy search:
- → Start with an initial value function or Q function
- → Nudge each feature weight up & down & see if your policy is better than before

———————— ✕ ———————— ✕ ————————