

LinUBC online implementation

Abhishek Bhatia

September 10, 2018

1 Problem statement and Motivation

The problem of personalized news article recommendation can be modeled as a multi-armed bandit problem with context information. Formally, a contextual-bandit algorithm A proceeds in discrete trials $t = 1, 2, 3, \dots$. In trial t :

1. The LinUSB algorithm observes the current user u_t and a set A_t of arms or actions together with their feature vectors $x_{t,a}$ for $a \in A_t$. This vector $x_{t,a}$ summarizes information of both u_t and arm a , and will be referred to as the context.
2. Based on observed payoffs in previous trials, the algorithm A chooses an arm $a_t \in A_t$, and consequently receives payoff $r_{t,a}$ whose expectation depends on both user u_t and the arm a_t .
3. The algorithm then improves its arm-selection technique with the new observation obtained, that is, $(x_{t,a}, a_t, r_{t,a_t})$.

When this problem is viewed in the context of article recommendation, articles are the arms. When the article that was presented is clicked, a reward of 1 is obtained; otherwise, the payoff is 0. The expected reward of an article is then its clickthrough rate (CTR), and by choosing an article with maximum CTR is the same as maximizing the total expected reward in previous mentioned bandit formulation.

2 Dataset

- The personalization dataset consists of $t = 1, \dots, 10000$ data points.

- The dataset contains the action(a_t) was performed at each timestep(t). The action $a_t \in \{1, \dots, 10\}$.
- For each time-step, the dataset contains the reward that was obtained $y_t \in \{0, 1\}$.
- There are 100 context features which is represented as a vector $x_t \in R^{100}$.

3 Algorithm

In Algorithm 1, we give the algorithm used. Here, we use a LinUCB with disjoint linear models which is modified to work in an online manner on the dataset given. We then step through the stream of logged events in the dataset one by one. If, given the current history h_{t-1} , it happens that the policy selects the same arm a as the one that was selected by the logging policy, then the corresponding event is retained. Otherwise, if the policy selects a different arm from the one that was taken by the logging policy, then the event is entirely ignored, and the algorithm proceeds to the next event without any other change in its state. Also, the context is same in our case for all actions.

Data: Inputs: $\alpha \in \mathbb{R}^+$, dataset

```

for  $t = 1, 2, 3 \dots T$  do
    Observe features of all arms  $a \in A_t$ :  $x_t, a \in \mathbb{R}^d$ ;
    for  $a$  in  $A_t$  do
        if  $a$  is new then
             $A_a \rightarrow I_d$  (d-dimensional identity matrix);
             $b_a \rightarrow 0_{d \times 1}$  (d-dimensional zero vector);
        end
         $\theta \rightarrow A_a^{-1} b_a$ ;
         $p_{t,a} \rightarrow \theta^T x_{t,a} + \alpha \sqrt{x_{t,a}^T A_a^{-1} x_{t,a}}$ ;
    end
    Choose arm  $a_t = \arg \max_{a \in A_t} p_{t,a}$  with ties broken arbitrarily;
    compute  $C(t)$  ;
    if  $a_t$  is  $a_t^{dataset}$  then
         $A_a \rightarrow A_{a_t} + x_{t,a_t} x_{t,a_t}^T$ ;
         $b_a \rightarrow b_{a_t} + r_t x_{t,a_t}$ ;
    end
end

```

Algorithm 1: linUSB

4 Metric for evaluation

In this report, we use cumulative take-rate of our actions for any time t as evaluation metric for the data-set. This is computed as go through the dataset.

Notation

π_{t-1} : algorithm trained on data upto time $t-1$.

$\pi_{t-1}(x_t)$: action that algorithm chooses for the context x_t which it observes at time t .

a_t : real action that taken in the data-set at time t .

y_t : real reward that was obtained at time t .

$$C(T) = \frac{\sum_{t=1}^T y_t \times \mathbf{1}[\pi_{t-1}(x_t) = a_t]}{\sum_{t=1}^T \mathbf{1}[\pi_{t-1}(x_t) = a_t]}$$

5 Experimental Setup and Optimizations

- We load the dataset and perform linUSB as given in 1 for a particular alpha strategy.
- Then we plot $C(t)$ for each time-step for each alpha strategy on a graph to compare.

