

50.039 Theory and Practice of Deep Learning

W12-S1 Introduction to Reinforcement Learning

Matthieu De Mari, Berrak Sisman



About this week (Week 12)

1. What is **Reinforcement Learning**?
2. What are the key ideas behind **reinforcement learning** and its **framework**?
3. What is the **exploration vs. exploitation tradeoff**?
4. How do we **train** an **RL agent** by exploring, then progressively exploiting?
5. What are some **advanced strategies** in **multi-arm bandit problems**?
6. What are the **Q** and **V functions** for a RL problem?
7. What is **Q-learning** and how can it be implemented in RL problems?

Reinforcement Learning, a definition

Definition (Reinforcement Learning - RL):

Reinforcement Learning (RL) is one of the many machine learning paradigms.

Where **Deep Learning** attempts to learn from **examples** (i.e. datasets), **RL** attempts to train an AI from **trial-and-error**.

There is no dataset, only a feedback on the action taken.

This paradigm challenges the idea that computers and AIs should be learning from datasets.

- What if the task at hand cannot be solved easily and no dataset can even be generated?
- What if the dataset is flawed?

Reinforcement Learning, a definition

- A Go AI would need to look at a board and produce the best move as an output.
- However, Go is a prime example of a problem for which no dataset (board state \rightarrow best action) can be generated.
- This has to do with the fact our understanding of the game of Go is pretty bad.



Reinforcement Learning, a definition

Definition (Reinforcement Learning - RL):

Reinforcement Learning (RL) is one of the many machine learning paradigms.

Where **Deep Learning** attempts to learn from **examples** (i.e. datasets), **RL** attempts to learn from **trial-and-error**.

There is no dataset, only a feedback on the action taken.

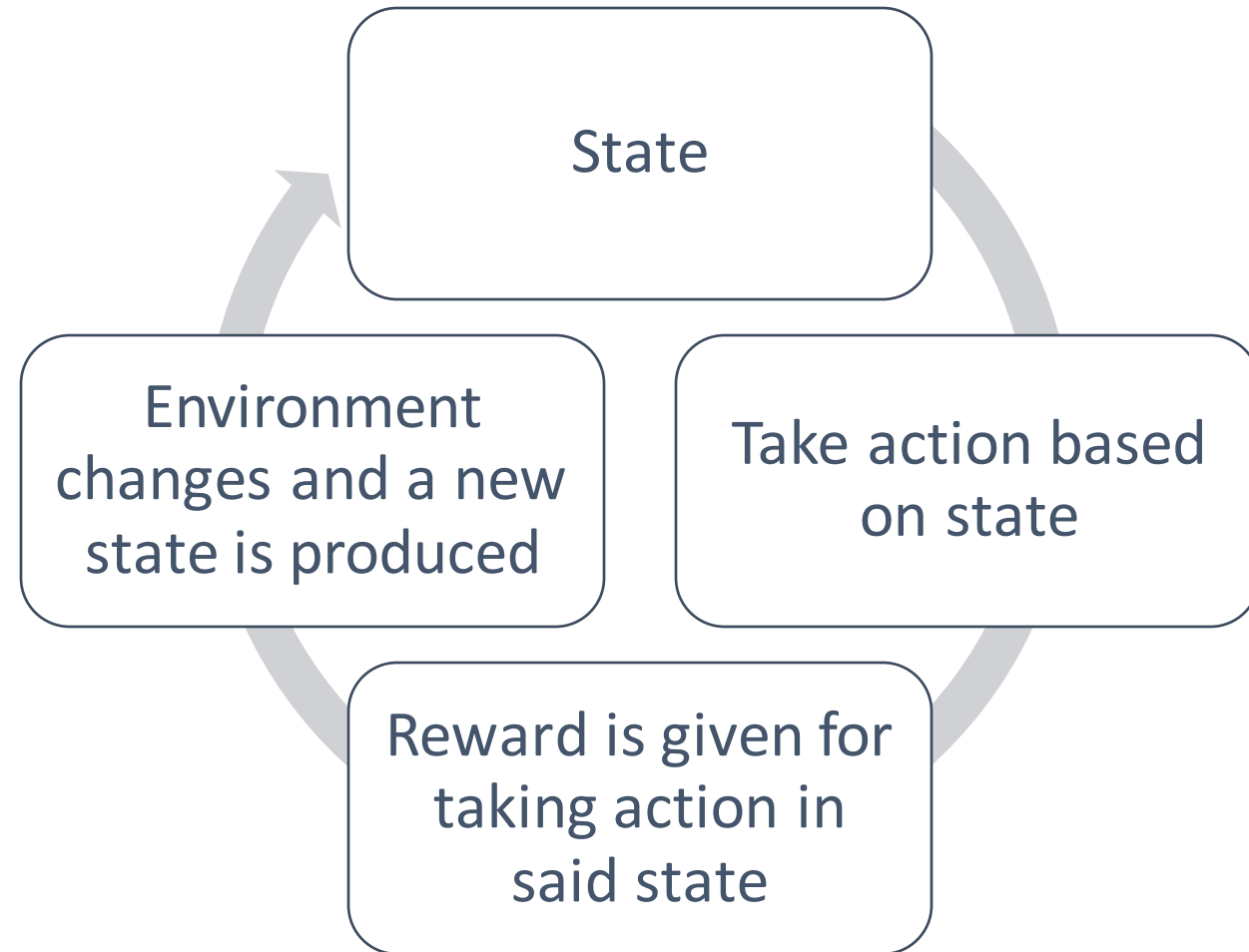
In RL, the AI learns by trying certain **actions** in given **states** and then receives a **feedback**, not always instantaneous, of whether or not the action turned out to be a good or bad decision.

The AI then learns, on-the-fly, from trying out combination of (**states**, **actions**, **feedback**) for a given problem.

A typical RL problem

A typical RL problem can be broken down in three parts.

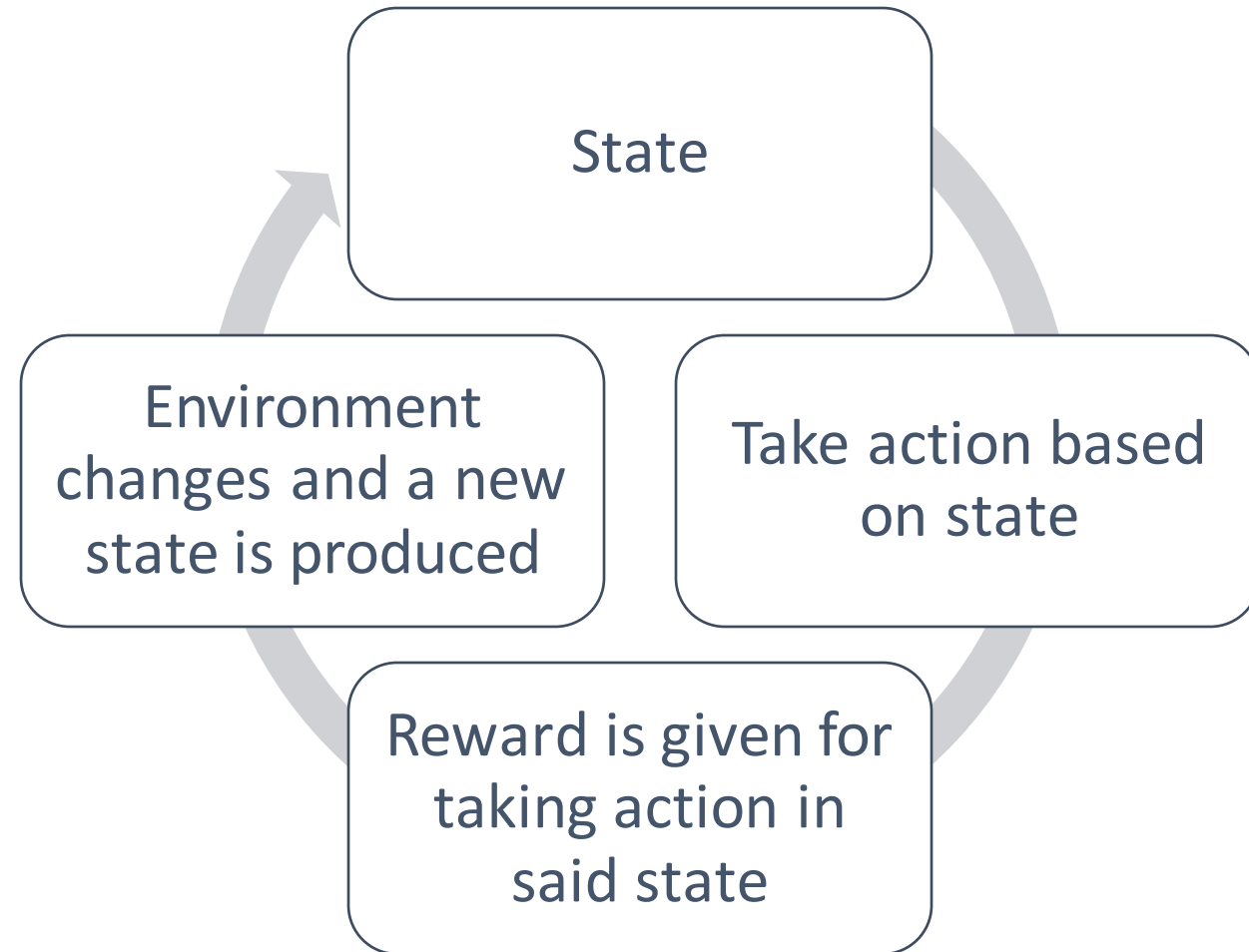
- The AI is given a problem, at time t , in a current **state** $s(t)$.
- The AI then takes an **action**, to answer the problem, in this given state, $a(t)$.
- The **action** has an effect, positive or negative, which is eventually measured by the AI, in terms of **feedback**, $R(t, a(t), s(t))$.



A typical RL problem

A typical RL problem can be broken down in three parts.

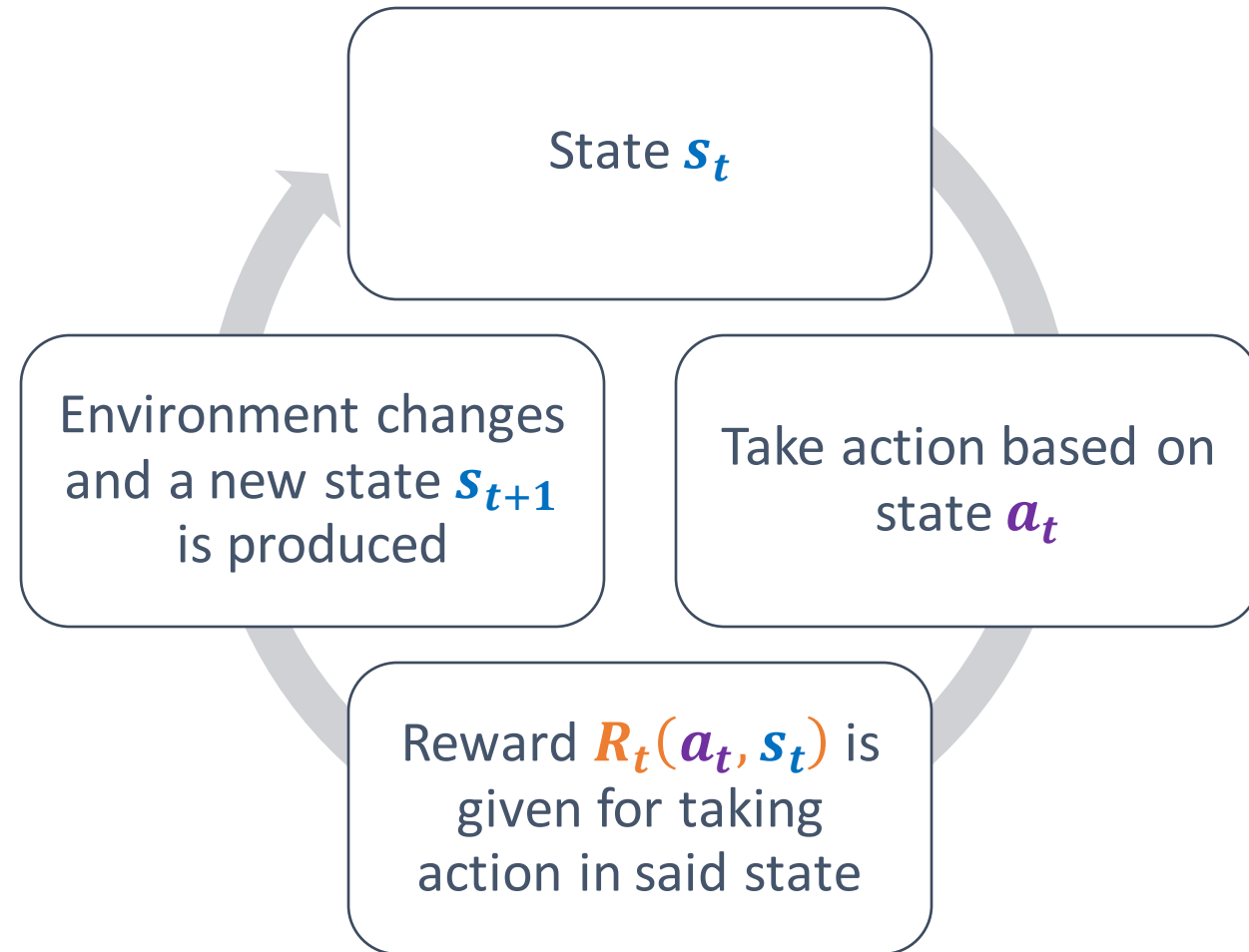
- The AI is given a problem, at time t , in a current **state** s_t .
- The AI then takes an **action**, to answer the problem, in this given state, a_t .
- The **action** has an effect, positive or negative, which is eventually measured by the AI, in terms of **feedback**, $R_t(a_t, s_t)$.



A typical RL problem

A typical RL problem can be broken down in three parts.

- The AI is given a problem, at time t , in a current **state s_t** .
- The AI then takes an **action**, to answer the problem, in this given state, a_t .
- The **action** has an effect, positive or negative, which is eventually measured by the AI, in terms of reward, $R_t(a_t, s_t)$.

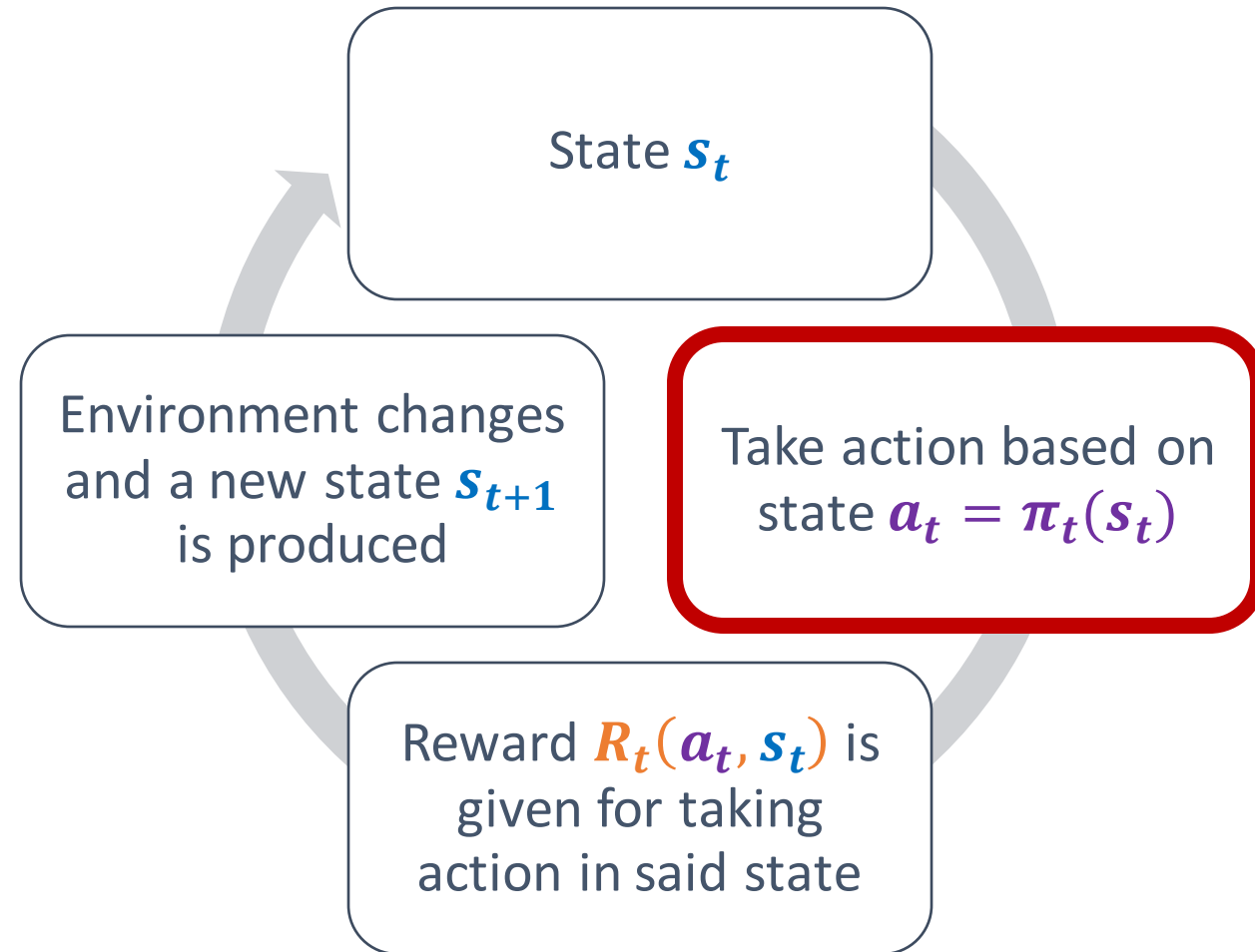


A typical RL problem

Definition (**policy**):

We call **policy**, the function, which returns the **action** a_t , among the set of possible actions A , to be taken in any given **state** s_t . It is often denoted as $a_t = \pi_t(s_t)$.

