# ISLab

**Universidade do Minho**
Escola de Engenharia
Departamento de Informática

# Mestrado Integrado em Engenharia Informática
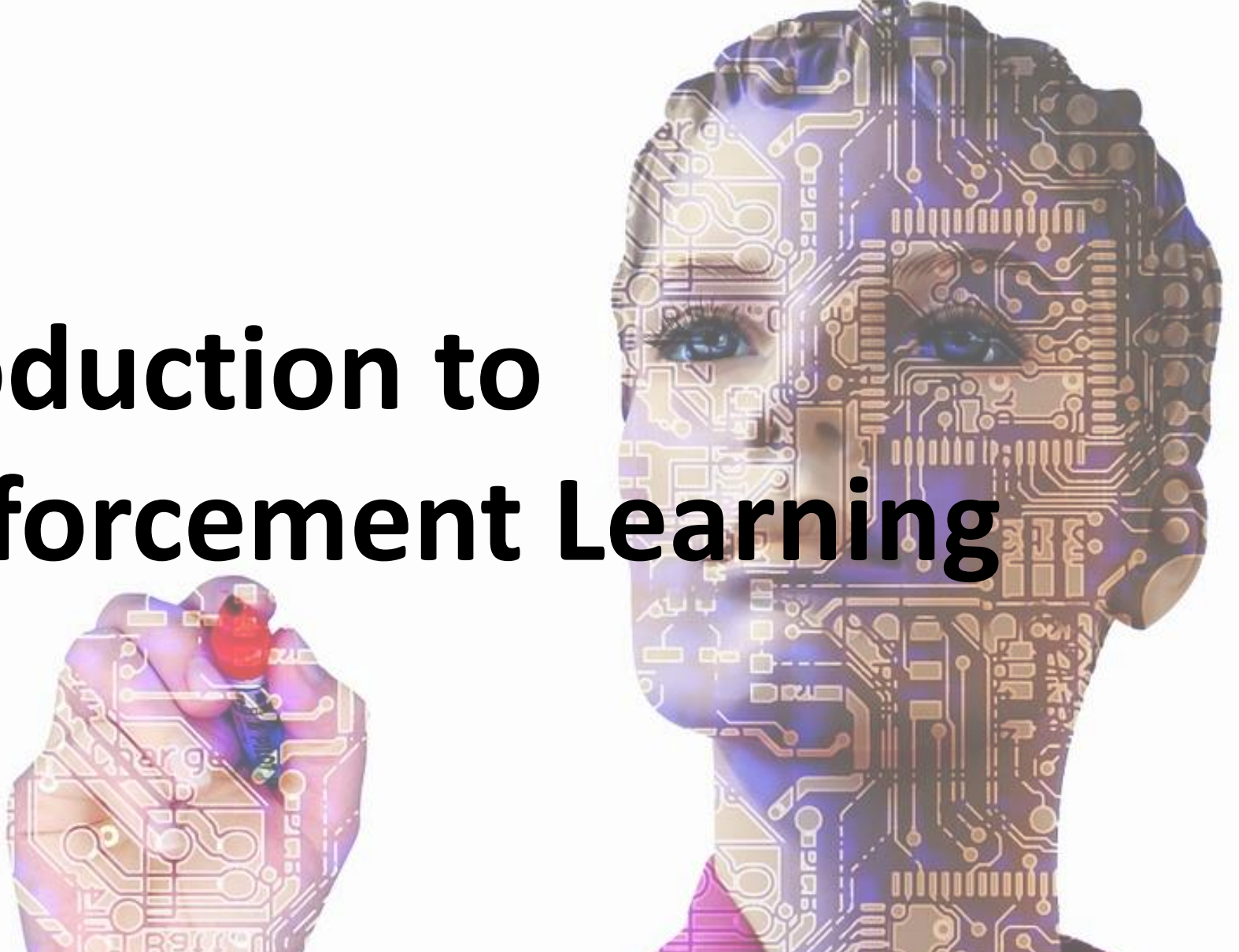# Computação Natural
# 2018/2019

Paulo Novais, Cesar Analide, Filipe Gonçalves

- Paulo Novais – pjon@di.uminho.pt

- Cesar Analide  –  cesar.analide@di.uminho.pt

- Filipe Gonçalves – fgoncalves@algoritmi.uminho.pt


- Departamento de Informática
  Escola de Engenharia
  Universidade do Minho

- Grupo ISLab – (Synthetic Intelligence Lab)

