

# RLD – Rapport TP 1: Bandits multi bras

Maxime DARRIN

17 septembre 2019

## 1 Baselines

```
1 def random_policy():
2     """
3     Return a random walk of the agent, taking uniformly each possible action
4     """
5     return [np.random.randint(0, 10) for i in range(5000)]
6
7 def staticbest_policy(click_rates):
8     """
9     Takes the action which maximises the score on that trajectory
10    """
11    a = np.argmax(np.sum(click_rates, axis=0))
12    return [a for i in range(5000)]
13
14
15 def opt_policy(click_rates):
16     """
17     At each timestep takes the best action
18     """
19    return np.argmax(click_rates, axis=1)
```

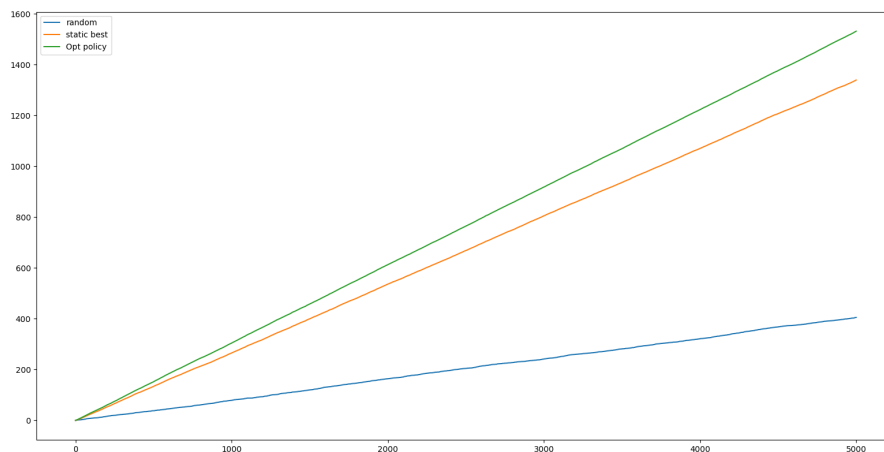


FIGURE 1 – Baselines – agents omniscients

## 2 UCB

```
1 def upper_bound(t, N, mu):
2     """
3     Computes the upper bound of the confidence interval of mean mu
4     """
5     return mu + np.sqrt(2 * np.log(t) / N)
6
7 def ucb_policy(click_rates):
8     """
9     Return the trajectory followed by the agent using ucb policy.
10    """
11
12    # Cumulative reward got by each actions
13    histo = np.array([click_rates[i][i] for i in range(10)])
14
15    # Number of times we took each action
16    counter = [1 for i in range(10)]
17
18    # List of the taken action
19    action_list = [i for i in range(10)]
20
21    for t in range(10, 5000):
22        action = np.argmax([upper_bound(t, counter[i], histo[i] / counter[i])
23                           for i in range(10)])
24        counter[action] += 1
25        histo[action] += click_rates[t][action]
26
27        action_list.append(action)
28    return action_list
```