We used the following alpha strategies:

- $\alpha(t) = \frac{1}{t}$
- $\alpha(t) = \frac{1}{\sqrt{t}}$
- $\alpha(t) = \frac{0.1}{\sqrt{t}}$
- $\alpha(t) = 0.1$
- $\alpha(t) = e^{-t}$

The strategy number 3 listed above was determined by doing a grid-search on the numerator value.

6 Results

6.1 Cumulative take-rate replay for different alpha strategies

In Table 1, we show the Cumulative take-rate replay($C(T)$) at last time step for different alpha strategies. We see that, the best $C(T)$ is 0.9458.

In the Figure 1, we plot $C(T)$ as time-steps progress for different alpha strategies. We notice $\alpha = \frac{0.1}{\sqrt{t}}$ achieves the best $C(T)$, followed by $\alpha = \frac{1}{\sqrt{t}}$ and $\alpha = \frac{1}{t}$. In comparison to these strategies, $\alpha = e^t$ and $\alpha = 0.1$ converge relatively slowly and achieve sub-optimal $C(T)$ s throughout and finally.

In addition, we can see that both $\alpha = \frac{1}{\sqrt{t}}$ and $\alpha = \frac{1}{t}$ achieve the same final $C(T)$ although $\alpha = \frac{1}{t}$ is noted to converge much faster.