- Centro ALGORITMI
  Universidade do Minho

# Introduction to Reinforcement Learning

ISLab
Synthetic Intelligence Lab

- The Multi-armed Bandit Problem (Exploration vs Exploitation Problem)



D1          D2          D3          D4          D5

ISLab
Synthetic Intelligence Lab

- The Multi-armed Bandit Problem (Exploration vs Exploitation Problem)



D1    D2    D3    D4    D5

ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection (Exploration vs Exploitation Problem)


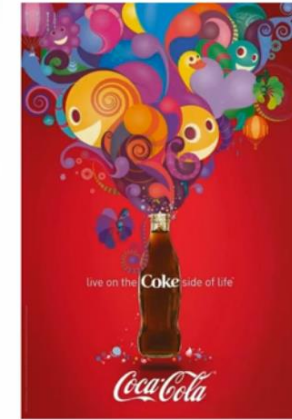
D1     D2     D3     D4     D5

ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection (Exploration vs Exploitation Problem)
    - We have d arms. For example, arms are ads that we display to users each time they connect to a web page
    - Each time a user connects to this web page, that makes a round
    - At each round n, we choose one ad to display to the user
    - At each round n, ad i gives reward:
        $r_i(n) \in \{0, 1\}$: $r_i(n) = 1$ if the user clicked on the ad i, 0 if the user didn't
    - Goal: maximize the total reward we get over many rounds

ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection
  - o Index: Person
  - o Column: Ads
    - 1: Person clicked on Ad
    - 0: Person ignored Ad

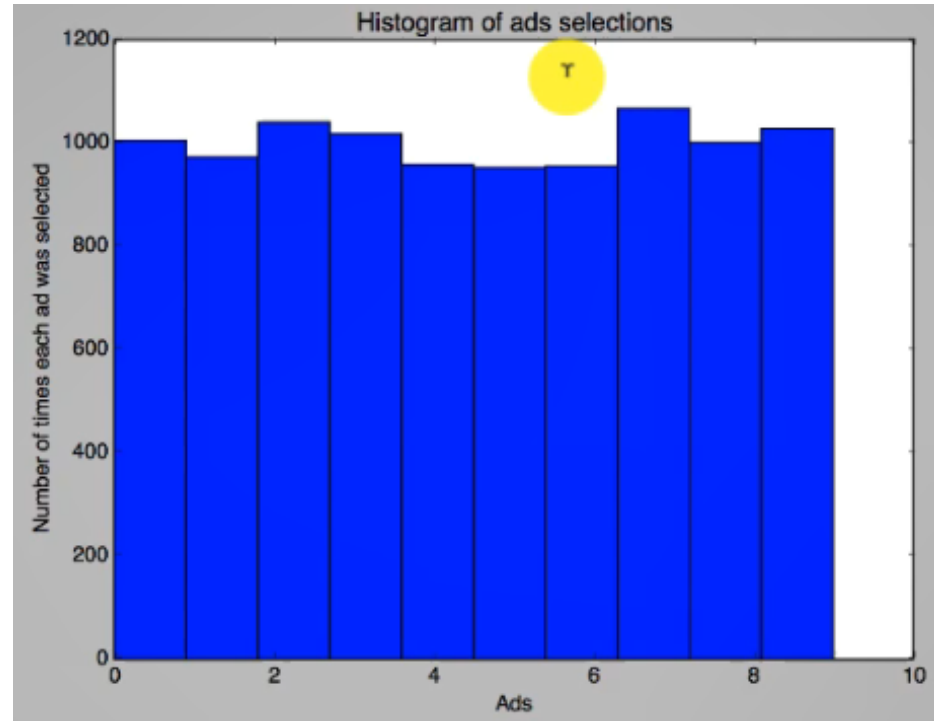| Index | Ad 1 | Ad 2 | Ad 3 | Ad 4 | Ad 5 | Ad 6 | Ad 7 | Ad 8 | Ad 9 | Ad 10 |
|-------|------|------|------|------|------|------|------|------|------|-------|
| 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 15 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 20 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 21 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |

ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection
  - ○ Random Selection

```python
1  # Random Selection
2
3  # Importing the libraries
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import pandas as pd
7
8  # Importing the dataset
9  dataset = pd.read_csv('Ads_CTR_Optimisation.csv')
10
11 # Implementing Random Selection
12 import random
13 N = 10000
14 d = 10
15 ads_selected = []
16 total_reward = 0
17 for n in range(0, N):
18     ad = random.randrange(d)
19     ads_selected.append(ad)
20     reward = dataset.values[n, ad]
21     total_reward = total_reward + reward
22
23 # Visualising the results — Histogram
24 plt.hist(ads_selected)
25 plt.title('Histogram of ads selections')
26 plt.xlabel('Ads')
27 plt.ylabel('Number of times each ad was selected')
28 plt.show()
```

