*Learning rate = 0.1*

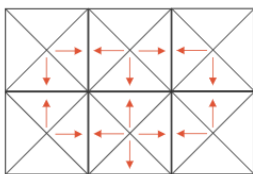*Epsilon = 1 (this decays each timestep i.e. reduces slowly towards the value 0)*

*Discount Factor = 0.99*

Environment

| S (*) | 1 | 0 |
|-------|-----|-----|
| 0 | -10 | 10 |

Initialise Q-Table (Q-Values) && 4 unique actions (up, right, down, left)

|  | UP | RIGHT | DOWN | LEFT |
|------|------|-------|------|------|
| S | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 |
| -10 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 |

For now, **our Q-table is useless**; we need **to train our Q-function using the Q-Learning algorithm**

## Training Timestep 1:

Choose an action using the Epsilon Greedy Strategy (Epsilon = 1, it will be decayed with time i.e. reduced)  Because epsilon is big (= 1.0), We can take a random action.

|  | * | 0 |
|-----|-----|-----|
| 0 | -10 | 10 |

Moving to the right gets a reward of 1

We then need to update our q-value for this State - Action pair i.e. moving to the right while at the starting state.

|  | UP | RIGHT | DOWN | LEFT |
|---|-----|-------|------|------|
| S | 0 | 0 | 0 | 0 |

To make this update, we use the Q-learning formula

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

| New Q-value estimation | Former Q-value estimation | Learning Rate | Immediate Reward | Discounted Estimate optimal Q-value of next state | Former Q-value estimation |
|---|---|---|---|---|---|

TD Target

TD Error

Therefore:

- Q(S, Right) is given by: **0 + 0.1 * [1 + 0.99 * 0 - 0] = 0.1**

|  | UP | RIGHT | DOWN | LEFT |
|---|---|---|---|---|
| **S** | 0 | 0.1 | 0 | 0 |

**Eq. from above**

*Learning rate*

*discount factor*

*the former Q value*

**0    + 0.1    * [1    + 0.99    *    0  -    0**

*the immediate reward*

*the former Q value*

*the current max value for this sate.
since all values were init. to zeros,
the current max == 0*

We first decay the epsilon slightly i.e. from 1 to 0.99

Because epsilon is still high (= .99), We can take another random action. e.g moving **down**

| | | |
|---|---|---|
| | | 0 |
| 0 | * | 10 |

Moving to the down gets a reward of -10

We then need to update our q-value for this State - Action pair i.e. moving to the right while at the starting state.

| | UP | RIGHT | DOWN | LEFT |
|---|---|---|---|---|
| S | 0 | 0.1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

Therefore:

- Q(1, Down) is given by: **0 + 0.1 * [-10 + 0.99 * 0 - 0]** = **-1**

| | UP | RIGHT | DOWN | LEFT |
|---|---|---|---|---|
| S | 0 | 0.1 | 0 | 0 |
| 1 | 0 | 0 | -1 | 0 |