When training an AI, with RL, the objective is define the **best policy**, so that the AI knows the best course of action $\pi_t(s_t)$ to use in response for any possible **state** s_t .

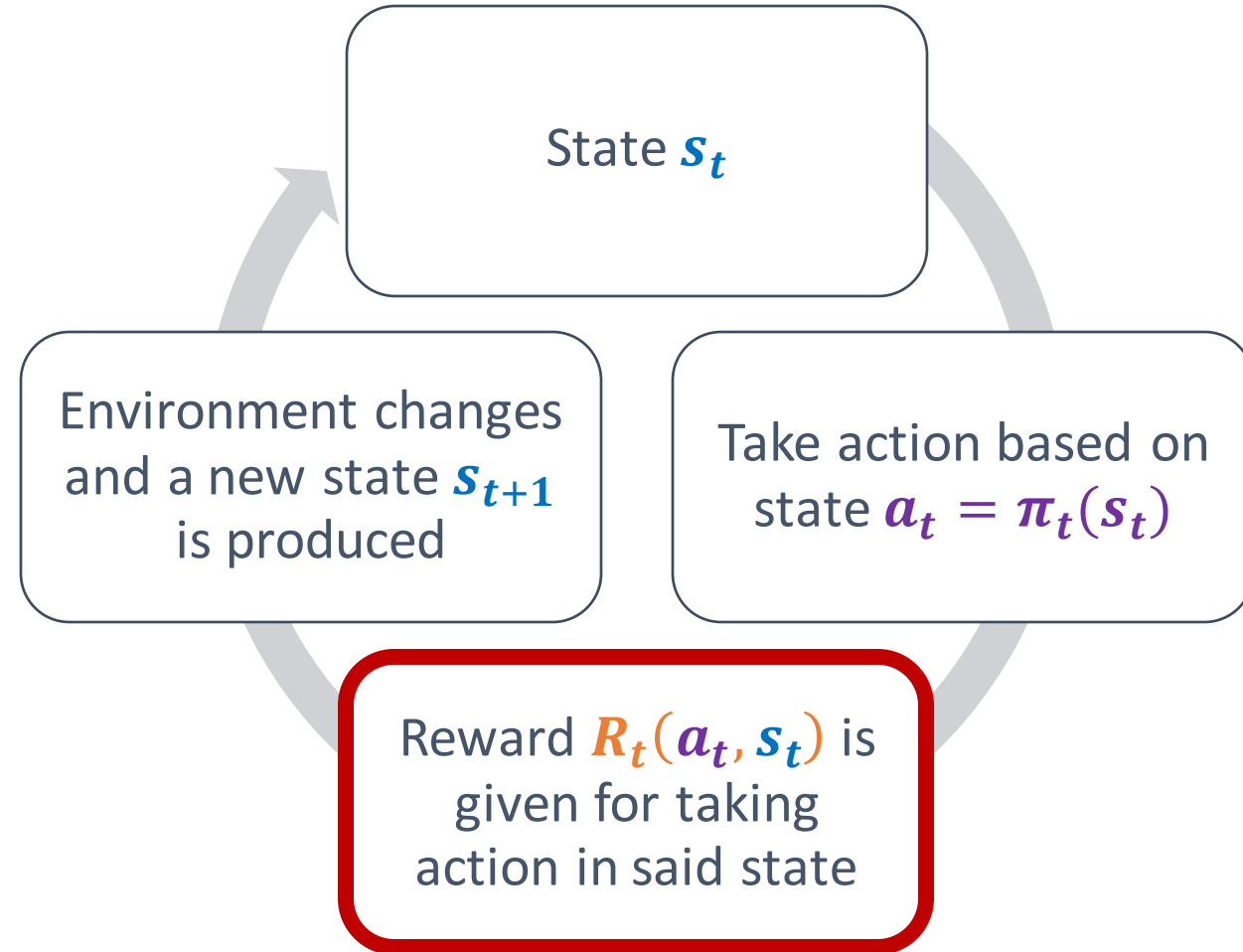


A typical RL problem

Definition (reward and return):

The **reward** function $R_t(a_t, s_t)$ is the “loss” function given to the AI.

The objective of the AI is then to find the best policy $\pi_t(s_t)$...



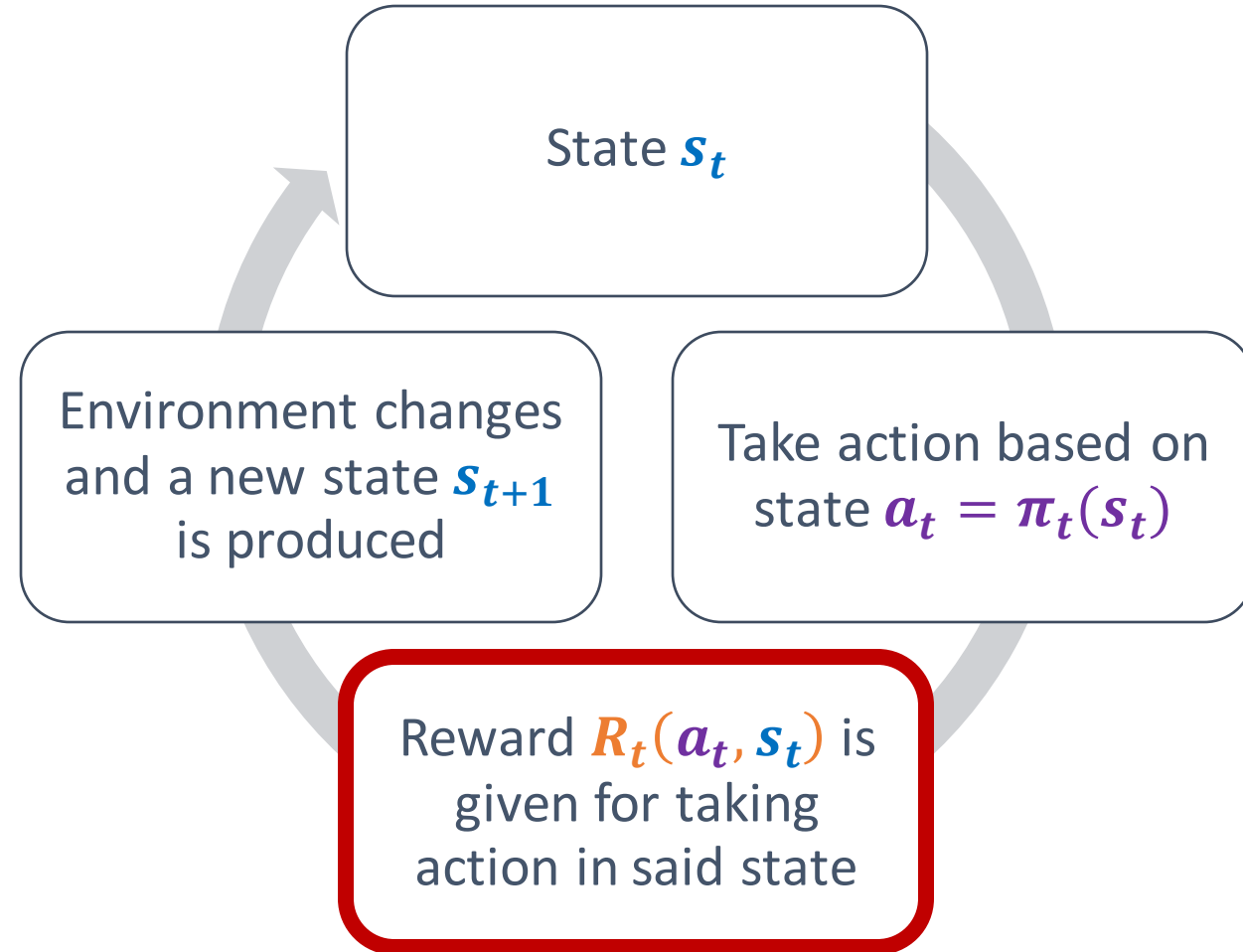
A typical RL problem

Definition (reward and return):

The **reward** function $R_t(a_t, s_t)$ is the “loss” function given to the AI.

The objective of the AI is then to find the best policy $\pi_t(s_t)$, that is, the policy, which maximizes the **return** G_t , i.e. the gains that can be expected in the future, if action $\pi_t(s_t)$ is taken in present state s_t .

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$$



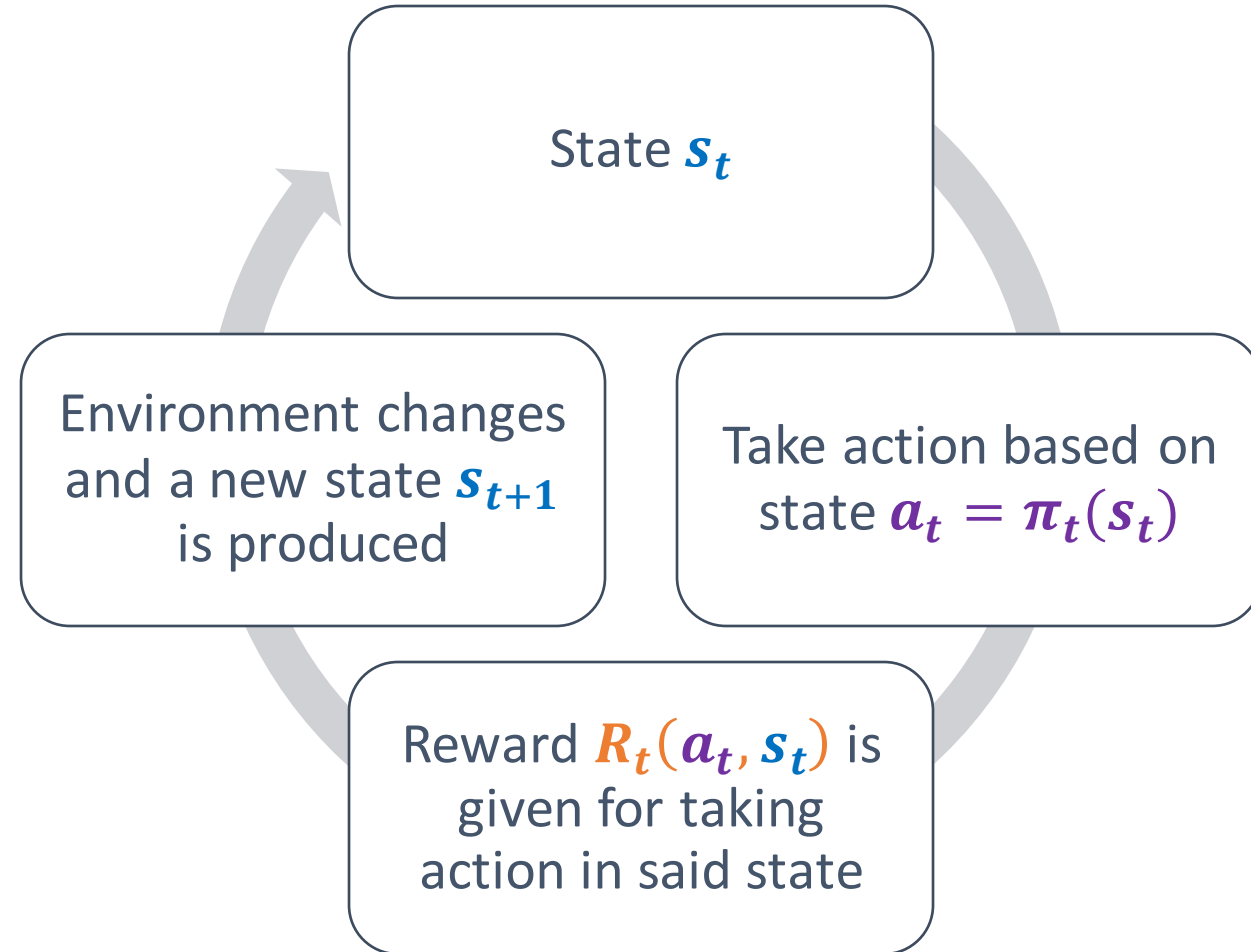
A typical RL problem

Definition (reward hypothesis):

The RL framework relies on the **reward hypothesis**.

This states that “Any goal can be formalized as the outcome of maximizing a cumulative reward (or gain) function G_t of some sort”.

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$$



A first toy problem

The Random Candy Machines problem

Candy Machines are **deterministic**

A standard candy machine often work as follows.

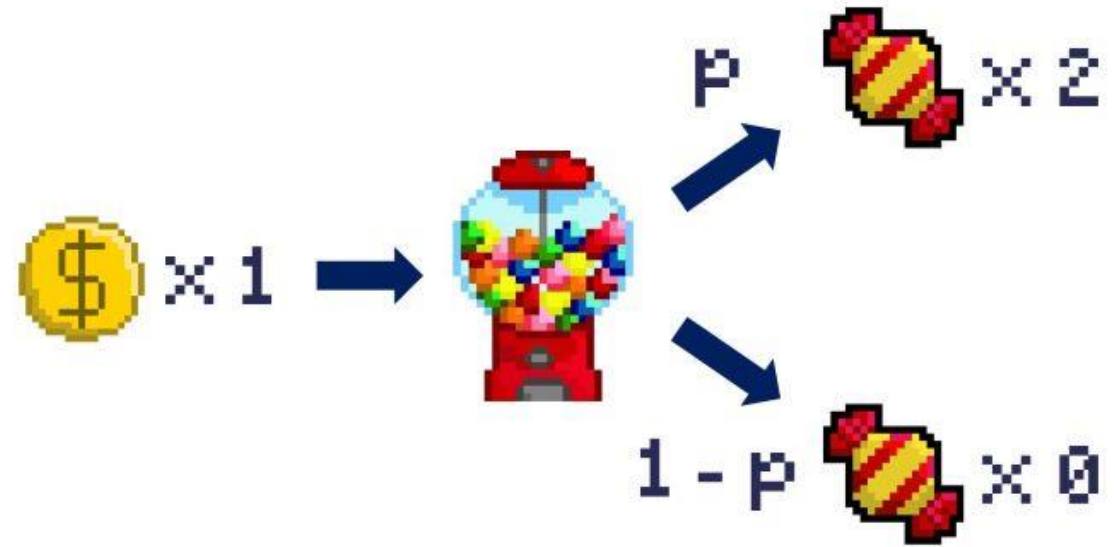
- You put one coin in the machine and obtain a candy as a result.
- Its behavior is **deterministic**, i.e. it is 100% predictable: it will ALWAYS return a candy, whenever you put a coin in the machine.



What if machines were **random**?

Let us now consider a world where candy machines are **random machines** instead.

- When you put a coin in, it might give you two candies... Or no candies at all.
- Depending on a probability p .
- For instance if $p = 80\%$, you will receive
 - two candies 80% of the time,
 - and no candies 20% of the time.



A simple problem?

 × 10000

#1



#2



#3



#4



#5





A simple problem?

Let us consider the following problem.

- You have **$N = 10000$ coins**.
- You are facing **5 machines**: machine #1 is a **deterministic** machine, machines #2-5 are **random** machines.
- Each **random** machine #2-5 has its own **winning probability** (p_2, p_3, p_4, p_5) and these will not change over time.
- **You do NOT know** what the probabilities (p_2, p_3, p_4, p_5) are, but you can use some of your coins to test the machines, and estimate the hidden probabilities (p_2, p_3, p_4, p_5) .
- **Big question**: what should be your **strategy** to **maximize** the number of candies you obtain with your 10000 coins?



Reusing the RL formalism

- Let us denote t , the **index of the coin to be played**.

$$t \in \{1, \dots, 10000\}$$

- The **set of actions** is simply

$$A = \{1, 2, 3, 4, 5\}$$

- We could define our **state**, at time t , as our **current estimates of the expectations** $e_i(t)$ of each machine $i \in A$.

$$s_t = (e_1, e_2, e_3, e_4, e_5)$$



Reusing the RL formalism

Ultimately, we get the feeling that the best strategy will follow this type of scheme...