ISLab
Synthetic Intelligence Lab

- **Marketing - Ads Selection**
  - ○ Random Selection

| i | Type | Size | Value |
|---|------|------|-------|
| 0 | int | 1 | 4 |
| 1 | int | 1 | 6 |
| 2 | int | 1 | 1 |
| 3 | int | 1 | 0 |
| 4 | int | 1 | 4 |
| 5 | int | 1 | 5 |
| 6 | int | 1 | 0 |
| 7 | int | 1 | 8 |
| 8 | int | 1 | 5 |
| 9 | int | 1 | 3 |
| 10 | int | 1 | 8 |
| 11 | int | 1 | 5 |
| 12 | int | 1 | 8 |
| 13 | int | 1 | 4 |
| 14 | int | 1 | 0 |



Histogram of ads selections

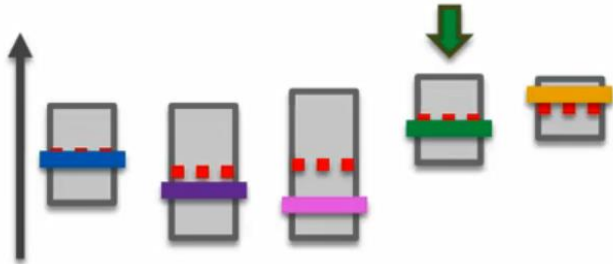| total_reward | int64 | 1 | 1193 |
|---|---|---|---|

# Upper Confidence Bound vs Thomson Sampling
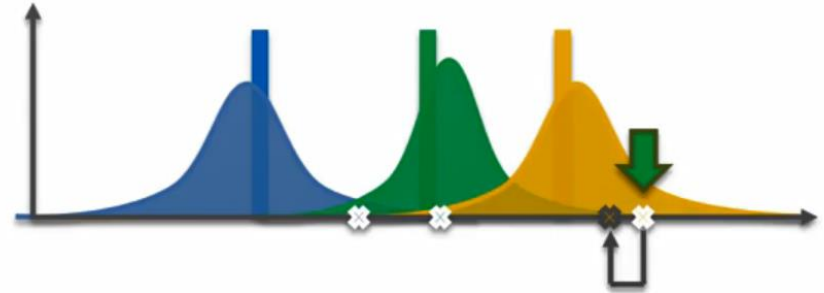
**UCB**
- Deterministic
- Requires update at every round

**Thompson Sampling**
- Probabilistic
- Can accommodate delayed feedback
- Better empirical evidence

- **Upper Confidence Bound Algorithm**

  **Step 1**. At each round $n$, we consider two numbers for each ad $i$:
  - $N_i(n)$ - the number of times the ad $i$ was selected up to round $n$,
  - $R_i(n)$ - the sum of rewards of the ad $i$ up to round $n$.

  **Step 2**. From these two numbers we compute:
  - the average reward of ad $i$ up to round $n$

  $$\bar{r}_i(n) = \frac{R_i(n)}{N_i(n)}$$

  - the confidence interval $[\bar{r}_i(n) - \Delta_i(n), \bar{r}_i(n) + \Delta_i(n)]$ at round $n$ with