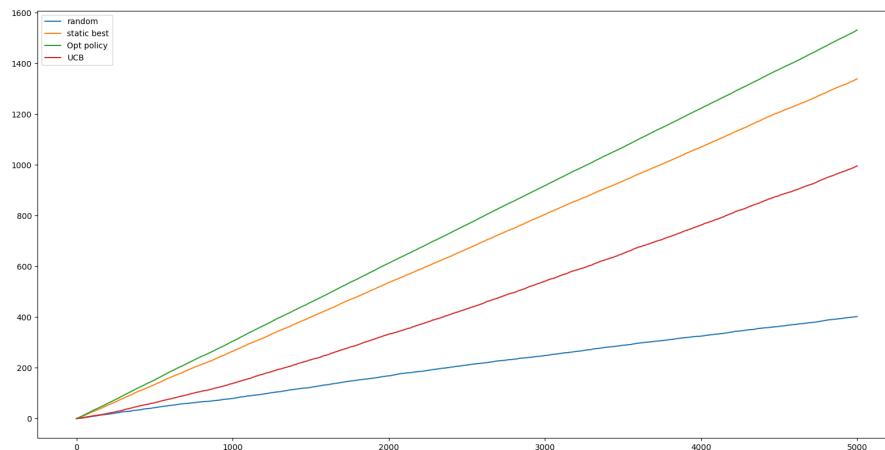


FIGURE 2 – Algorithmme UCB VS Baselines

## 3 LinUCB

```
1 def linucb_policy(alpha, articles, click_rates):
2     A = [np.identity(5, dtype=float) for i in range(10)]
```

```

3      b = [np.zeros((5, 1), dtype=float) for i in range(10)]
4
5      theta = [None for i in range(10)]
6      pt = [None for i in range(10)]
7
8      actions_list = []
9
10     for t in range(0, 5000):
11         for i in range(10):
12             theta[i] = np.dot(np.linalg.inv(A[i]), b[i])
13
14             pt[i] = (np.dot(np.transpose(theta[i]), articles[t]) + alpha * np.sqrt(
15                 np.dot(
16                     np.dot(np.transpose(articles[t]),
17                             np.linalg.inv(A[i])),
18                     articles[t]))) [0]
19
20             at = np.argmax(pt)
21             rt = click_rates[t][at]
22
23             A[at] = A[at] + np.dot(np.transpose(articles[t]), articles[t])
24             b[at] = b[at] + rt * articles[t]
25
26             actions_list.append(at)
27
28     return actions_list

```

Dans figure suivante on compare *Lin UCB* aux *baselines* en faisant varier  $\alpha$ . En particulier on vérifie bien qu'un  $\alpha$  bas correspond à une forte exploration (ici  $\alpha = 0.01$  donne des résultats proches de l'aléatoire) tandis que des valeurs plus importantes permettent un meilleur compromis exploration-exploitation.

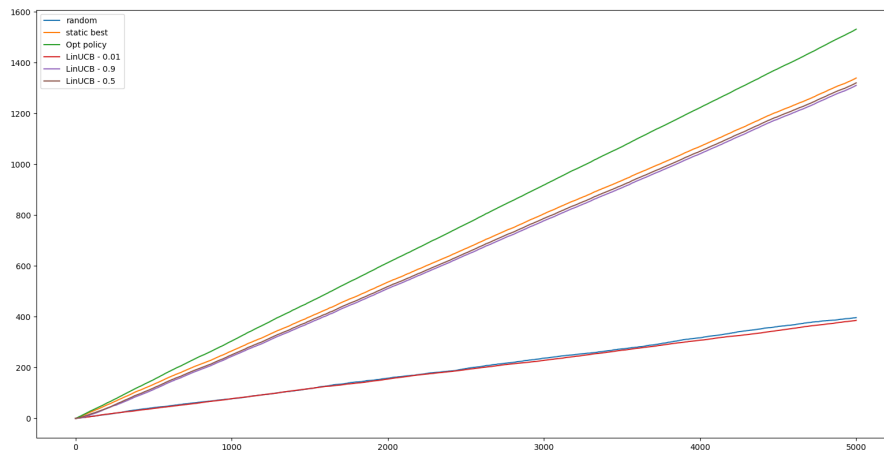


FIGURE 3 – Algorithme LinUCB VS Baselines

Par ailleurs, on observe bien que *lin UCB* est clairement meilleur que *UCB*. En effet, en utilisant le contexte pour prendre des décisions plus averties.

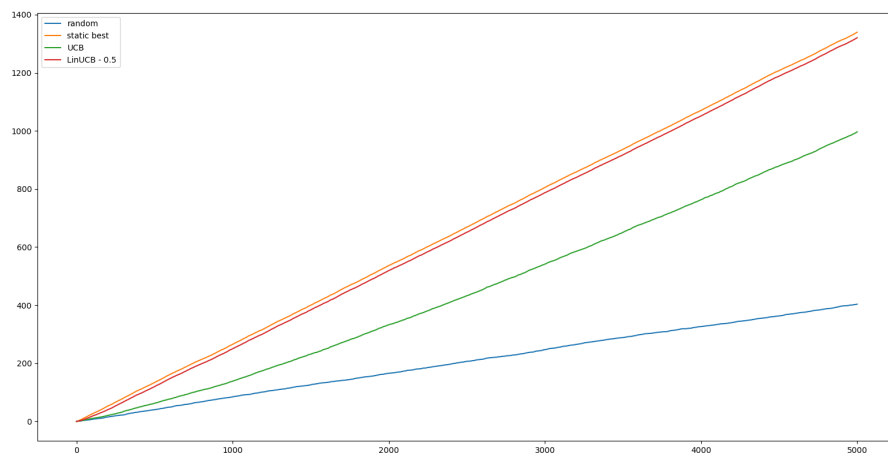


FIGURE 4 – Algorithmme LinUCB VS UCB