- We will probably **try our machines to guess/estimate which one is the best.**

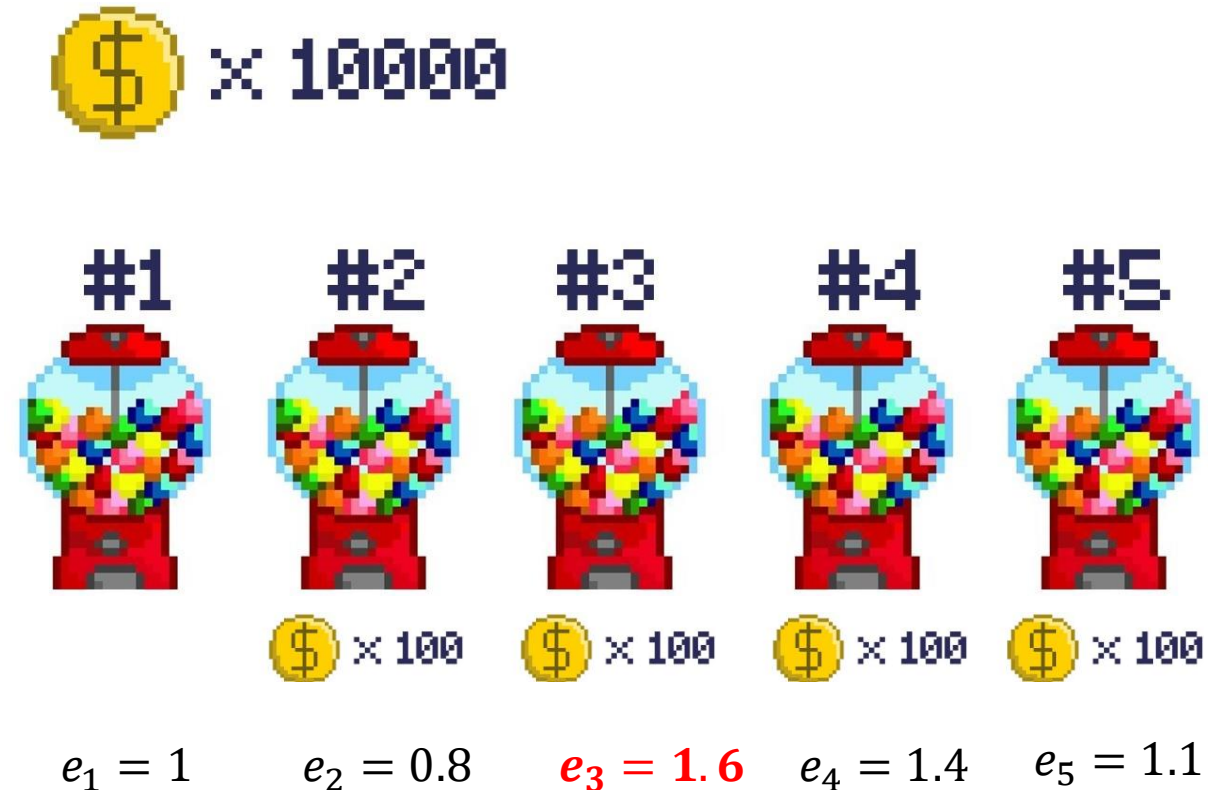
 $\times 10000$



Reusing the RL formalism

Ultimately, we get the feeling that the best strategy will follow this type of scheme...

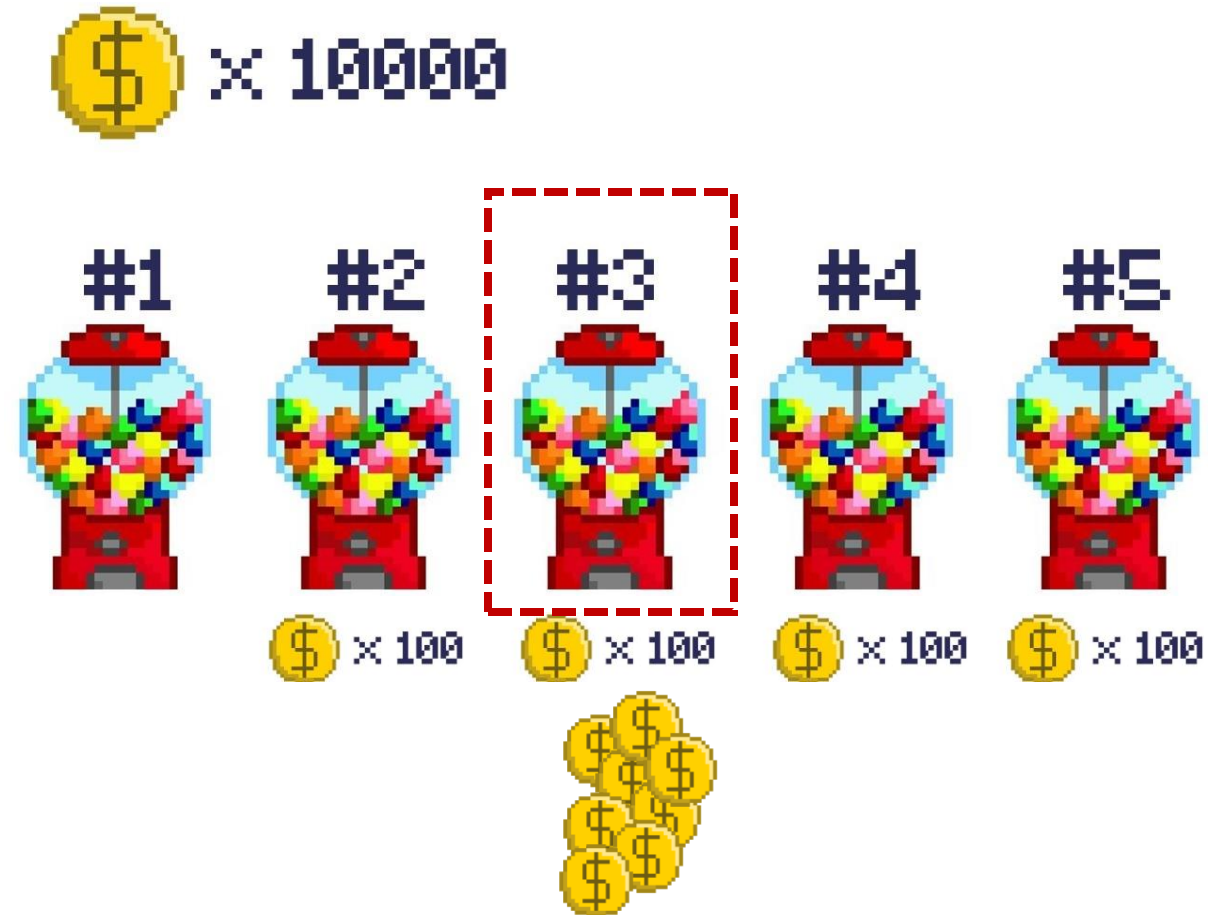
- We will probably **try our machines to guess/estimate** which one is the best.



Reusing the RL formalism

Ultimately, we get the feeling that the best strategy will follow this type of scheme...

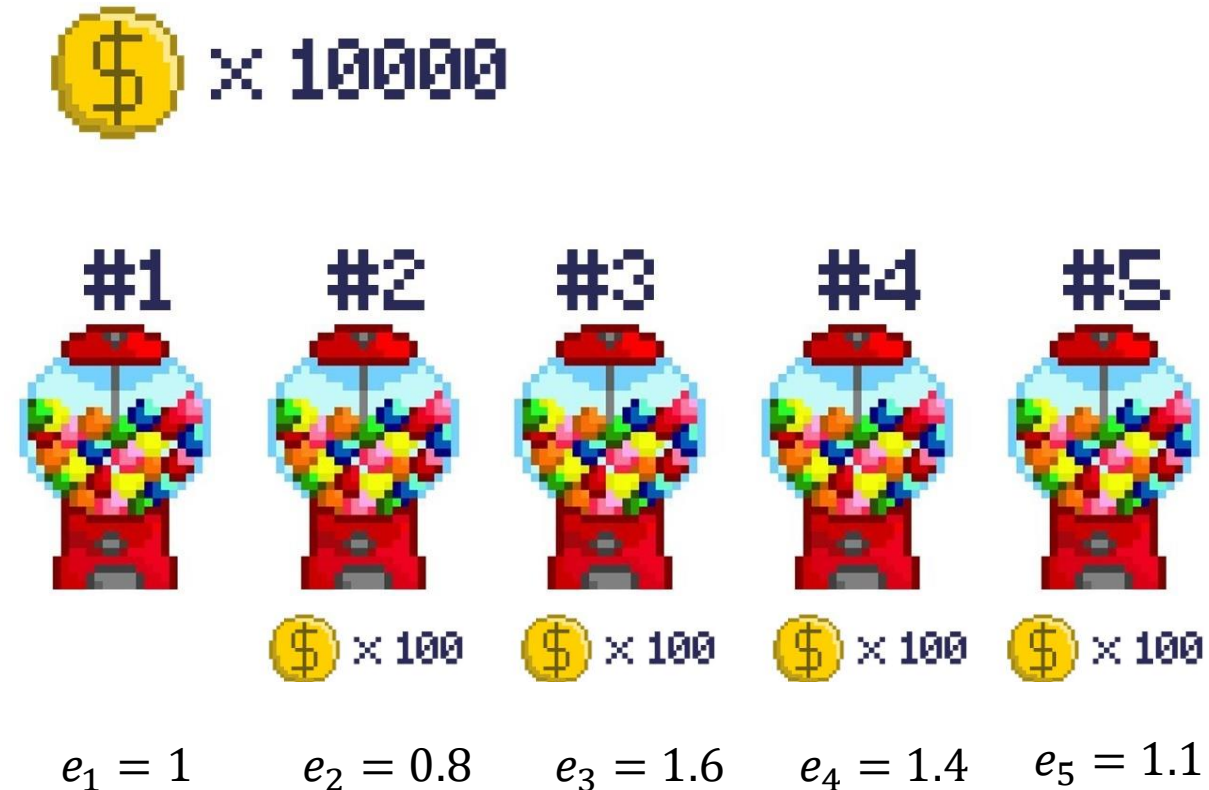
- We will probably **try our machines to guess/estimate which one is the best.**
- And eventually, we will **use all our coins on that “best” machine**, to maximize our number of candies.



Exploration and exploitation tradeoff

Definition (exploration and exploitation):

The phase during which, you try out things to acquire knowledge about the problem is called **exploration**.

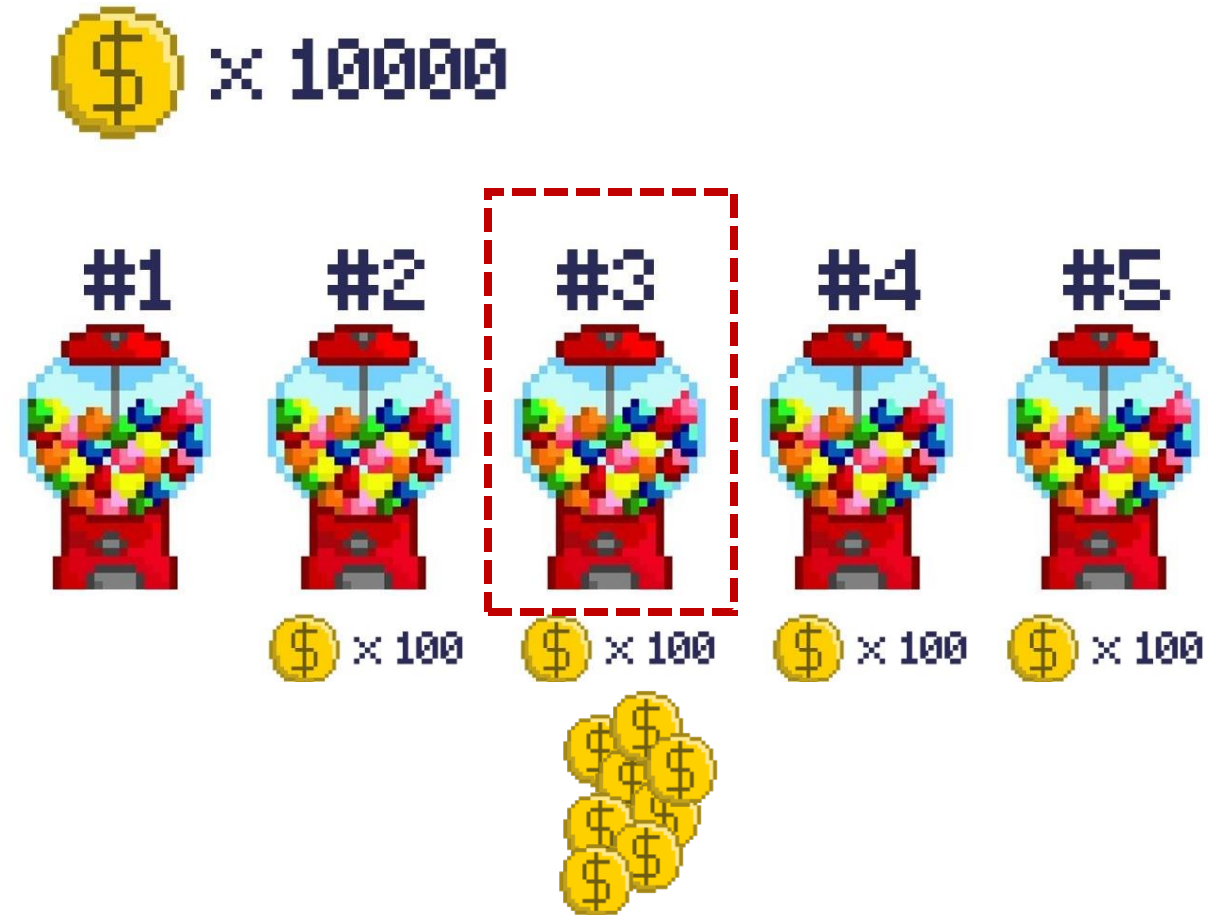


Exploration and exploitation tradeoff

Definition (exploration and exploitation):

The phase during which, you try out things to acquire knowledge about the problem is called **exploration**.

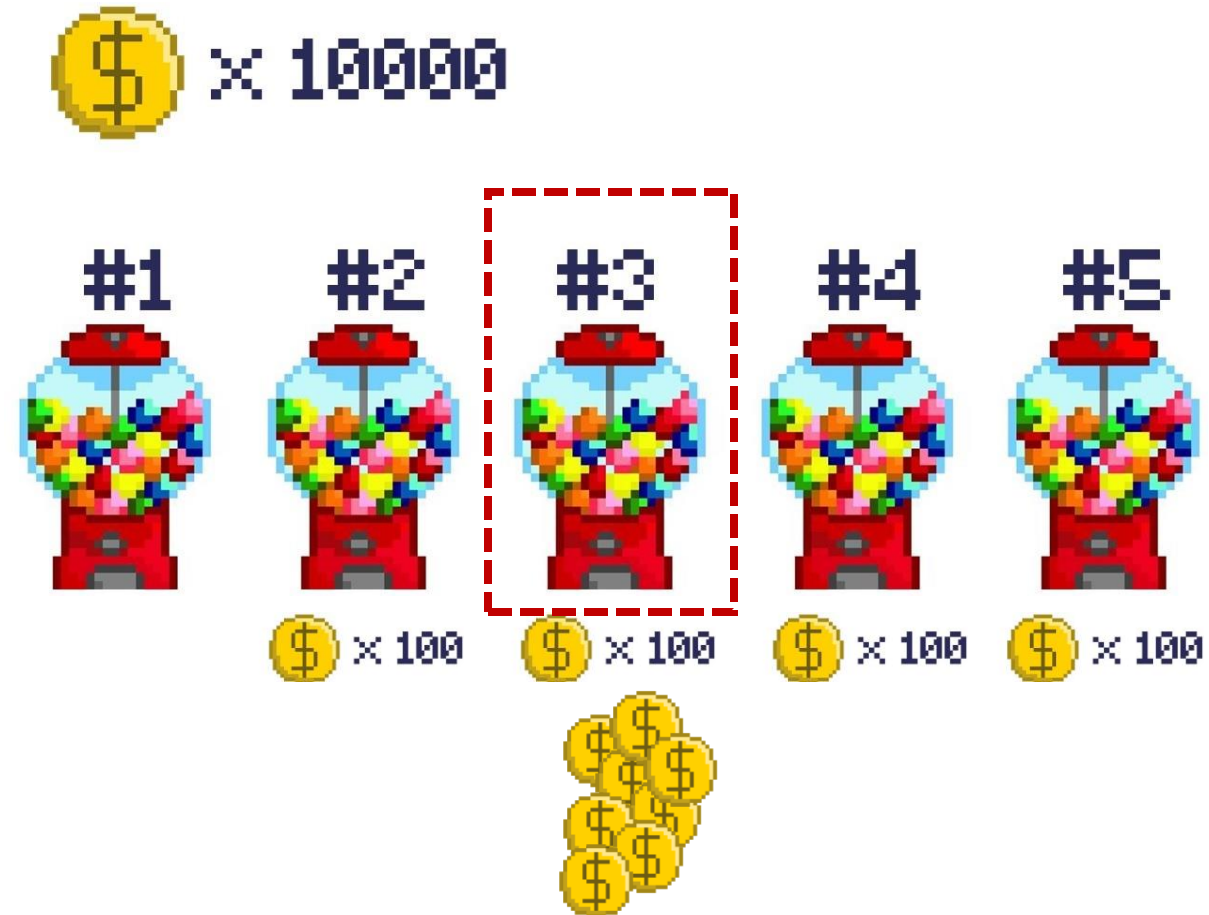
The second phase, where you rely on your acquired knowledge, and play what you feel is the best move, is called **exploitation**.



Exploration and exploitation tradeoff

Definition (exploration vs. exploitation tradeoff):

A good RL-based AI, needs to smartly combine **exploration** and **exploitation** phases.