  $$\Delta_i(n) = \sqrt{\frac{3}{2}\frac{\log(n)}{N_i(n)}}$$

  **Step 3**. We select the ad $i$ that has the maximum UCB $\bar{r}_i(n) + \Delta_i(n)$.
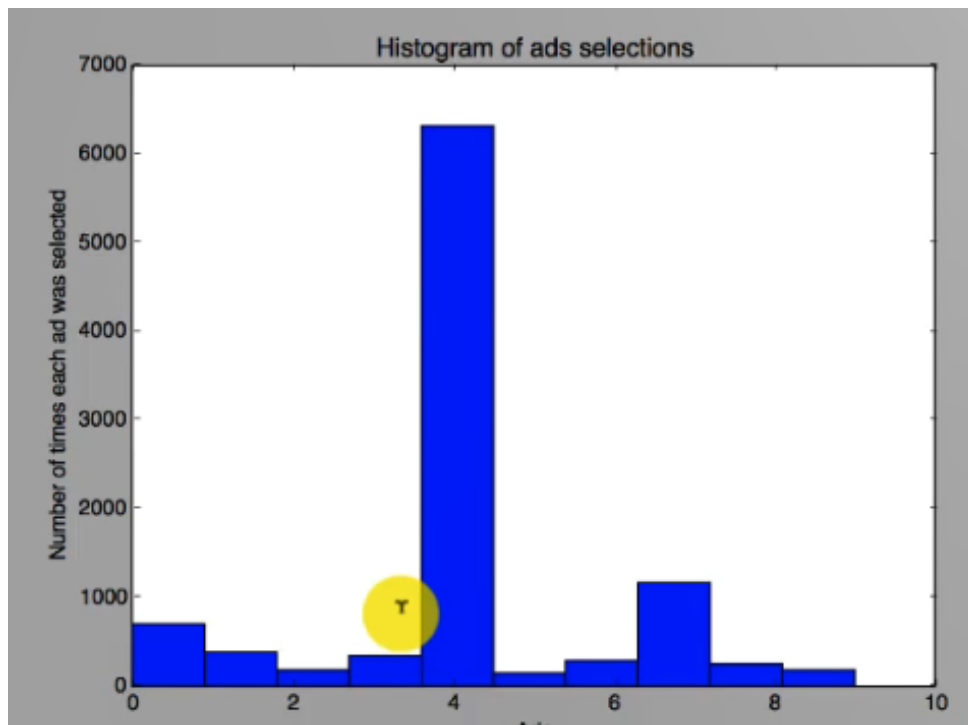
ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection
  - o Upper Confidence Bound

```
1 # Upper Confidence Bound
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Ads_CTR_Optimisation.csv')
10
11 # Implementing UCB
12 import math
13 N = 10000
14 d = 10
15 ads_selected = []
16 numbers_of_selections = [0] * d
17 sums_of_rewards = [0] * d
18 total_reward = 0
19 for n in range(0, N):
20     ad = 0
21     max_upper_bound = 0
22     for i in range(0, d):
23         if (numbers_of_selections[i] > 0):
24             average_reward = sums_of_rewards[i] / numbers_of_selections[i]
25             delta_i = math.sqrt(3/2 * math.log(n + 1) / numbers_of_selections[i])
26             upper_bound = average_reward + delta_i
27         else:
28             upper_bound = 1e400
29         if upper_bound > max_upper_bound:
30             max_upper_bound = upper_bound
31             ad = i
32     ads_selected.append(ad)
33     numbers_of_selections[ad] = numbers_of_selections[ad] + 1
34     reward = dataset.values[n, ad]
35     sums_of_rewards[ad] = sums_of_rewards[ad] + reward
36     total_reward = total_reward + reward
37
```

# ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection
  - o Upper Confidence Bound

| 1 ▲ | Type | Size | Value |
|---|---|---|---|
| 9985 | int | 1 | 4 |
| 9986 | int | 1 | 4 |
| 9987 | int | 1 | 4 |
| 9988 | int | 1 | 4 |
| 9989 | int | 1 | 4 |
| 9990 | int | 1 | 4 |
| 9991 | int | 1 | 4 |
| 9992 | int | 1 | 4 |
| 9993 | int | 1 | 4 |
| 9994 | int | 1 | 4 |
| 9995 | int | 1 | 4 |
| 9996 | int | 1 | 4 |
| 9997 | int | 1 | 4 |
| 9998 | int | 1 | 4 |
| 9999 | int | 1 | 4 |



Histogram of ads selections

| total_reward | int64 | 1 | 2178 |
|---|---|---|---|
| upper_bound | float64 | 1 | 0.31017236647899182 |

- Thompson Sampling Algorithm

    **Step 1**. At each round $n$, we consider two numbers for each ad $i$:

    - $N_i^1(n)$ - the number of times the ad $i$ got reward 1 up to round $n$,
    - $N_i^0(n)$ - the number of times the ad $i$ got reward 0 up to round $n$.

    **Step 2**. For each ad $i$, we take a random draw from the distribution below:

    $$\theta_i(n) = \beta(N_i^1(n) + 1, N_i^0(n) + 1)$$

    **Step 3**. We select the ad that has the highest $\theta_i(n)$.

- Thompson Sampling Algorithm
  - Ad $i$ gets rewards $\mathbf{y}$ from Bernoulli distribution $p(\mathbf{y}|\theta_i) \sim \mathcal{B}(\theta_i)$.
  - $\theta_i$ is unknown but we set its uncertainty by assuming it has a uniform distribution $p(\theta_i) \sim \mathcal{U}([0,1])$, which is the prior distribution.
  - Bayes Rule: we approach $\theta_i$ by the posterior distribution

$$\underbrace{p(\theta_i|\mathbf{y})}_{\text{posterior distribution}} = \frac{p(\mathbf{y}|\theta_i)p(\theta_i)}{\int p(\mathbf{y}|\theta_i)p(\theta_i)d\theta_i} \propto \underbrace{p(\mathbf{y}|\theta_i)}_{\text{likelihood function}} \times \underbrace{p(\theta_i)}_{\text{prior distribution}}$$

