

All gists GitHub

New gist



anonymous / bandit\_simulations.py

Created 3 years ago

bandit\_simulations.py

```
1  import numpy as np
2  from matplotlib import pylab as plt
3  #from mpltools import style # uncomment for prettier plots
4  #style.use(['ggplot'])
5
6  '''
7  function definitions
8  '''
9
10 # generate all bernoulli rewards ahead of time
11 def generate_bernoulli_bandit_data(num_samples,K):
12     CTRs_that_generated_data = np.tile(np.random.rand(K), (num_samples,1))
13     true_rewards = np.random.rand(num_samples,K) < CTRs_that_generated_data
14     return true_rewards,CTRs_that_generated_data
15
16 # totally random
17 def random(estimated_beta_params):
18     return np.random.randint(0,len(estimated_beta_params))
19
20 # the naive algorithm
21 def naive(estimated_beta_params,number_to_explore=100):
22     totals = estimated_beta_params.sum(1) # totals
23     if np.any(totals < number_to_explore): # if have been explored less than specified
24         least_explored = np.argmin(totals) # return the one least explored
25         return least_explored
26     else: # return the best mean forever
27         successes = estimated_beta_params[:,0] # successes
28         estimated_means = successes/totals # the current means
29         best_mean = np.argmax(estimated_means) # the best mean
30         return best_mean
31
32 # the epsilon greedy algorithm
33 def epsilon_greedy(estimated_beta_params,epsilon=0.01):
34     totals = estimated_beta_params.sum(1) # totals
35     successes = estimated_beta_params[:,0] # successes
36     estimated_means = successes/totals # the current means
37     best_mean = np.argmax(estimated_means) # the best mean
38     be_exporatory = np.random.rand() < epsilon # should we explore?
39     if be_exporatory: # totally random, excluding the best_mean
40         other_choice = np.random.randint(0,len(estimated_beta_params))
41         while other_choice == best_mean:
42             other_choice = np.random.randint(0,len(estimated_beta_params))
43         return other_choice
44     else: # take the best mean
45         return best_mean
46
47 # the UCB algorithm using
48 # (1 - 1/t) confidence interval using Chernoff-Hoeffding bound)
49 # for details of this particular confidence bound, see the UCB1-TUNED section, slide 18, of:
50 # http://lane.compbio.cmu.edu/courses/slides_ucb.pdf
51 def UCB(estimated_beta_params):
52     t = float(estimated_beta_params.sum()) # total number of rounds so far
53     totals = estimated_beta_params.sum(1)
54     successes = estimated_beta_params[:,0]
55     estimated_means = successes/totals # sample mean
56     estimated_variances = estimated_means - estimated_means**2
57     UCB = estimated_means + np.sqrt( np.minimum( estimated_variances + np.sqrt(2*np.log(t)/totals), 0.25 ) * np.log(t)/totals )
58     return np.argmax(UCB)
59
60 # the UCB algorithm - using fixed 95% confidence intervals
61 # see slide 8 for details:
62 # http://dept.stat.lsa.umich.edu/~kshedden/Courses/Stat485/Notes/binomial_confidence_intervals.pdf
63 def UCB_bernoulli(estimated_beta_params):
64     totals = estimated_beta_params.sum(1) # totals
65     successes = estimated_beta_params[:,0] # successes
```

```

65     estimated_means = successes/totals # sample mean
66     estimated_variances = estimated_means - estimated_means**2
67     UCB = estimated_means + 1.96*np.sqrt(estimated_variances/totals)
68     return np.argmax(UCB)
69
70
71 # the bandit algorithm
72 def run_bandit_dynamic_alg(true_rewards,CTRs_that_generated_data,choice_func):
73     num_samples,K = true_rewards.shape
74     # seed the estimated params (to avoid )
75     prior_a = 1. # aka successes
76     prior_b = 1. # aka failures
77     estimated_beta_params = np.zeros((K,2))
78     estimated_beta_params[:,0] += prior_a # allocating the initial conditions
79     estimated_beta_params[:,1] += prior_b
80     regret = np.zeros(num_samples) # one for each of the 3 algorithms
81
82     for i in range(0,num_samples):
83         # pulling a lever & updating estimated_beta_params
84         this_choice = choice_func(estimated_beta_params)
85
86         # update parameters
87         if true_rewards[i,this_choice] == 1:
88             update_ind = 0
89         else:
90             update_ind = 1
91
92         estimated_beta_params[this_choice,update_ind] += 1
93
94         # updated expected regret
95         regret[i] = np.max(CTRs_that_generated_data[i,:]) - CTRs_that_generated_data[i,this_choice]
96
97     cum_regret = np.cumsum(regret)
98
99     return cum_regret
100
101 '''
102 main code
103 '''
104 # define number of samples and number of choices
105 num_samples = 10000
106 K = 5 # number of arms
107 number_experiments = 100
108
109 regret_accumulator = np.zeros((num_samples,5))
110 for i in range(number_experiments):
111     print "Running experiment:", i+1
112     true_rewards,CTRs_that_generated_data = generate_bernoulli_bandit_data(num_samples,K)
113     regret_accumulator[:,0] += run_bandit_dynamic_alg(true_rewards,CTRs_that_generated_data,random)
114     regret_accumulator[:,1] += run_bandit_dynamic_alg(true_rewards,CTRs_that_generated_data,naive)
115     regret_accumulator[:,2] += run_bandit_dynamic_alg(true_rewards,CTRs_that_generated_data,epsilon_greedy)
116     regret_accumulator[:,3] += run_bandit_dynamic_alg(true_rewards,CTRs_that_generated_data,UCB)
117     regret_accumulator[:,4] += run_bandit_dynamic_alg(true_rewards,CTRs_that_generated_data,UCB_bernoulli)
118
119 plt.semilogy(regret_accumulator/number_experiments)
120 plt.title('Simulated Bandit Performance for K = 5')
121 plt.ylabel('Cumulative Expected Regret')
122 plt.xlabel('Round Index')
123 plt.legend(('Random', 'Naive', 'Epsilon-Greedy', '(1 - 1/t) UCB', '95% UCB'),loc='lower right')
124 plt.show()

```