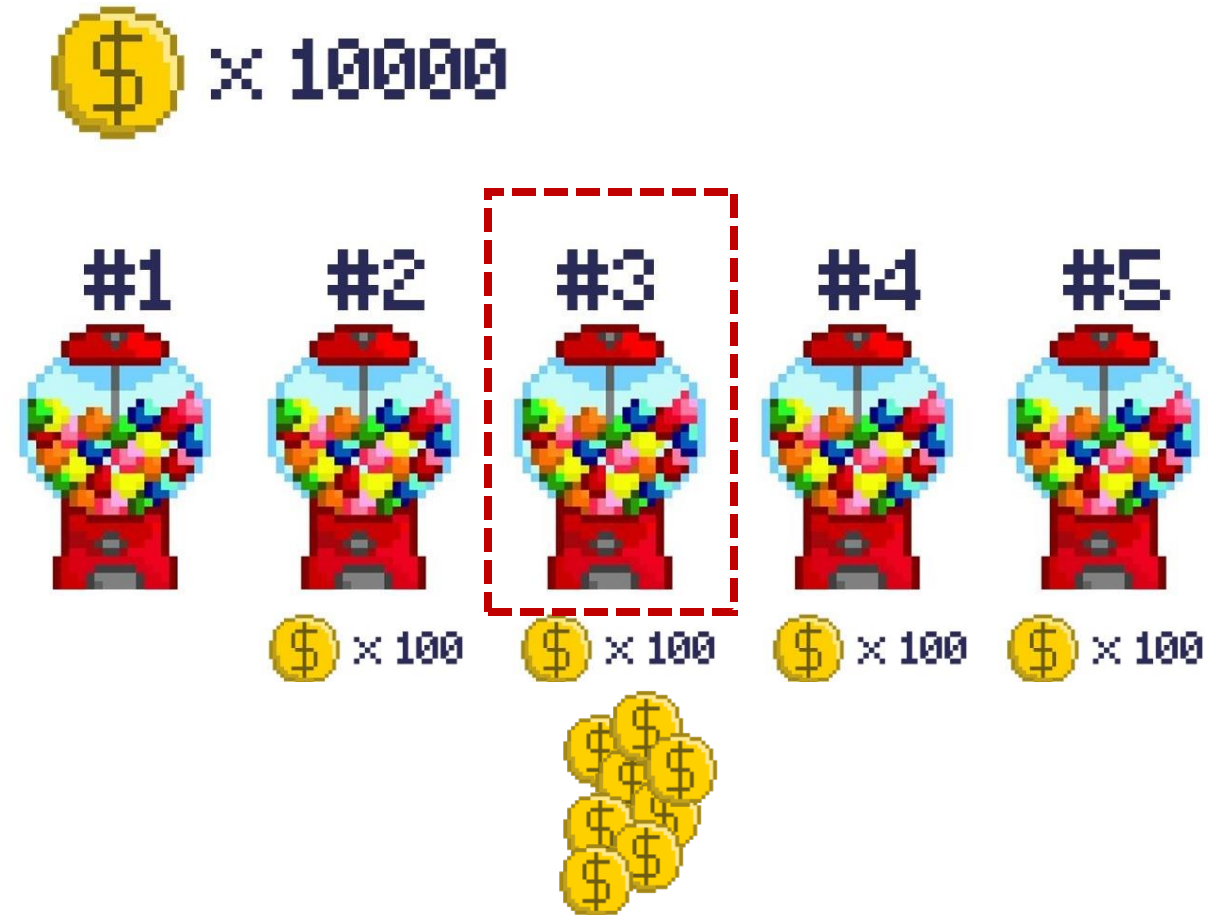


Exploration and exploitation tradeoff

Definition (**exploration vs. exploitation tradeoff**):

A good RL-based AI, needs to smartly combine **exploration** and **exploitation** phases.

- **Too much exploration?** You have wasted coins trying out bad machines.

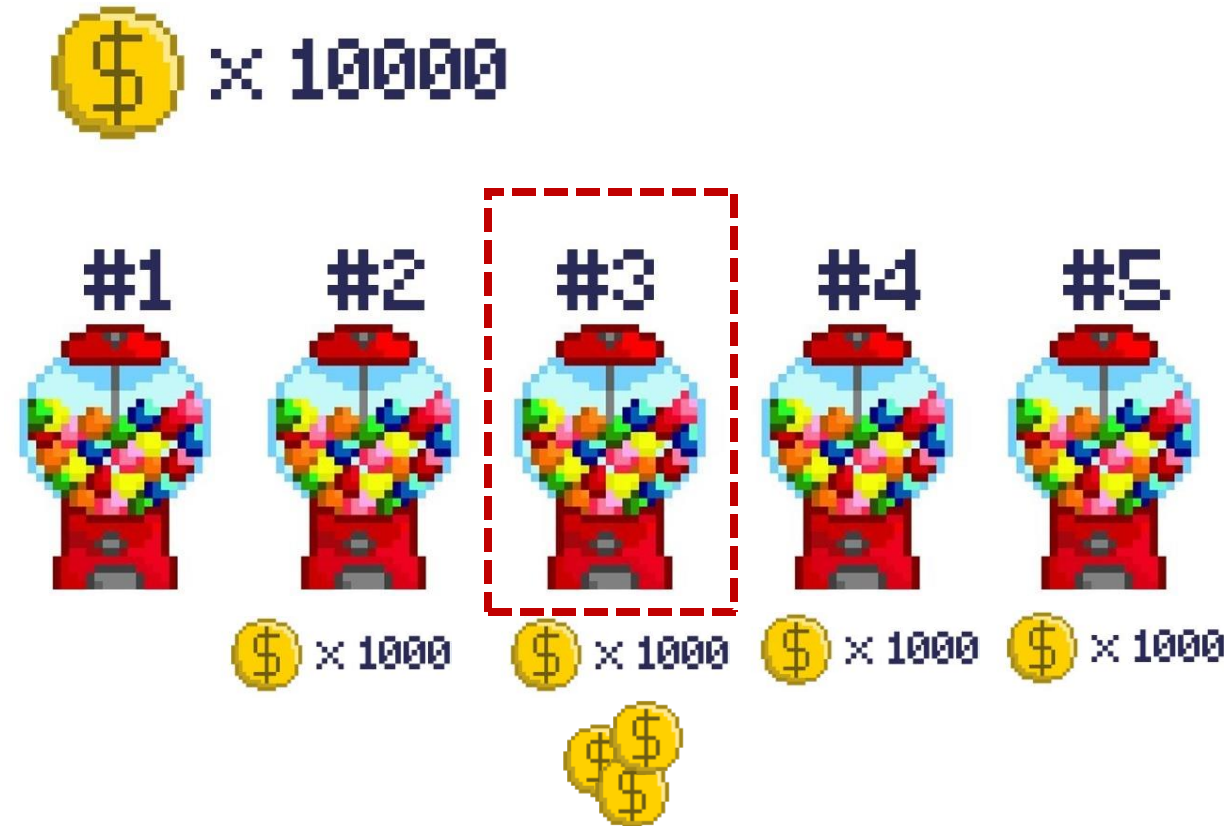


Exploration and exploitation tradeoff

Definition (**exploration vs. exploitation tradeoff**):

A good RL-based AI, needs to smartly combine **exploration** and **exploitation** phases.

- **Too much exploration?** You have wasted coins trying out bad machines.

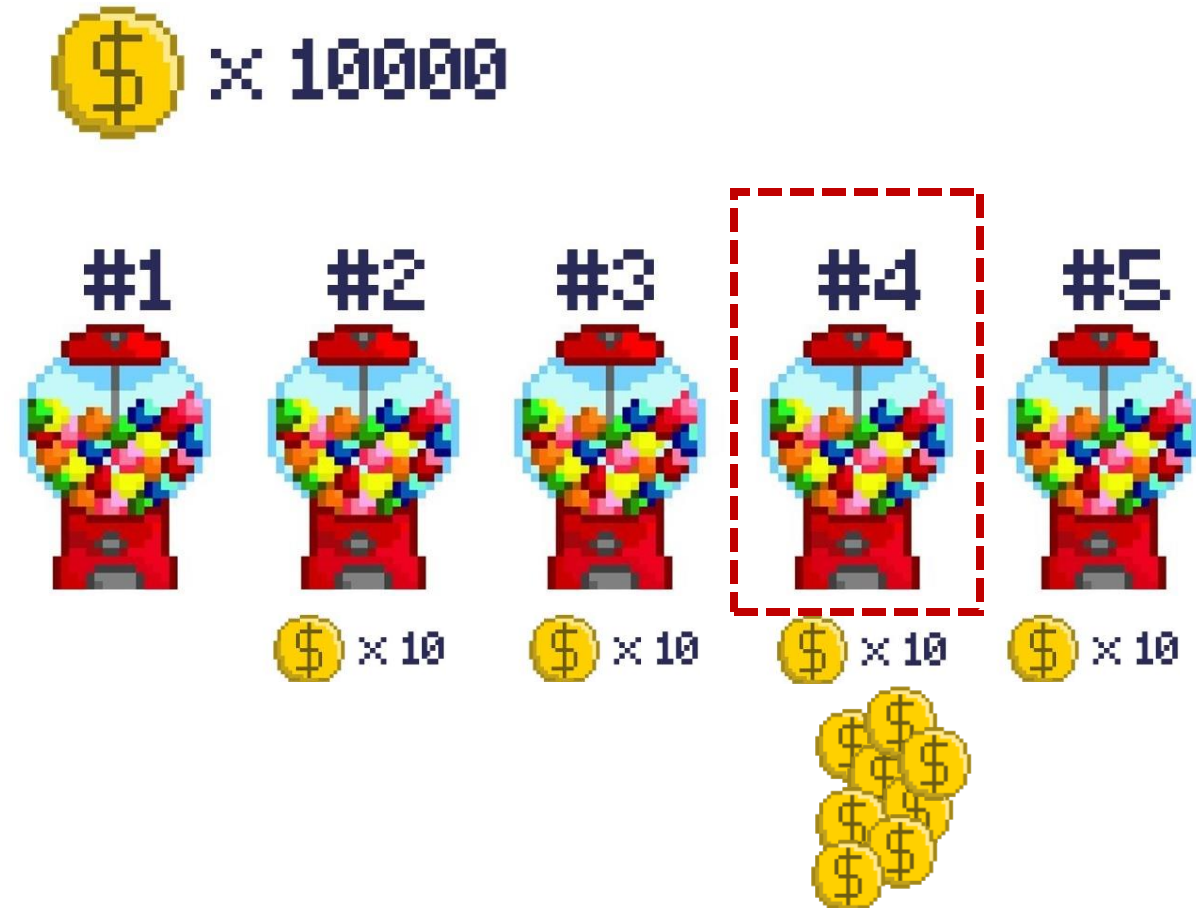


Exploration and exploitation tradeoff

Definition (**exploration vs. exploitation tradeoff**):

A good RL-based AI, needs to smartly combine **exploration** and **exploitation** phases.

- **Too much exploration?** You have wasted coins trying out bad machines.
- **No enough exploration?** You might end up choosing the wrong machine as the “best” one.

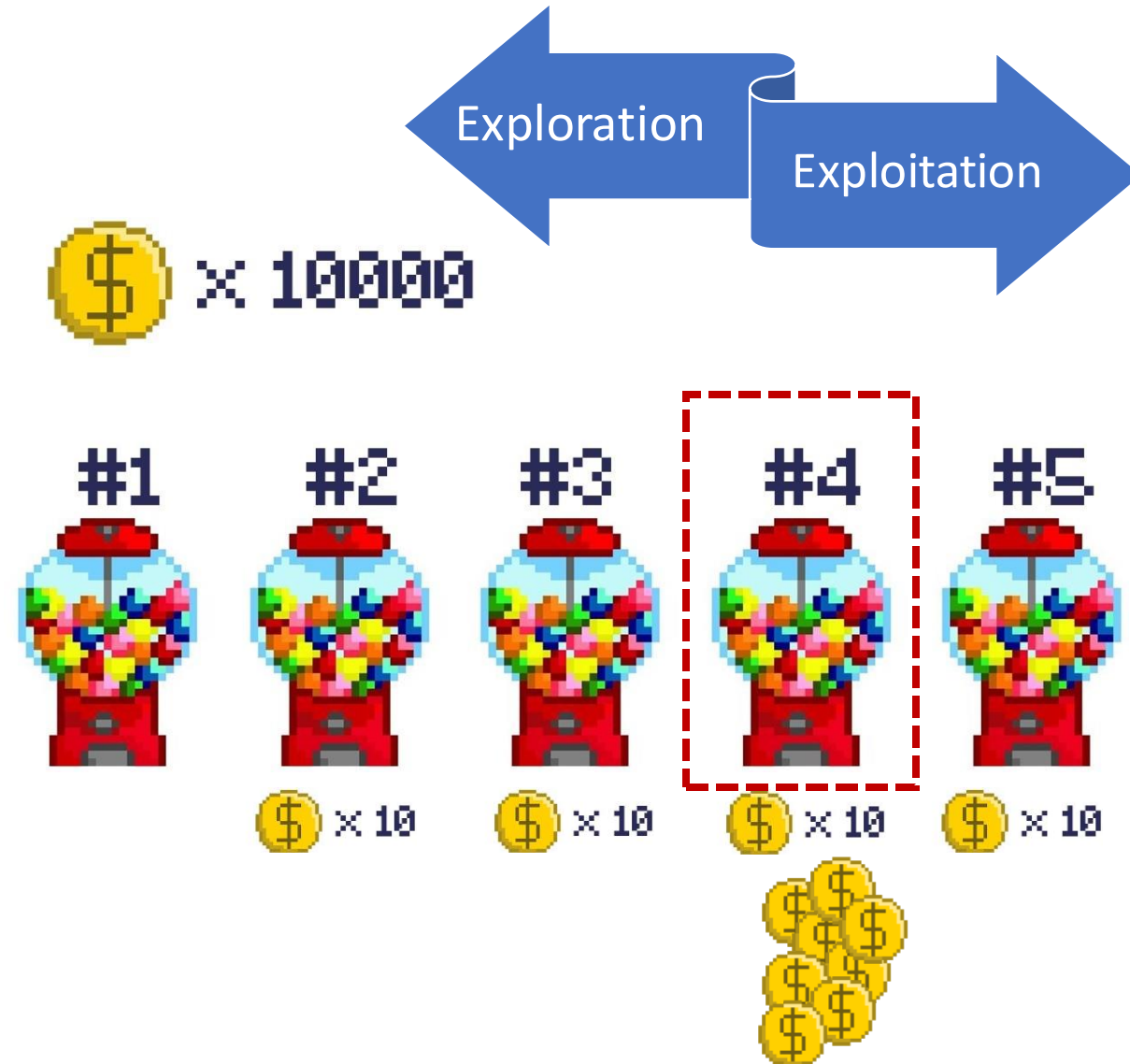


Exploration and exploitation tradeoff

Definition (**exploration vs. exploitation tradeoff**):

A good RL-based AI, needs to smartly combine **exploration** and **exploitation** phases.

- **Too much exploration?** You have wasted coins trying out bad machines.
- **No enough exploration?** You might end up choosing the wrong machine as the “best” one.



Reusing the RL formalism

- Let us denote t , the **index of the coin to be played**.

$$t \in \{1, \dots, 10000\}$$

- The **set of actions** is simply

$$A = \{1, 2, 3, 4, 5\}$$

- We could define our **state**, at time t , as our **current estimates of the expectations** $e_i(t)$ of each machine $i \in A$.

$$s_t = (e_1, e_2, e_3, e_4, e_5)$$

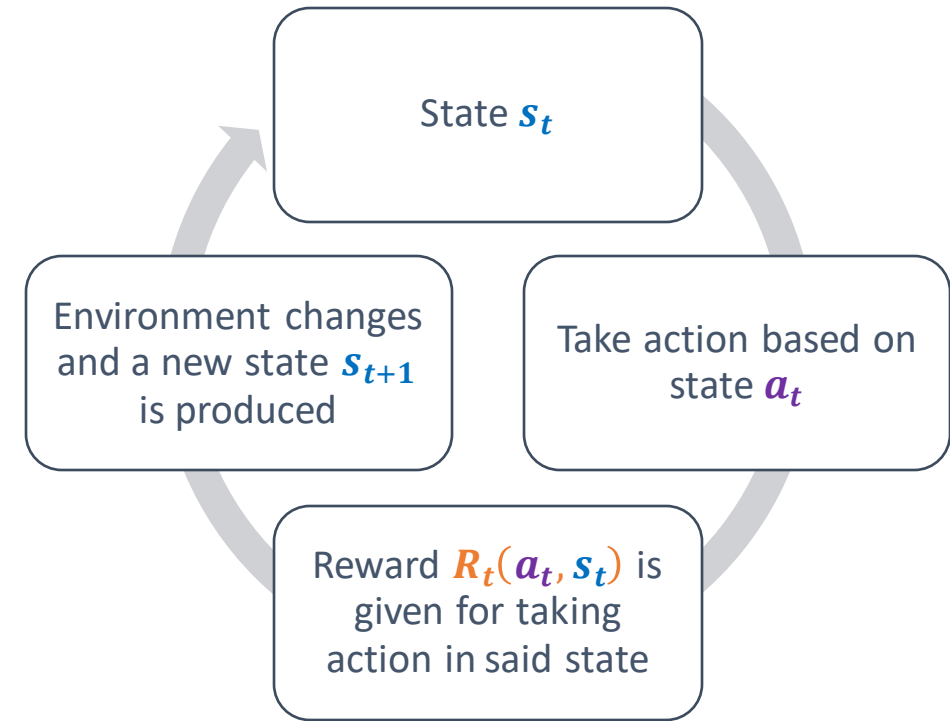


Reusing the RL formalism

Definition (**agent**):

In RL, we refer to the AI, as an **agent**. At each step, the **agent**:

- Looks at the current **state** s_t ,
- Then takes an **action**, in this given state, a_t .
- The **action** has an effect, which is eventually measured in terms of reward, R_t .



