

Basic Knowledge:

Core Components:

State. Action. Award

unchanged once set

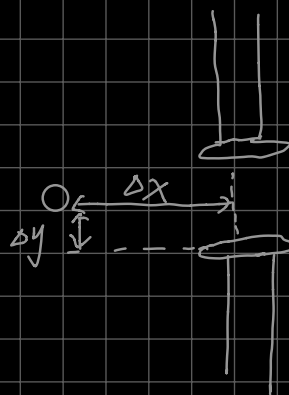
Q-learning:

Q : action-utility function, 用于评价在状态 x 下采取动作 y

s 下采取动作 a 的优劣。 \Rightarrow 相当于 agent 的记忆, 随着训练过程而改变。 (即 Q-table)

Entity: Flappy Bird

State $\Rightarrow (\Delta x, \Delta y)$



Action \Rightarrow fly. fall

Award \Rightarrow alive through 1 tube: +100, dead: -1000

Q /action-utility function \Rightarrow $m \times n$ table, each grid in the table represents utility value of an action

fly/fall upon a state $(\Delta x_i, \Delta y_j)$.

State	fly	fall
$(\Delta x_1, \Delta y_1)$	+1	+20
$(\Delta x_1, \Delta y_2)$	-100	2
\vdots		
$(\Delta x_m, \Delta y_n)$	50	-200

$\Rightarrow Q$

Train goal: To get a perfect table Q , which makes the bird get maximum utility on current state and will be alive forever. How to Train \Rightarrow

def train():

Init $Q = \{\}$:

while Q is still Diverged:

Init state $S \leftarrow S_0$

while $S \neq \text{Dead}$:

Get core components

Get action a through strategy $\pi(S_0)$ according to Q . choose max utility.

Get award $R(S, a)$ and new state S' .

Update

update $Q[S, A]$

$$Q[S, A] \leftarrow (1 - \alpha) * Q[S, A] + \alpha * (R(S, a) + \gamma * \max_a Q[S', a])$$

past training learning rate Current utility memorial utility discount factor

update State

$S = S'$

△ α 越大, 过去训练保留的效果越少

γ 越大, agent 越重视过去的经验, 而忽视眼前的利益.

△ 缺点: ① 适用的 state space. action space 非常小.

② 没有预测能力/泛化能力 \Rightarrow 当一个 state 未出现过, Q learning 无法处理

How to solve \Rightarrow

Solve ①:

Policy Based method.

Solve ②:

DQN (Deep Q-learning Network)