7 Code

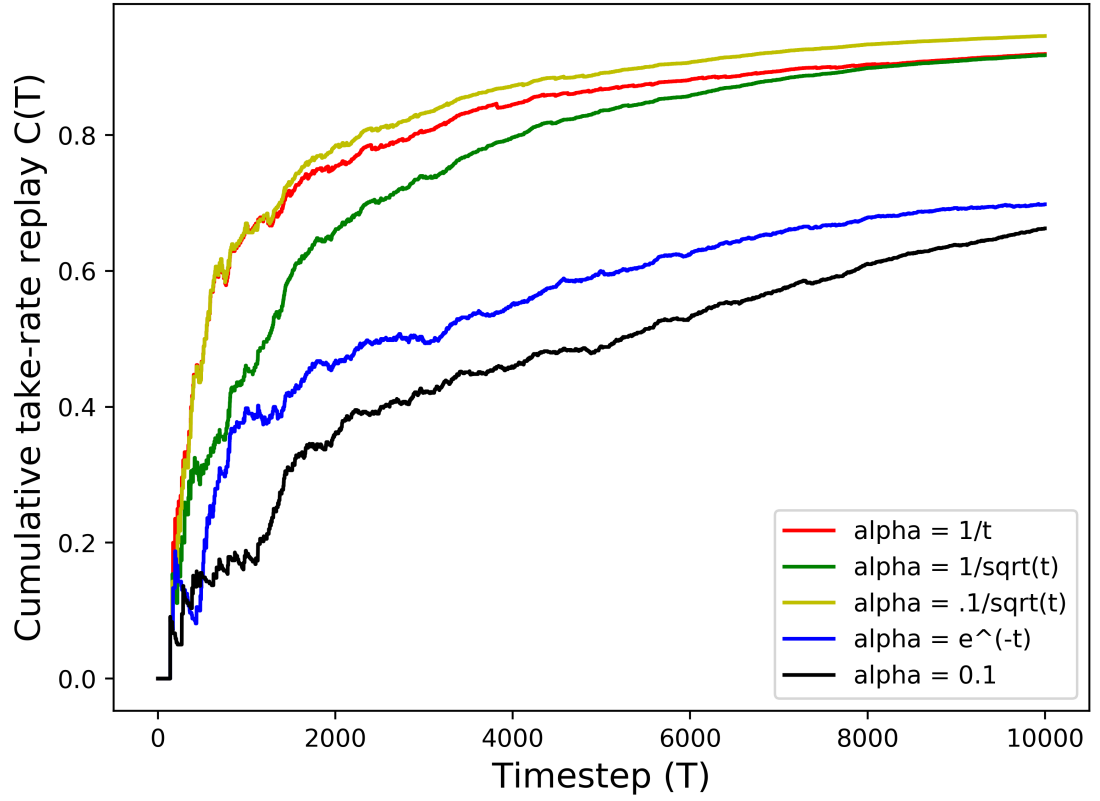


Figure 1: Cumulative take-rate replay with timesteps for different strategies of α

α	$C(T)$
$\frac{1}{t}$	0.919417
$\frac{1}{\sqrt{t}}$	0.917636
$\frac{0.1}{\sqrt{t}}$	0.945894
e^{-t}	0.698150
0.1	0.662687

Table 1: Final Cumulative take-rate replay for different strategies of α

```

1 import pandas as pd
2 import numpy as np
3 from numpy.linalg import inv
4 from math import sqrt
5 import ipdb
6 from tqdm import tqdm
7 import matplotlib.pyplot as plt
8 import math
9
10 def read_data(filename):
11     '''
12     Input: path to csv file
13     Output: numpy arrays of action, reward and context
14     Description: Read csv to dataframe. Parses the dataframe
15                 and converts to numpy arrays.
16     '''
17     df = pd.read_csv(filename, header=None, sep="\s+")
18
19     action = df.iloc[:,0].values
20     reward = df.iloc[:,1].values
21     x_ta = df.iloc[:,2:].values
22
23     return action, reward, x_ta
24
25
26 def initialize_params(d,num_actions):
27     '''
28     Input: time-steps and number of dimensions
29     Output: A_a and b_a
30     Description: Initialize A_a to identity matrix for each action.
31                 Initialize b_a to 0 d-dimensional vector
32     '''
33     A = [np.identity(d) for i in range(num_actions)]
34     b_as = [np.zeros(d) for i in range(num_actions)]
35
36     return A , b_as
37
38 def linUCB(action, reward, x_ta, A, b_as, alphas):
39     '''
40     Input: action, reward and context features for each timestep in the dataset.
41           Initialized values for A and b_as.
42           Alphas are each timestep for a particular alpha strategy.
43     Output: C(T) as an array of size Tx1, where T is the number of
44            time-steps in the dataset.
45     Description: perform a LinUCB with disjoint linear models which is
46                 modified to work in an online manner on the dataset given.
47                 And compute C(t) at each timestep t.
48     '''
49
50     c_t_num = []

```

```

51     c_t_den = []
52
53     for i in tqdm(range(time_steps)):
54         alpha = alphas[i]
55         x_ta_i = x_ta[i,:]
56         r_i = reward[i]
57         action_i = action[i]
58
59         p_t_a=[0.]*10
60         for a_i in range(10):
61             inv_a = inv(A[a_i])
62             theta = np.dot(inv_a,b_as[a_i])
63             p_t_a[a_i] = np.dot(theta.T,x_ta_i) + alpha * \
64                 sqrt(np.dot(np.dot(x_ta_i.T,inv_a),x_ta_i))
65
66         # choose the best action
67         a_t = np.argmax(p_t_a)
68
69         if i>=1:
70             # update
71             if a_t+1 == action_i:
72
73                 x_ta_i= x_ta_i.reshape(100,1)
74                 A[a_t] += np.dot(x_ta_i,x_ta_i.T)
75                 b_as[a_t] += r_i*x_ta_i.flatten()
76
77                 c_t_num.append(r_i)
78                 c_t_den.append(1.)
79             else:
80                 c_t_num.append(0.)
81                 c_t_den.append(0.)
82
83         c_t_num = np.cumsum(c_t_num)
84         c_t_den = np.cumsum(c_t_den)
85
86         np.seterr(all='ignore')
87         c_t = np.nan_to_num(np.divide(c_t_num,c_t_den, dtype=float))
88
89     return c_t
90
91 def plot_ct(c_t_1_t, c_t_sqrt,c_01sqrt, c_e_t, c_point01):
92     '''
93     Input: C(T) for alpha strategies
94     Description: Plot the C(T) for each alpha strategy as time-steps progress.
95                 The resulting plot is saved on disk
96     '''
97     p1, = plt.plot(range(1,len(c_t_1_t)+1),c_t_1_t,'r',
98                     label="alpha = 1/t")
99     p2, = plt.plot(range(1,len(c_t_sqrt)+1),c_t_sqrt,'g',
100                    label="alpha = 1/sqrt(t)")