  - We get $p(\theta_i|\mathbf{y}) \sim \beta(\text{number of successes}+1, \text{number of failures}+1)$
  - At each round $n$ we take a random draw $\theta_i(n)$ from this posterior distribution $p(\theta_i|\mathbf{y})$, for each ad $i$.
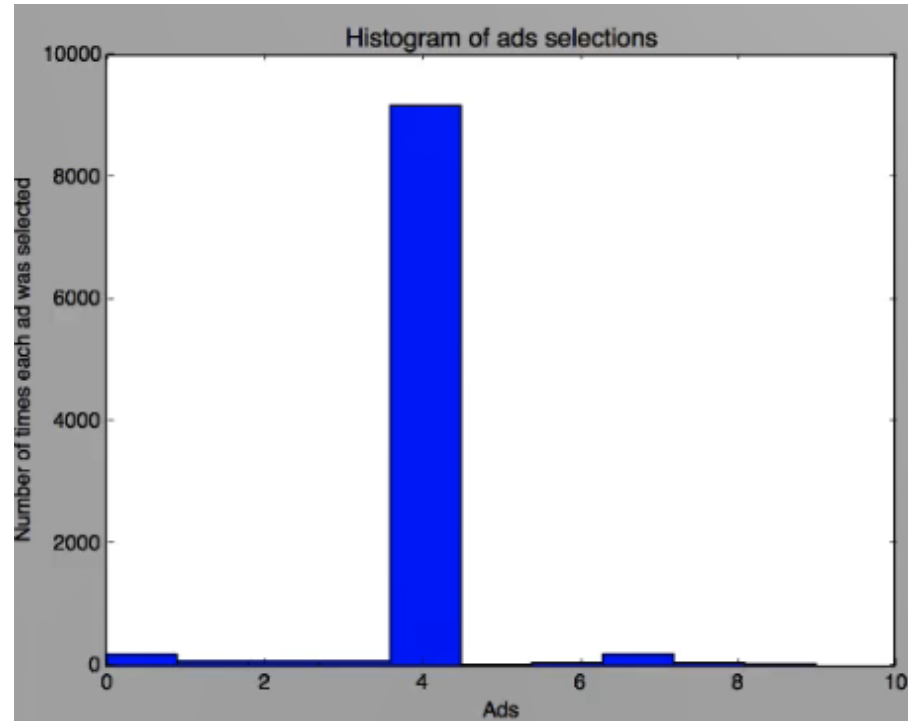  - At each round $n$ we select the ad $i$ that has the highest $\theta_i(n)$.

ISLab
Synthetic Intelligence Lab

- Marketing - Ads Selection
  - Thompson Sampling

```python
1  # Thompson Sampling
2
3  # Importing the libraries
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import pandas as pd
7
8  # Importing the dataset
9  dataset = pd.read_csv('Ads_CTR_Optimisation.csv')
10
11 # Implementing Thompson Sampling
12 import random
13 N = 10000
14 d = 10
15 ads_selected = []
16 numbers_of_rewards_1 = [0] * d
17 numbers_of_rewards_0 = [0] * d
18 total_reward = 0
19 for n in range(0, N):
20     ad = 0
21     max_random = 0
22     for i in range(0, d):
23         random_beta = random.betavariate(numbers_of_rewards_1[i] + 1, numbers_of_rewards_0[i] + 1)
24         if random_beta > max_random:
25             max_random = random_beta
26             ad = i
27     ads_selected.append(ad)
28     reward = dataset.values[n, ad]
29     if reward == 1:
30         numbers_of_rewards_1[ad] = numbers_of_rewards_1[ad] + 1
31     else:
32         numbers_of_rewards_0[ad] = numbers_of_rewards_0[ad] + 1
33     total_reward = total_reward + reward
34
```

- Marketing - Ads Selection
  - Thompson Sampling

**Universidade do Minho**
Escola de Engenharia
Departamento de Informática

# Mestrado Integrado em Engenharia Informática
# Computação Natural
# 2018/2019

## Paulo Novais, Cesar Analide, Filipe Gonçalves