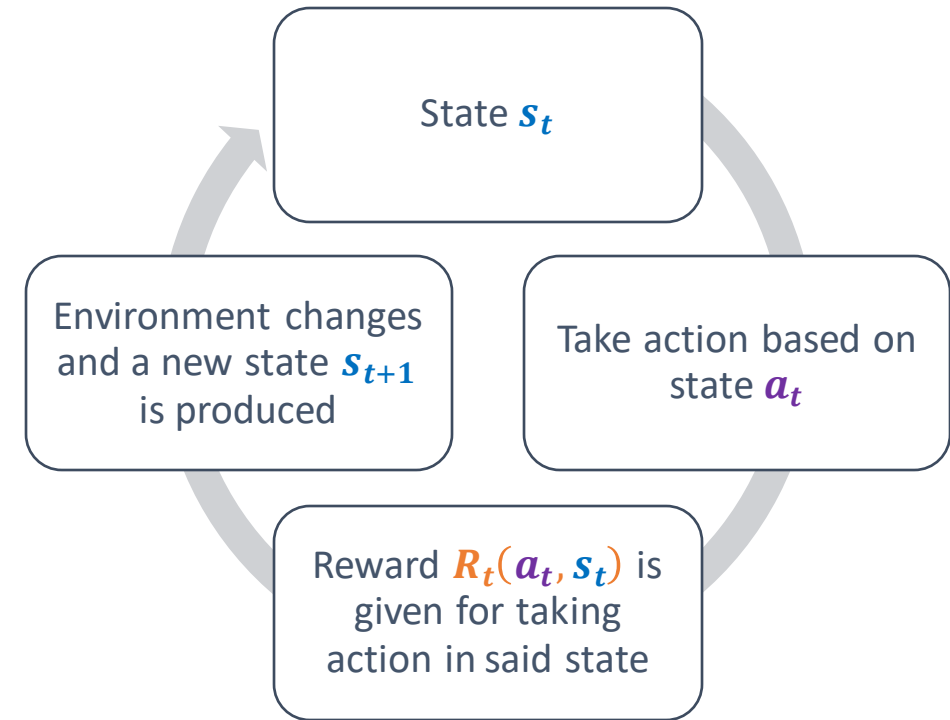
 $\times 10000$



Reusing the RL formalism

Action and reward:

- The **action** $a_t = i$, here, simply consists of **using coin t , on machine i** .
- The **reward** R_t is then defined as **the amount of candies obtained with coin t** , randomly defined with the hidden winning probability of that machine.



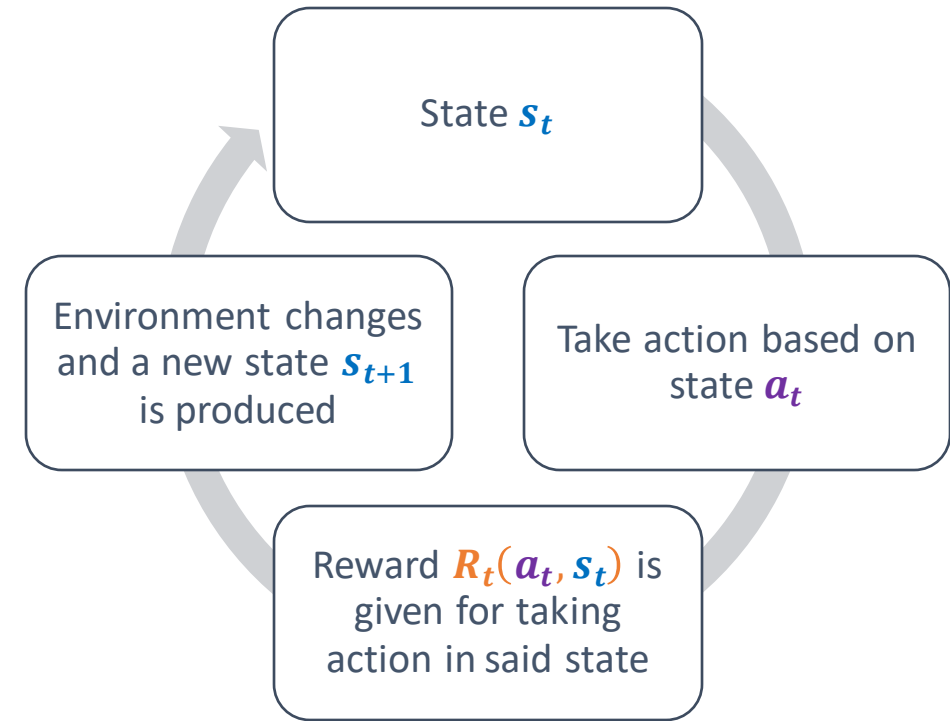
Reusing the RL formalism

State update:

- On each action, update the probability estimates, as

$$\forall i \in A, e_i = \frac{\sum_{k=1}^t \delta_{a_t=i} R_t}{\sum_{k=1}^t \delta_{a_t=i}}$$

This mechanism is used to update the **state** $s_t \rightarrow s_{t+1}$.



Reusing the RL formalism

Million dollar question: what should then be the optimal policy π_t^* , which maximizes the number of candies we hope to obtain?

$$\pi^* = \arg \max_{\pi} \left(E \left[\sum_{t=1}^{10000} R_t(a_t = \pi(t), s_t) \right] \right)$$

Starting with three reference strategies

Three reference strategies can be considered to address this problem.

- **The deterministic strategy:** we always play the deterministic machine, and ignore the random ones.

$$\forall t, \pi_t = 1$$

In a sense, that is the lowest-risk strategy, but misses the opportunity to play some “good” random machines.

Starting with three reference strategies

Three reference strategies can be considered to address this problem.

- **The naive random strategy:** we randomly decide on the machine to play, with uniform probability (20% for each of the five machine).

$$\forall t, \pi_t = \dots ?$$

Starting with three reference strategies

Definition (**stochastic policy**):

We can define a **stochastic policy**, where $\pi_t(a|s)$ is the probability to play action a in state s .

$$\pi_t(a|s) = P[a_t = a | s_t = s]$$

Starting with three reference strategies

Three reference strategies can be considered to address this problem.

- **The naive random strategy:** we randomly decide on the machine to play, with uniform probability (20% for each of the five machine).

Following the stochastic policy definition, $\pi_t(a|s)$ is then defined as:

$$\forall s, \forall a, \quad \pi_t(a|s) = \frac{1}{5}$$

Starting with three reference strategies

Three reference strategies can be considered to address this problem.

- **The deterministic strategy (v2):** we always play the deterministic machine, and ignore the random ones.
- We can rewrite the deterministic strategy using the stochastic policy definition as well.
- (Any deterministic policy can be rewritten as a stochastic one).

$$\forall s, \quad \pi_t(a|s) = \begin{cases} 1 & \text{if } a = 1 \\ 0 & \text{else} \end{cases}$$

Starting with three reference strategies

Three reference strategies can be considered to address this problem.

- **The perfect knowledge strategy:** Let us assume we know about the hidden probability.

In that case, we always play the machine, with the highest expectation in terms and candies.

Using the stochastic policy formalism,

$$\forall s, \quad \pi_t(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{i \in A} (1, 2p_2, 2p_3, 2p_4, 2p_5) \\ 0 & \text{else} \end{cases}$$

This is the upper bound performance strategy, as we cannot do better than this. It is however “cheating”, as it requires prior knowledge.

Performance of reference strategies

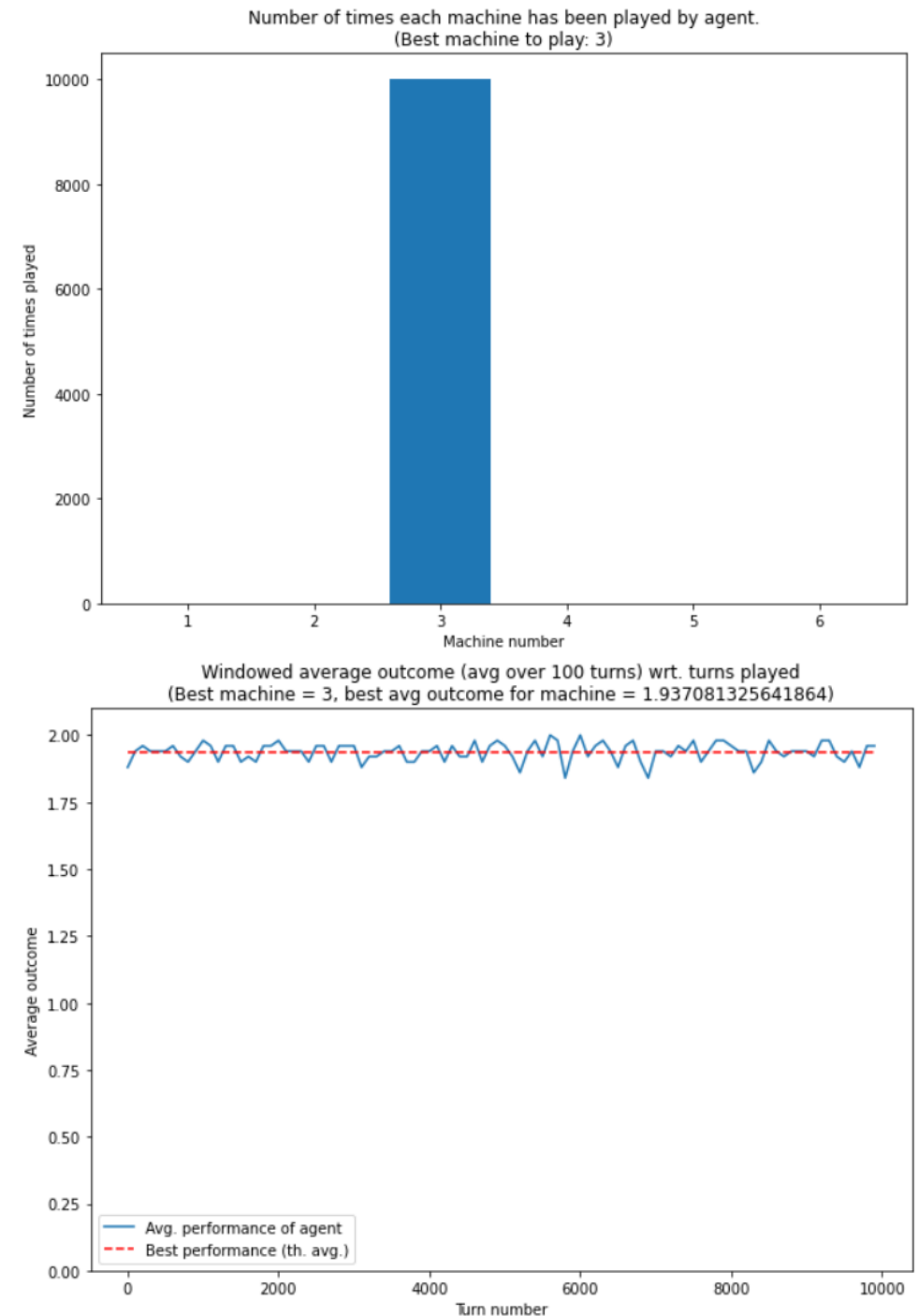
- Consider the six machines (one deterministic and five random machines) below.

Machine_number	Machine_type	Cost	Return_win	Return_loss	Win_probability
1	deterministic	1	1	1	1
2	random	1	2	0	0.873429
3	random	1	2	0	0.968541
4	random	1	2	0	0.869195
5	random	1	2	0	0.530856
6	random	1	2	0	0.232728

- The best machine to play is obviously **machine 3**, here.

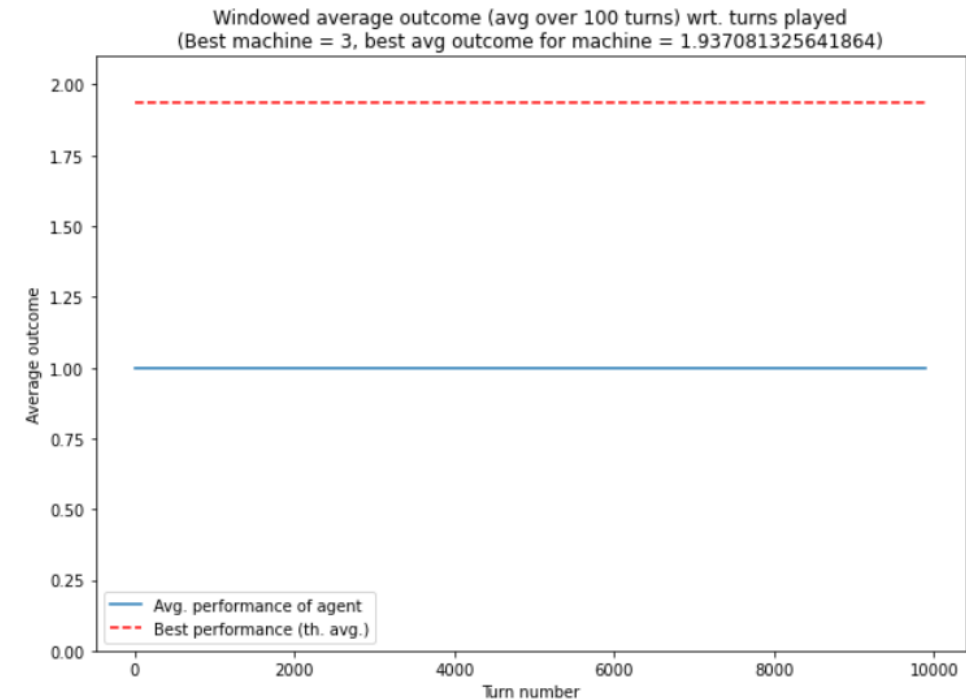
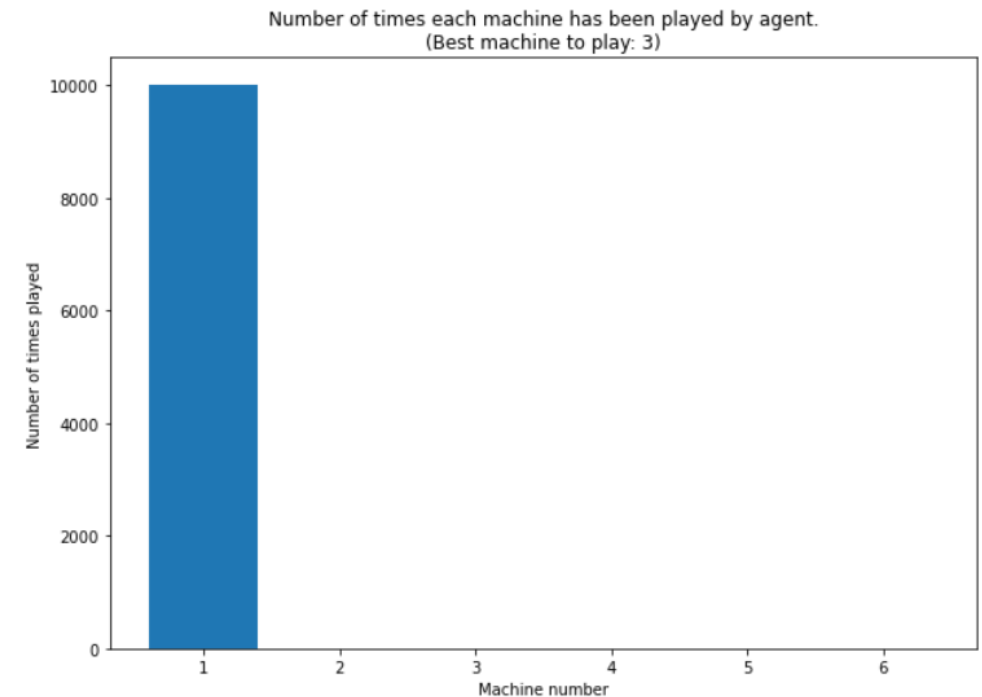
Performance of reference strategies

- The **perfect knowledge strategy** then always plays machine 3.
- It defines the **theoretical upper bound** for all strategies, i.e. the maximal number of candies one can hope to obtain.
- Unfortunately, the perfect knowledge strategy assumes it knows the hidden probabilities.
- All other strategies will not have access to such knowledge!