```

```

101     p3, = plt.plot(range(1,len(c_01sqrt)+1),c_01sqrt,'y',
102                     label="alpha =  $.1/\sqrt{t}$ ")
103     p4, = plt.plot(range(1,len(c_e_t)+1),c_e_t,'b',
104                     label="alpha =  $e^{-t}$ ")
105     p5, = plt.plot(range(1,len(c_point01)+1),c_point01,'k',
106                     label="alpha = 0.1")
107
108     plt.legend(handles=[p1, p2, p3, p4, p5])
109
110
111     plt.xlabel("Timestep (T)",fontsize=14)
112     plt.ylabel("Cumulative take-rate replay C(T)",fontsize=14)
113     plt.tight_layout()
114     plt.savefig("plots", dpi=500)
115     plt.close()
116
117 # read the dataset
118 action, reward, x_ta = read_data("dataset.txt")
119
120 # determine dataset specific values
121 d = x_ta.shape[1]
122 time_steps = action.shape[0]
123 num_actions = 10
124
125 # use different alpha strategies
126 alphas_sqrt = [1./sqrt(i+1.0) for i in range(time_steps)]
127 alphas_01sqrt = [0.1/sqrt(i+1.0) for i in range(time_steps)]
128 alphas_1_t = [1./(i+1.0) for i in range(time_steps)]
129 alphas_e_t = [math.exp(-i) for i in range(time_steps)]
130 alphas_point01 = [0.1 for i in range(time_steps)]
131
132 # run linUCB for different alphas
133 A_init , b_as_init = initialize_params(d,num_actions)
134 c_t_1_t = linUCB(action, reward, x_ta, A_init , b_as_init, alphas_1_t)
135 print("Final C(T) %s %f " %("alpha = 1/t", c_t_1_t[-1]))
136
137 A_init , b_as_init = initialize_params(d,num_actions)
138 c_t_sqrt = linUCB(action, reward, x_ta, A_init , b_as_init, alphas_sqrt)
139 print("Final C(T) %s %f " %("alpha =  $1/\sqrt{t}$ ", c_t_sqrt[-1]))
140
141 A_init , b_as_init = initialize_params(d,num_actions)
142 c_e_t = linUCB(action, reward, x_ta, A_init , b_as_init, alphas_e_t)
143 print("Final C(T) %s %f " %("alpha =  $e^{-t}$ ", c_e_t[-1]))
144
145 A_init , b_as_init = initialize_params(d,num_actions)
146 c_point01 = linUCB(action, reward, x_ta, A_init , b_as_init, alphas_point01)
147 print("Final C(T) %s %f " %("alpha = 0.1", c_point01[-1]))
148
149
150

```



```

151 | A_init , b_as_init = initialize_params(d,num_actions)
152 | c_01sqrt = linUCB(action, reward, x_ta, A_init , b_as_init, alphas_01sqrt)
153 | print("Final C(T) %s %f " %("alpha = .1/sqrt(t)", c_01sqrt[-1]))
154 |
155 | # plot the result for C(T)
156 | plot_ct(c_t_1_t, c_t_sqrt,c_01sqrt, c_e_t, c_point01)

```