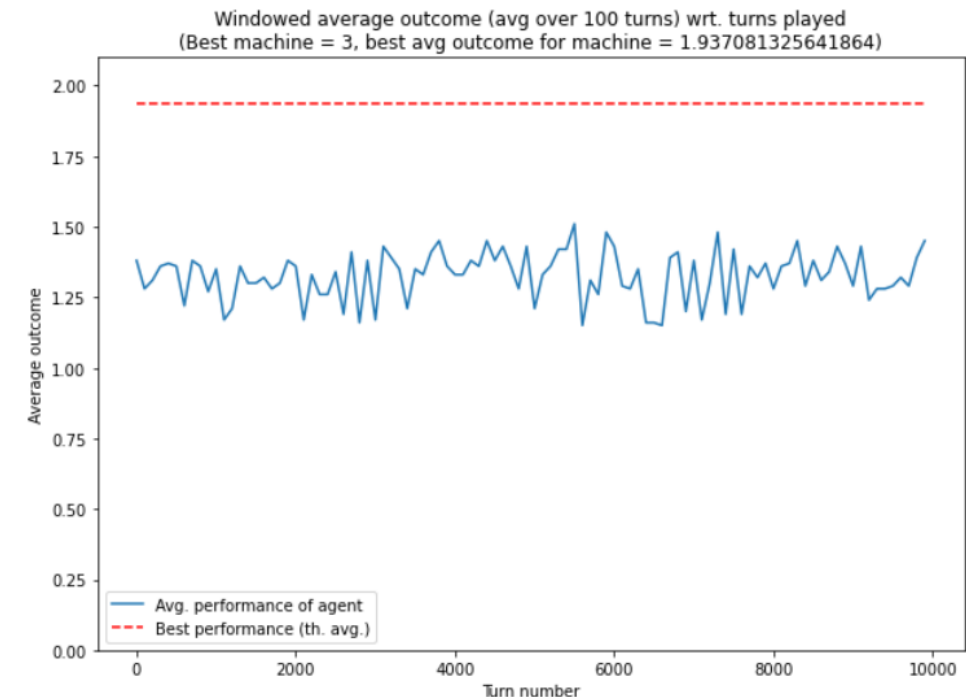
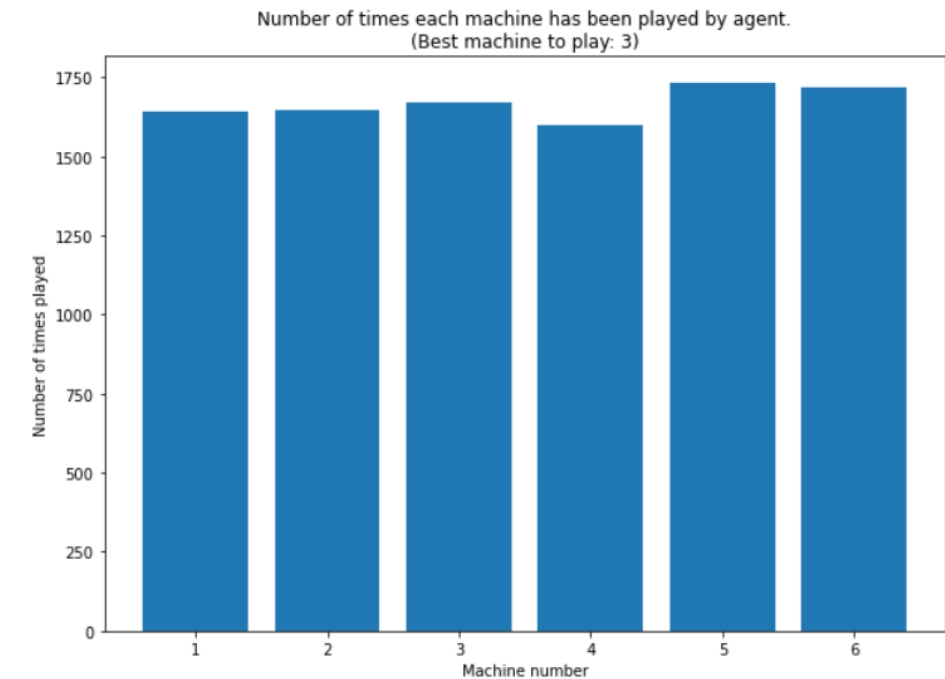
Performance of reference strategies

- The **deterministic strategy** always plays on machine 1.
- Its performance is very low compared to the perfect knowledge one.
- It misses on the opportunity of winning more candies by playing some random machines.
- It defines a **lower bound** for the performance.



Performance of reference strategies

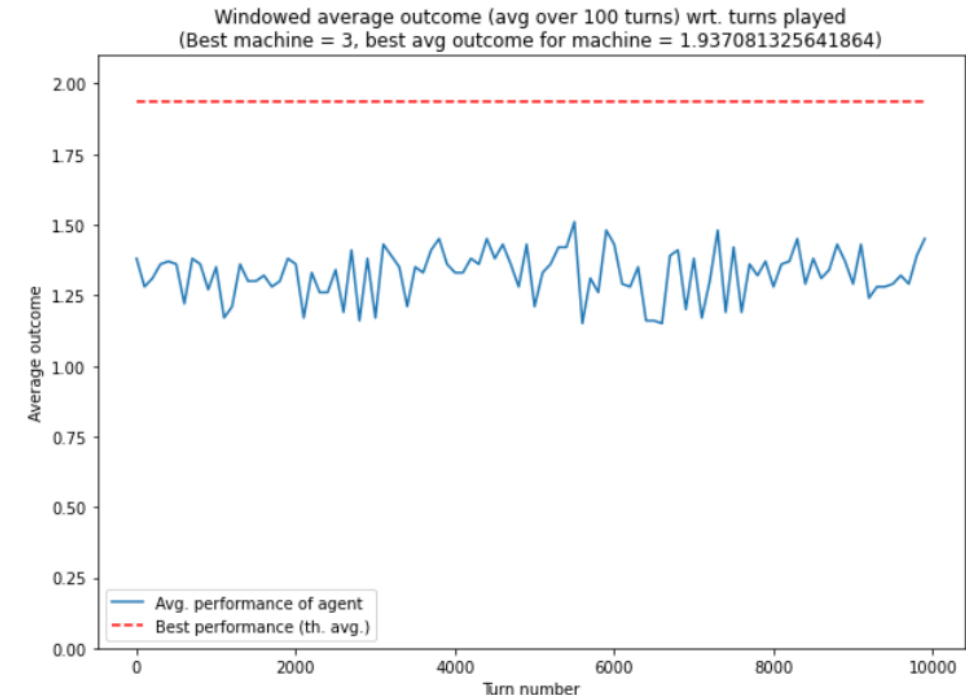
- The **naive random strategy**, plays all machine, roughly the same number of times.
- Note: it performs slightly better than the deterministic one. But this is due to the fact that the random machines were - on that particular run - initialized with probabilities roughly better than the deterministic one.



Performance of reference strategies

- The **naive random strategy**, plays all machine, roughly the same number of times.
- **Note:** it performs slightly better than the deterministic one. But this is due to the fact that the random machines were - on that particular run - initialized with probabilities roughly better than the deterministic one.

Machine_number	Machine_type	Cost	Return_win	Return_loss	Win_probability
1	deterministic	1	1	1	1
2	random	1	2	0	0.873429
3	random	1	2	0	0.968541
4	random	1	2	0	0.869195
5	random	1	2	0	0.530856
6	random	1	2	0	0.232728



Performance of reference strategies

Overall, our deterministic and random strategies struggle because

- The deterministic one refuses to explore before exploiting, as it plays it “safe”, using the deterministic machine only.
- It then misses on the opportunity to play random machines, which might have better performance.
- The naive random strategy explores all machines, but is unable to exploit any knowledge it might obtain from the exploration.
- More advanced strategies will have to combine exploration and exploitation in an intelligent manner.

Exploring and exploiting (strategy #1)

Definition (ϵ -first strategy):

In the ϵ -first strategy, we have two phases.

- **Exploration:** use a proportion $\epsilon \in [0,1]$ of our 10000 coins, randomly distributed over all the random machines, with a uniform distribution.
- **Exploitation:** then, use all remaining coins on the machine with the highest estimate $e_i(t)$.

$$\forall i \in A, e_i(t) = \frac{\sum_{k=1}^t \delta_{a_t=i} R_t}{\sum_{k=1}^t \delta_{a_t=i}}$$

Exploring and exploiting (strategy #1)

Definition (ϵ -first strategy):

The ϵ -first strategy, has the following stochastic policy.

$$\forall s, \forall a \in A, \quad \pi_t(a|s) = \begin{cases} 1/5 & \text{if } t \leq 10000\epsilon \\ 1 & \text{if } t > 10000\epsilon \\ & \text{and } a = \arg \max_{i \in A} (e_i(t)) \\ 0 & \text{else} \end{cases}$$

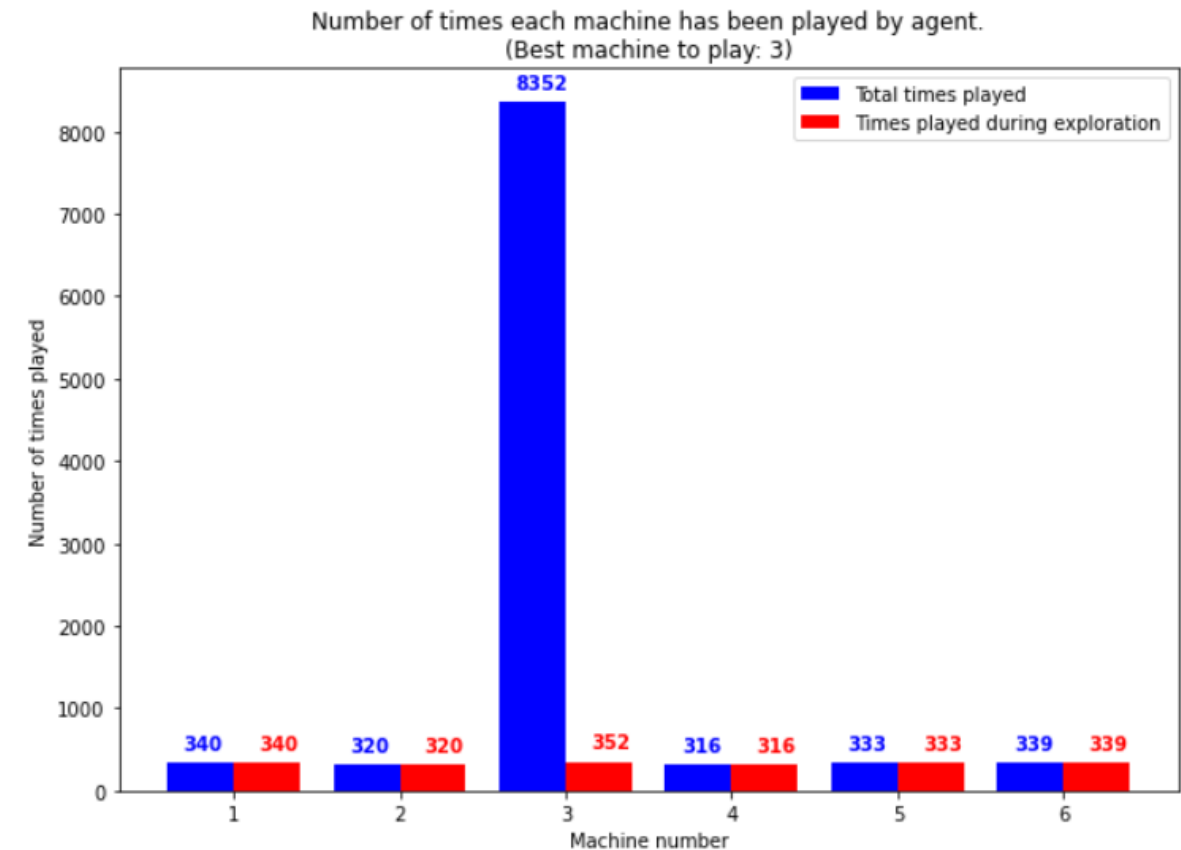
$$\text{with, } \forall i \in A, e_i(t) = \frac{\sum_{k=1}^t \delta_{a_t=i} R_t}{\sum_{k=1}^t \delta_{a_t=i}}$$

Exploring and exploiting (strategy #1)

Definition (ϵ -first strategy):

In the ϵ -first strategy, we have two phases.

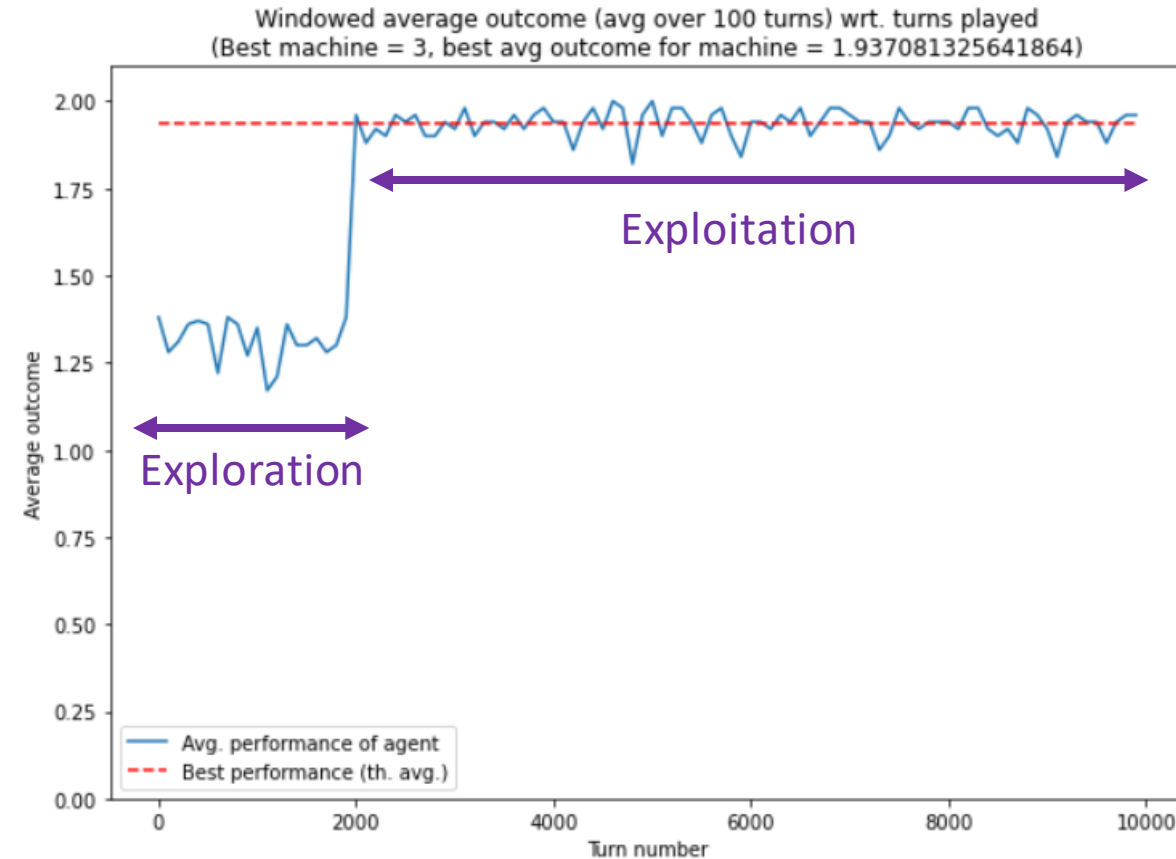
- **Exploration:** use a proportion $\epsilon \in [0,1]$ of our 10000 coins, randomly distributed over all the random machines, with uniform distribution.
- **Exploitation:** then, use all remaining coins on the machine with the highest estimate $e_i(t)$.



Exploring and exploiting (strategy #1)

The **ϵ -first strategy**, explores by using 2000 coins.

After the exploration phase, it is then able to identify the best machine immediately and play it with the remaining coins.

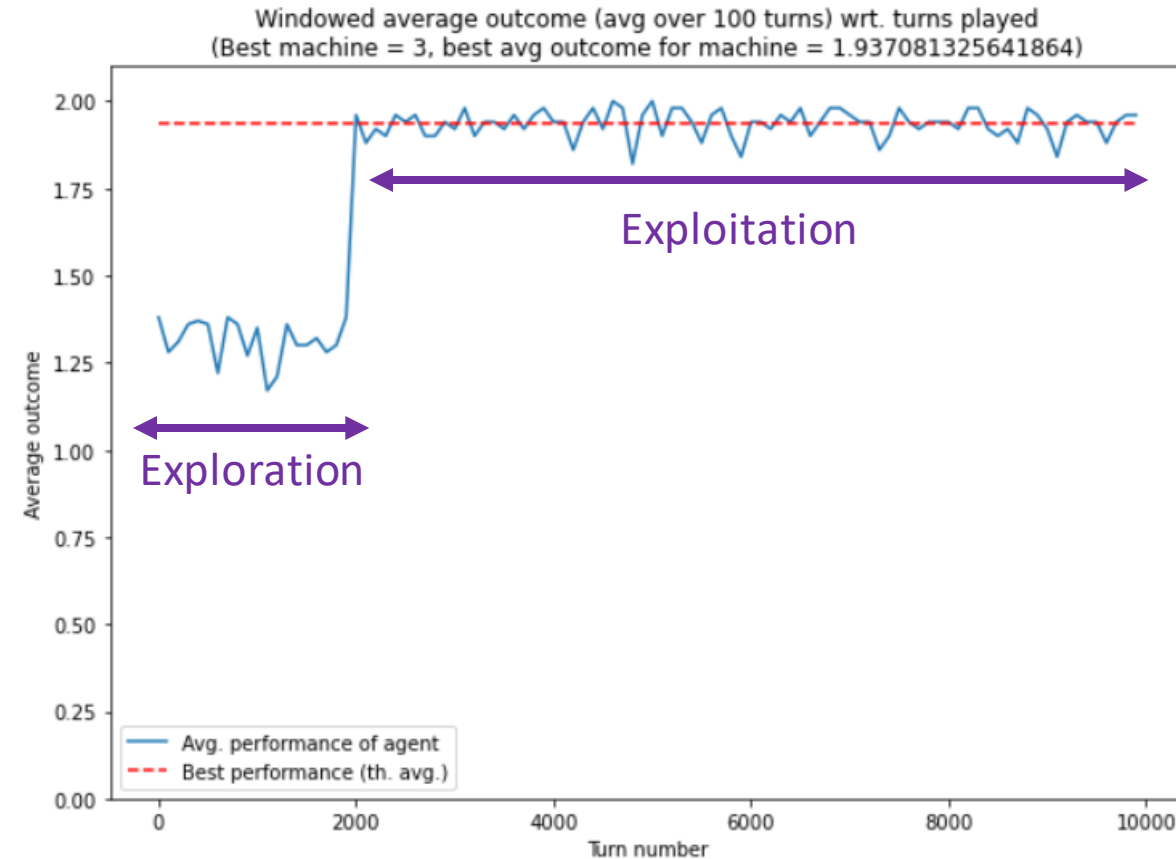


Exploring and exploiting (strategy #1)

The **ϵ -first strategy**, explores by using 2000 coins.

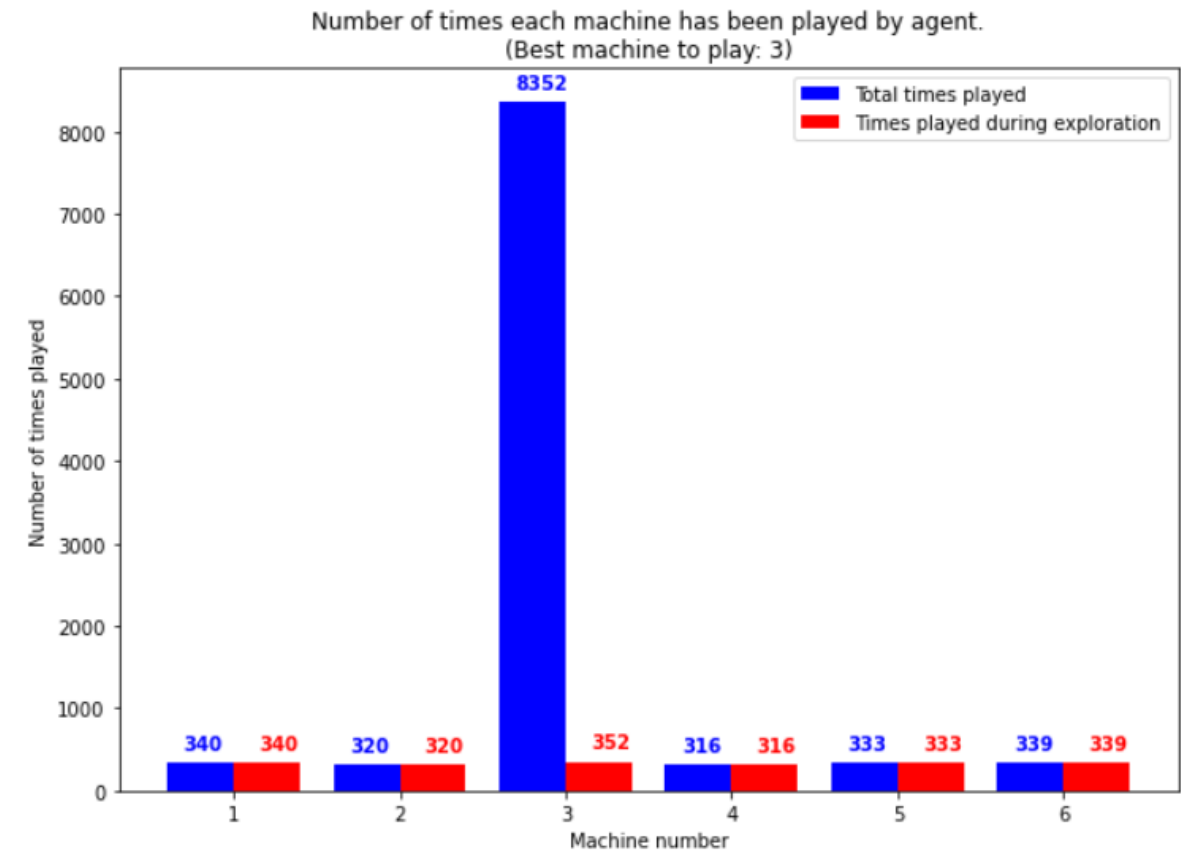
After the exploration phase, it is then able to identify the best machine immediately and play it with the remaining coins.

Ultimately, the objective is to explore for a while, and then go into exploitation, eventually matching the performance of the upper bound strategy.



Exploring and exploiting (strategy #1)

- **Problem:** during the exploration phase, we keep on using coins on machines, which have shown poor results.
- It would be preferable to explore, by giving more coins to machines which have shown good results.
- Dropping machines, which feel bad, therefore not wasting any more coins on them.



Exploring and exploiting (strategy #2)

Definition (ϵ -first strategy with softmax exploration):

In the ϵ -first strategy with softmax exploration, we have two phases.

- **Exploration:** use a proportion $\epsilon \in [0,1]$ of our 10000 coins, randomly distributed over all the random machines, with softmax distribution.

$$\forall t \leq 10000 \cdot \epsilon, \forall i \in A, \quad \pi_t(i|s) = \frac{\exp(e_i(t))}{\sum_{k \in A} \exp(e_k(t))}$$

- **Note:** if we have never played the machine before, initialize it as $e_i(t) = 1$, which is the same value as the deterministic machine.

Exploring and exploiting (strategy #2)

Definition (ϵ -first strategy with softmax exploration):

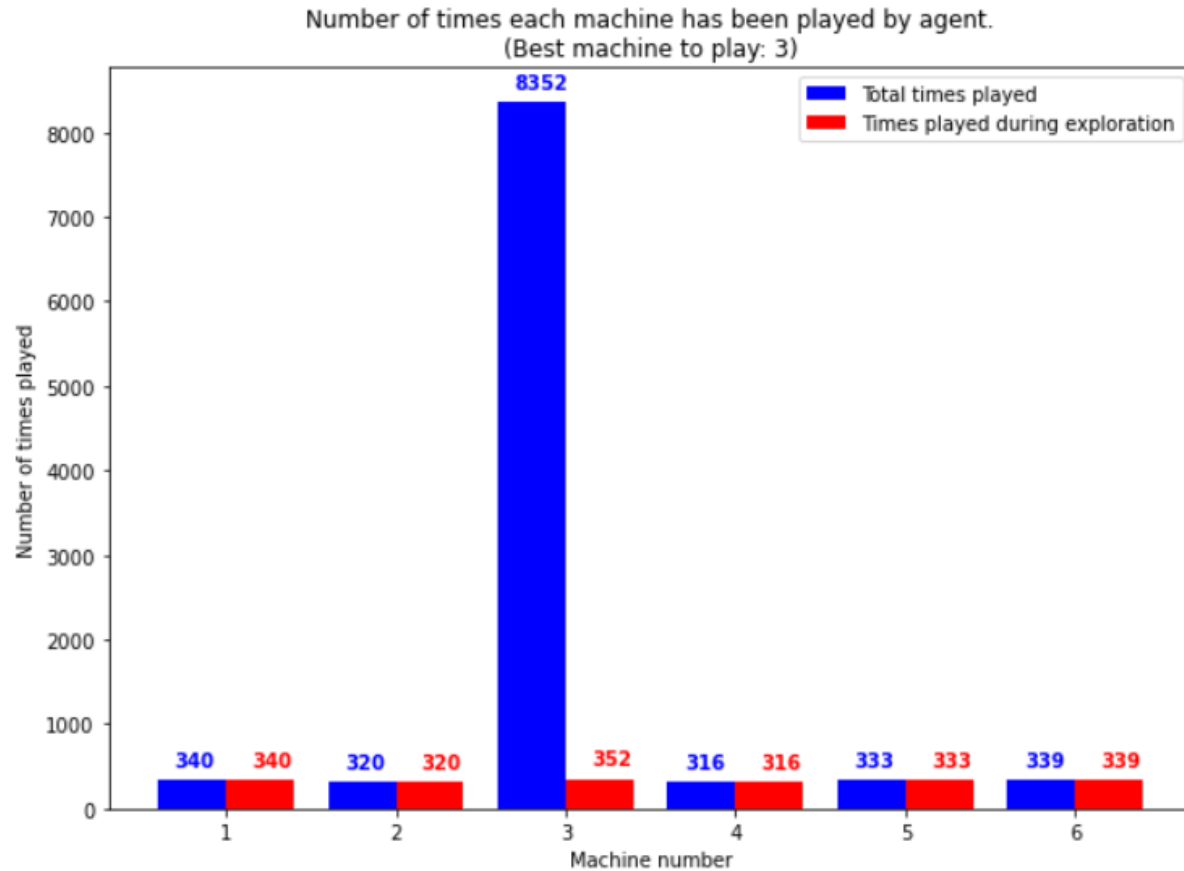
In the ϵ -first strategy with softmax exploration, we have two phases.

- **Exploration:** use a proportion $\epsilon \in [0,1]$ of our 10000 coins, randomly distributed over all the random machines, with softmax distribution.

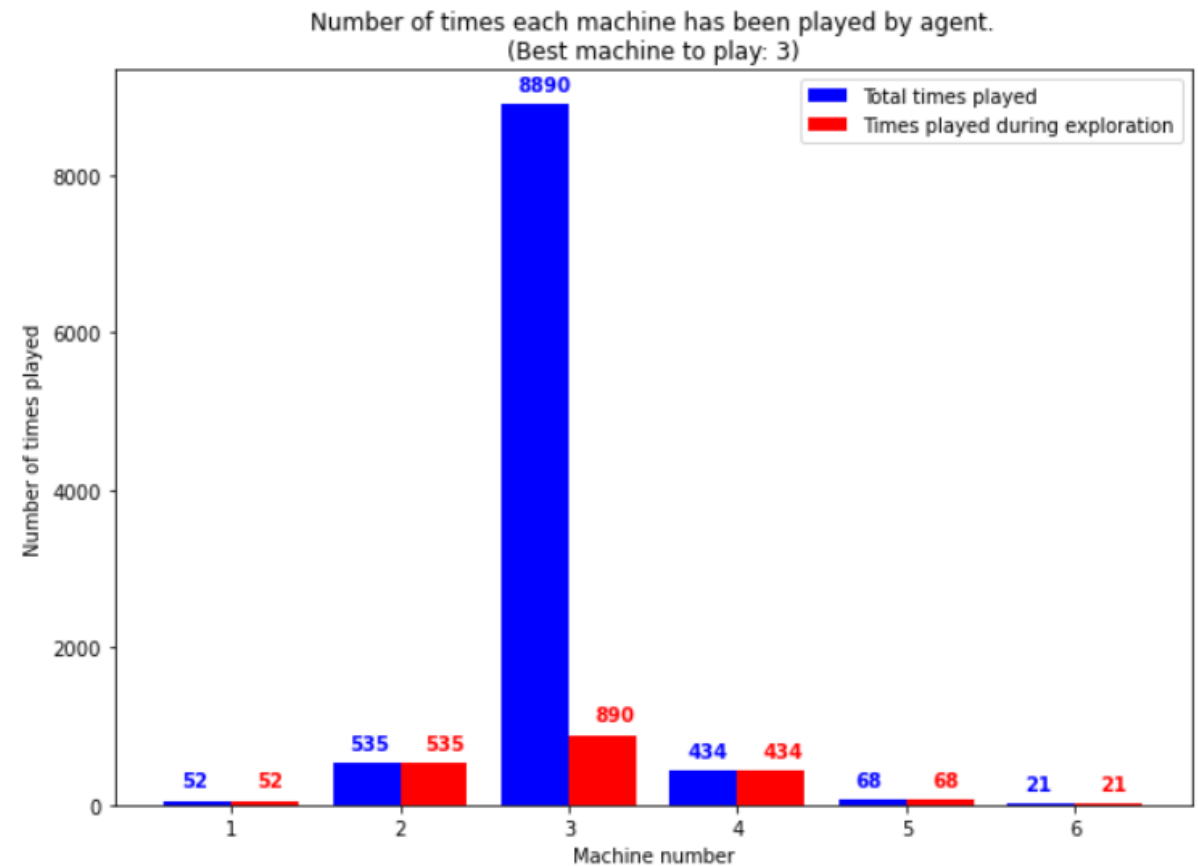
$$\forall t \leq 10000 \cdot \epsilon, \forall i \in A, \quad \pi_t(i|s) = \frac{\exp(e_i(t))}{\sum_{k \in A} \exp(e_k(t))}$$

- **Exploitation:** then, use all remaining coins on the machine with the highest estimate $e_i(t)$.

Exploring and exploiting (strategy #2)

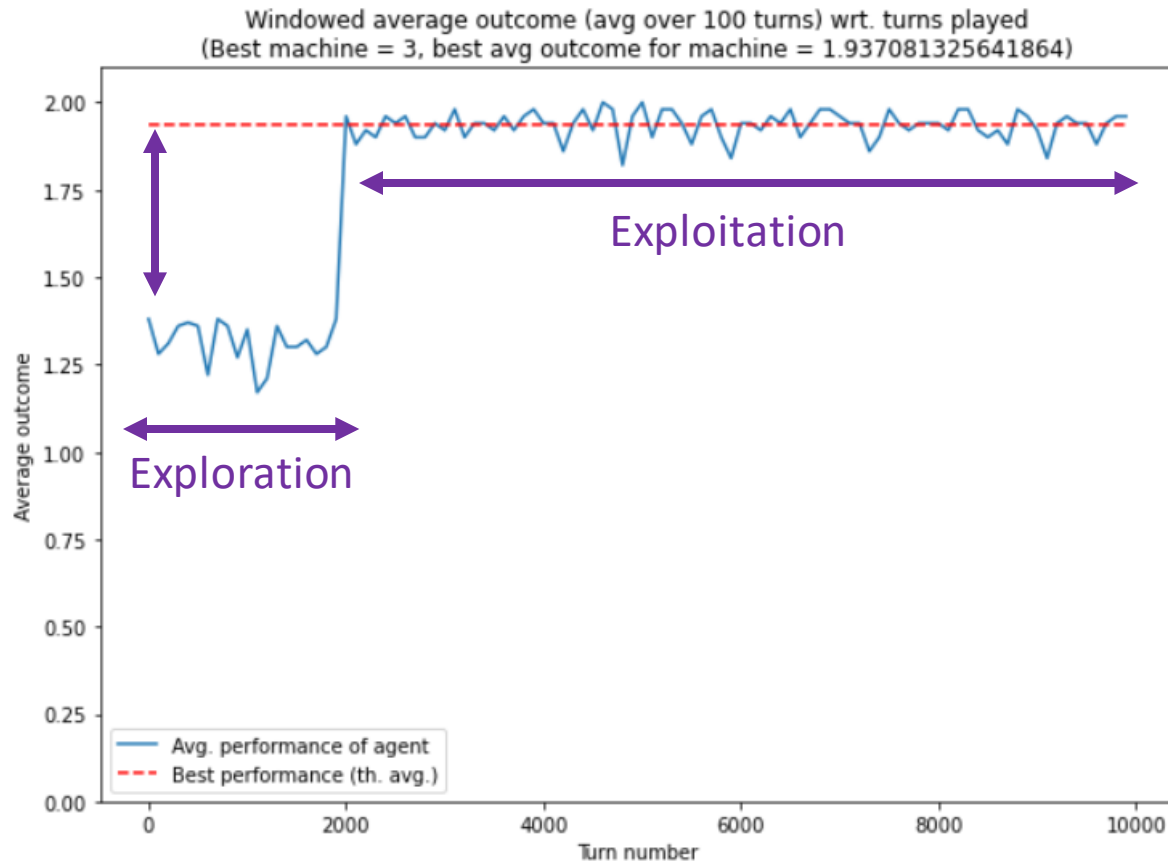


ϵ -first strategy

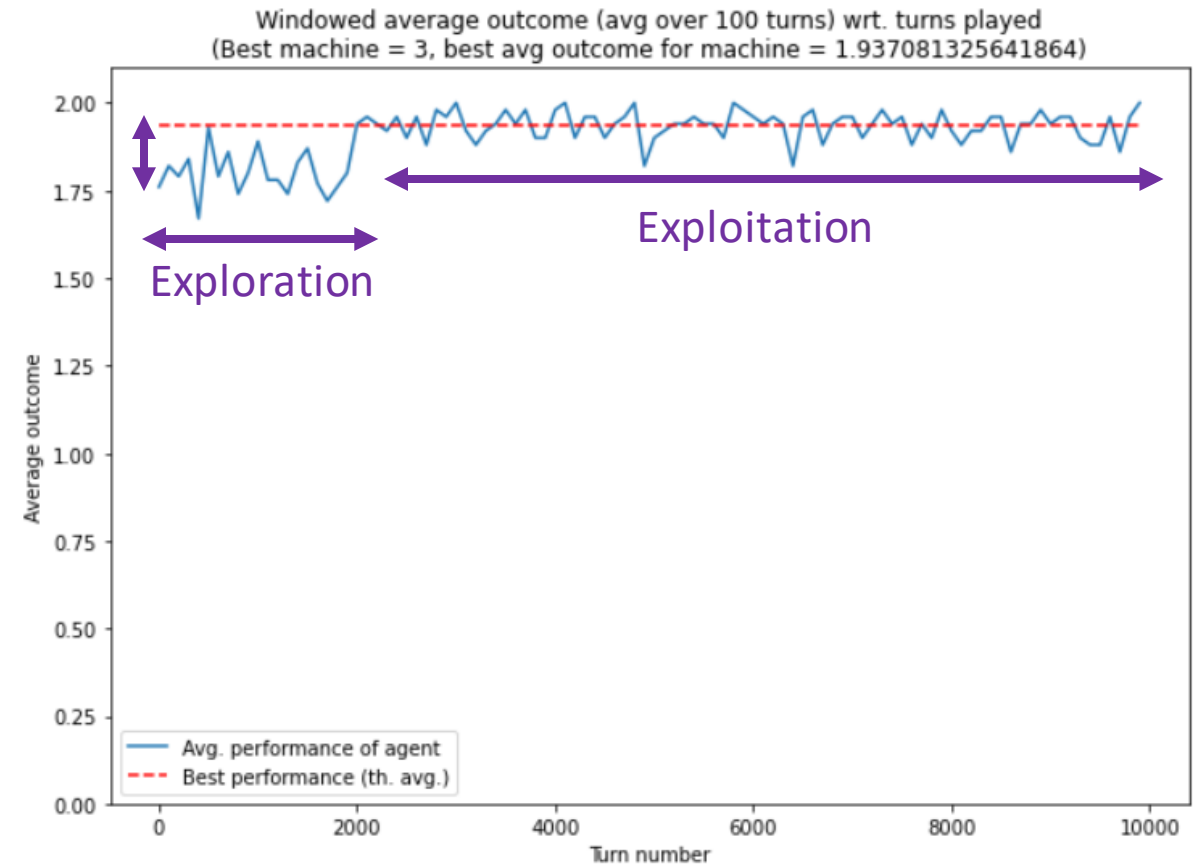


ϵ -first strategy with
softmax exploration

Exploring and exploiting (strategy #2)



ϵ -first strategy



ϵ -first strategy with
softmax exploration

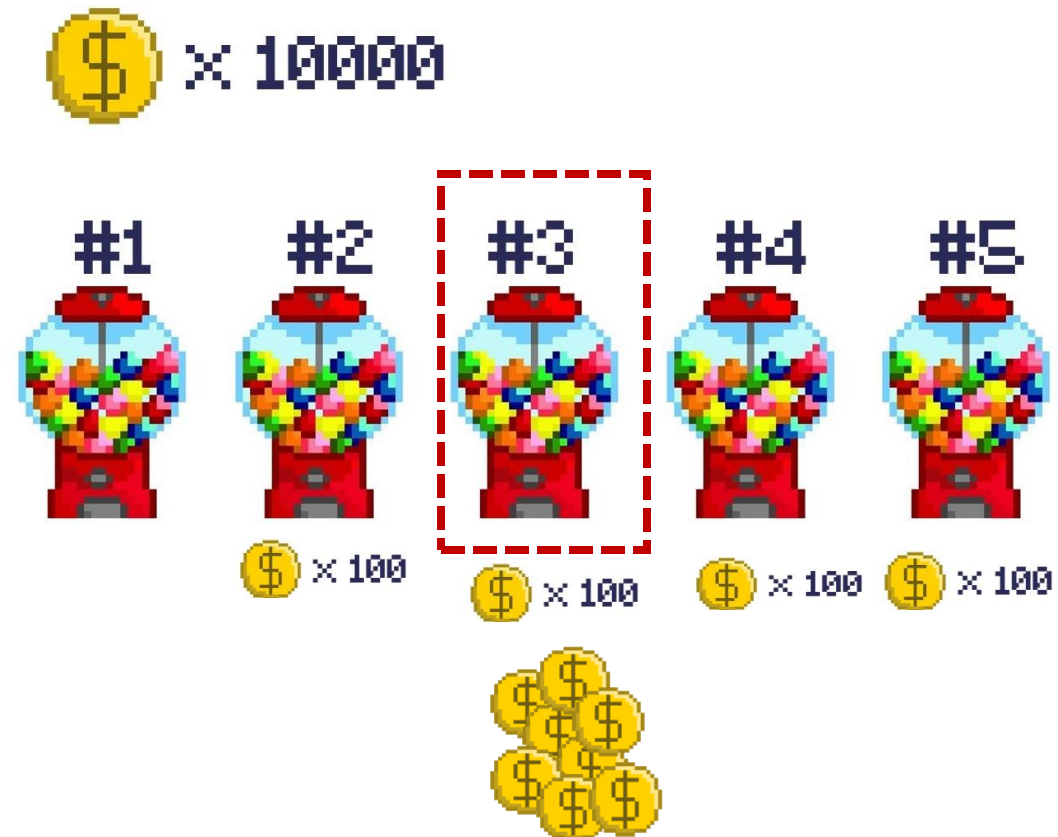
On regret and “cost of knowledge”

Definition (**regret**):

Unfortunately, we will always be wasting coins to understand that some machines are bad and should not be played.

We define the **regret**, as the **quantity of candies we have missed**, because we have not played the best machine, but another one.

It is the **performance difference between a given strategy and the perfect knowledge one**.



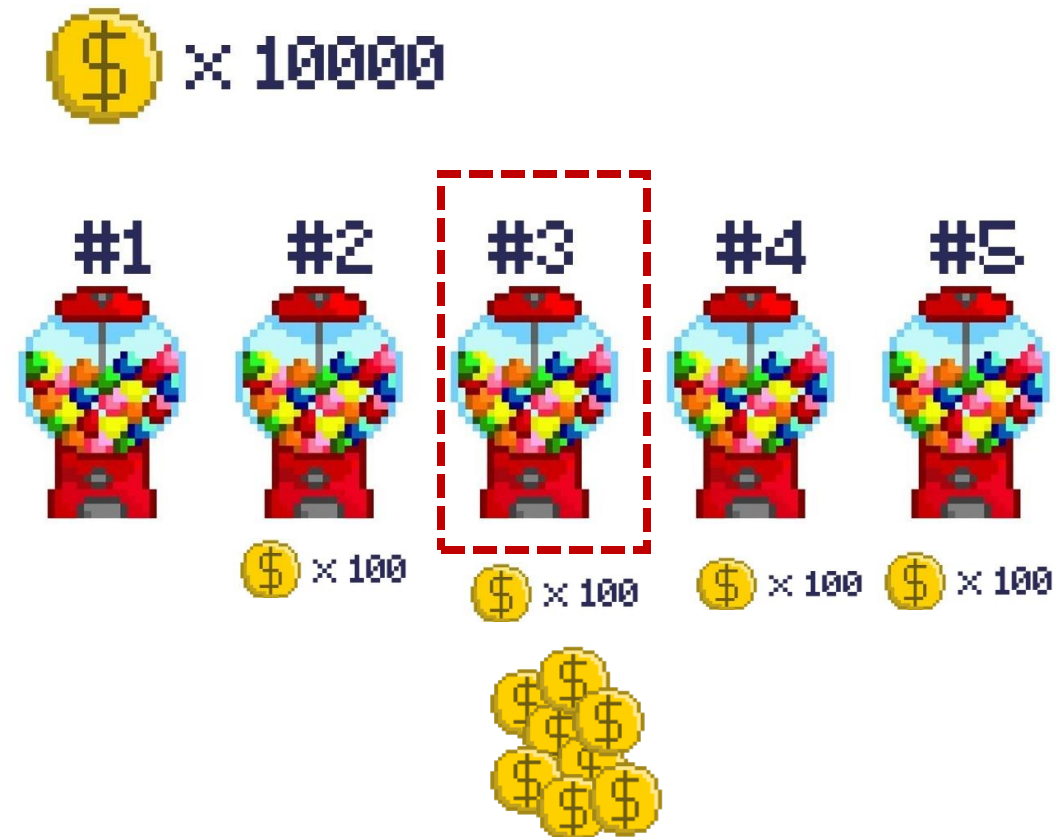
On regret and “cost of knowledge”

Definition (“cost of knowledge”):

The **regret** is the amount of candies we have to **waste to obtain information** about the hidden probabilities of the machines.

Objective: we want to find the strategy, which minimizes the **regret**.

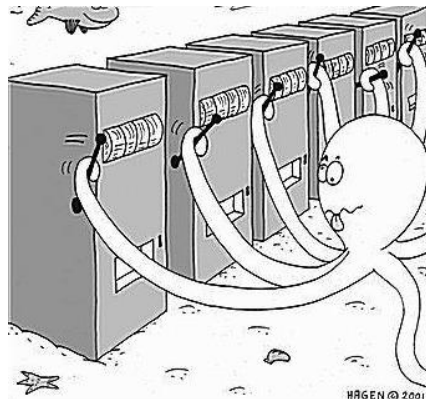
The **regret** cannot be zero, and its **minimal value** is the “**cost of knowledge**”, what one must pay to acquire the missing knowledge.



Some more advanced concepts and strategies

Definition (**multi-armed bandit problem**):

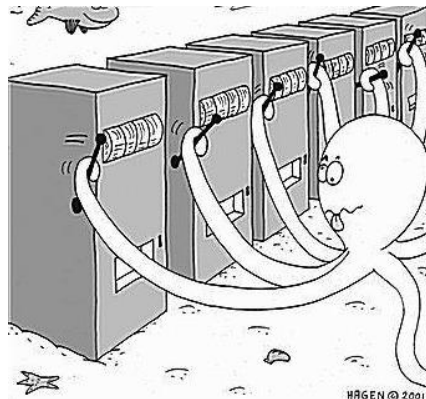
Our random candy machine problem is a well-known mathematical problem, commonly referred to as the **multi-armed bandit problem**.



Some more advanced concepts and strategies

Definition (**multi-armed bandit problem**):

Our random candy machine problem is a well-known mathematical problem, commonly referred to as the **multi-armed bandit problem**.



- It consists of is a problem in which a fixed limited set of resources must be allocated between alternative choices in a way that maximizes their expected gain.
- Each choice's properties are only partially known at the time of allocation, and may become better understood as time passes or by allocating resources to each possible choice.

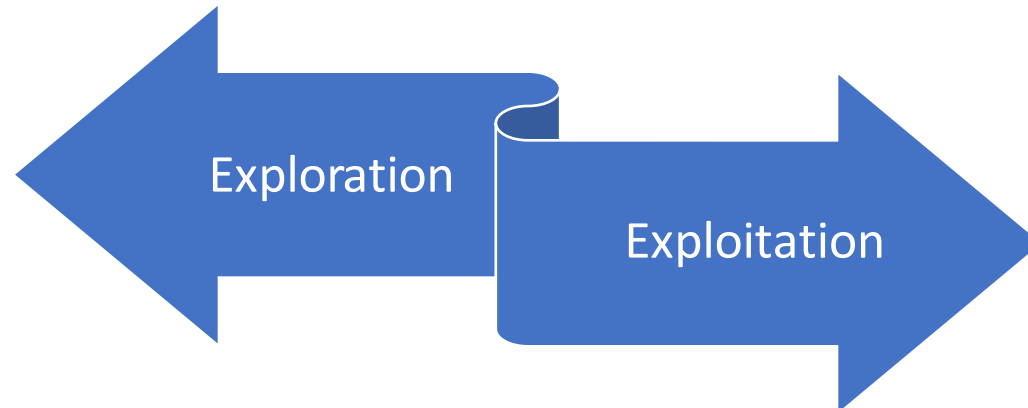
Some more advanced concepts and strategies

Some advanced questions

- **How to decide of the value ϵ , in the ϵ -first strategies?**

If ϵ is large, waste of coins, high regret during exploration phase.

If too low, risk of missing the “best” machine, and regret will build up during exploitation phase.



Some more advanced concepts and strategies

Some advanced questions

- **Should the exploration phase duration be decided beforehand or on-the-fly, based on what is seen during the exploration?**

If we manage to identify the “best” machine early, no point in exploring machines anymore, we should transition immediately into exploitation.

If we realize we are not confident about the “best” machine to use at the end of the fixed exploration phase, maybe better to do another round of exploration?

Or even better, smoothly transition between exploration and exploitation.

Exploring and exploiting (strategy #3)

Definition (ϵ -greedy strategy):

In the ϵ -greedy strategy, we define a value $\epsilon \in [0,1]$, which decreases over time. At each coin t , we decide on exploration/exploitation as:

$$\phi_t = \begin{cases} \text{exploration} & \text{with prob. } \epsilon \\ \text{exploitation} & \text{with prob. } 1 - \epsilon \end{cases}$$

- **Exploration:** choose action among all the random machines, with softmax distribution over the estimates $e_i(t)$.

Exploring and exploiting (strategy #3)

Definition (ϵ -greedy strategy):

In the ϵ -greedy strategy, we define a value $\epsilon \in [0,1]$, which decreases over time. At each coin t , we decide on exploration/exploitation as:

$$\phi_t = \begin{cases} \text{exploration} & \text{with prob. } \epsilon \\ \text{exploitation} & \text{with prob. } 1 - \epsilon \end{cases}$$

- **Exploitation:** use coin on the machine with the current highest estimate $e_i(t)$.

Exploring and exploiting (strategy #3)

Definition (ϵ -greedy strategy):

In the ϵ -greedy strategy, we define a value $\epsilon \in [0,1]$, which decreases over time. At each coin t , we decide on exploration/exploitation as:

$$\phi_t = \begin{cases} \text{exploration} & \text{with prob. } \epsilon \\ \text{exploitation} & \text{with prob. } 1 - \epsilon \end{cases}$$

- **Decay on ϵ over time:** the value of ϵ starts with 1, and progressively decays over time, progressively becoming 0.

Exploring and exploiting (strategy #4)

Definition (Upper Confidence Bound strategy):

In the **Upper Confidence Bound strategy**, we do not use the estimates of the expectation of the machine $e_i(t)$.

$$\forall i \in A, e_i(t) = \frac{\sum_{k=1}^t \delta_{a_t=i} R_t}{\sum_{k=1}^t \delta_{a_t=i}}$$

Exploring and exploiting (strategy #4)

Definition (**Upper Confidence Bound strategy**):

In the **Upper Confidence Bound strategy**, we do not use the estimates of the expectation of the machine $e_i(t)$.

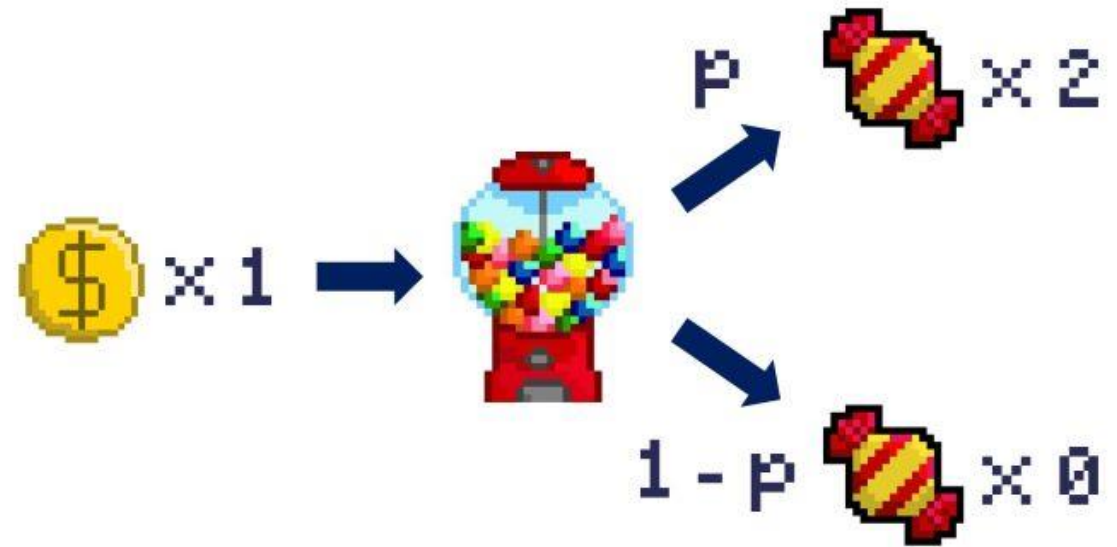
Instead, we choose the machine i , with the highest upper bound for the confidence interval of the statistical evaluation for machine i .

$$\forall t, a_t = \arg \max_{i \in A} \left[e_i(t) + c \sqrt{\frac{\log(t)}{N_t(i)}} \right]$$

With c a parameter defining the amount of exploration and $N_t(i)$ the amount of time we have played machine i , before time t .

Non-stationarity in problems

- So far, we have assumed that all parameters, describing the system (e.g. the hidden probabilities of machines) were **stationary**. This means that they do not change over time.
- But... What if they did?
- **Example:** machines increase their win probability when you lose and vice-versa.



→ Would then have to learn the dynamics of the probabilities (how they vary) on top of everything else!

Not

12-2012

Examining Slot Machine Play with Varying Percentages of Losses Disguised as Wins

Johnna L. Dunning
jld13731@siu.edu

- So far, parameters of the system produced statistically significant results.
- But, the results are not as clear as they seem.
- Example: the results show that the system is not as effective as it appears to be.



earn the
ties
:



Conclusion

1. What is Reinforcement Learning?
2. What are the key ideas behind reinforcement learning and its framework?
3. What is the exploration vs. exploitation tradeoff?
4. How do we train an RL agent by exploring, then progressively exploiting?
5. What are some advanced strategies in multi-arm bandit problems?
6. What are the Q and V functions for a RL problem?
7. What is Q-learning and how can it be implemented in RL problems?

Learn more about these topics

Out of class, for those of you who are curious

- [TheBibleOfRL] R. **Sutton** et al., “Reinforcement learning: An Introduction, 2nd edition”, 2018.