

ETC algorithm

Explore-then-commit in Python

June 22, 2022

Explore-then-commit

1. **Input** m .
2. In round t choose action

$$A_t = \begin{cases} (t \bmod k) + 1 & \text{if } t \leq mk \\ \operatorname{argmax}_i \hat{\mu}_i(mk) & \text{if } t > mk \end{cases}$$

Code

```
1 def get_reward(rew_avg) -> np.ndarray:  
2     # Add epsilon (sub-gaussian noise) to reward.  
3     mean = np.zeros(rew_avg.size)  
4     cov = np.eye(rew_avg.size)  
5     epsilon = np.random.multivariate_normal(mean, cov)  
6     reward = rew_avg + epsilon  
7  
8     return reward
```

均值 $E(X_i) = 0, \quad \forall 1 \leq i \leq n$

协方差 $\text{cov}(X, Y) = 0$

方差 $D(X_i) = 1, \quad \forall 1 \leq i \leq n$

随机生成服从 1-subgaussian 的 epsilon, reward 也服从 1-subgaussian

在 *run_algo* 函数中一些变量初始化:

```
1 regret = np.zeros((num_trial, num_iter))
2 k = rew_avg.size
3 max_arm = np.argmax(rew_avg)
```

对于每一次 trial 中变量初始化

```
1 means = np.zeros(rew_avg.size)
2 num = np.zeros(rew_avg.size)
3 cum = [0]
```

means 记录每一个 arm 的均值, num 记录每一个 arm 被选择的次数 (方便算均值), cum 记录每一次增加后 regret 的累积量

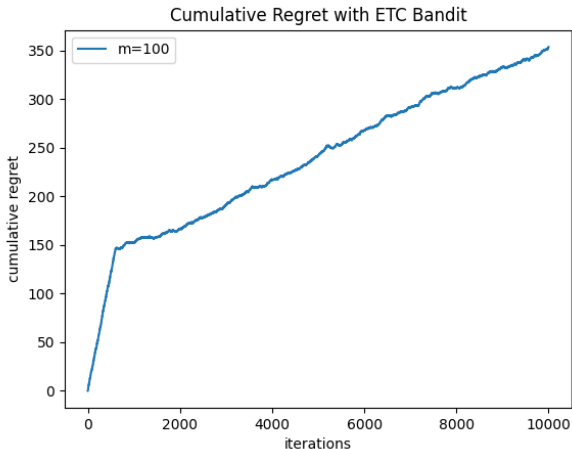
Code

```
1  for i in range(num_trial):
2      for t in range(num_iter - 1):
3          rew = get_reward(rew_avg)
4          if t <= m * k:
5              chosen_arm = t % k
6              num[chosen_arm] += 1
7              means[chosen_arm] += (rew[chosen_arm] - means[
                  chosen_arm]) / num[chosen_arm]
8          else:
9              chosen_arm = np.argmax(means)
10
11         reg = rew[max_arm] - rew[chosen_arm]
12         reg += cum[-1]
13         cum.append(reg)
14         regret[i, :] = np.asarray(cum)
15  return regret
```

绘制 cumulative regret 随着轮数增加的变化趋势的折线图

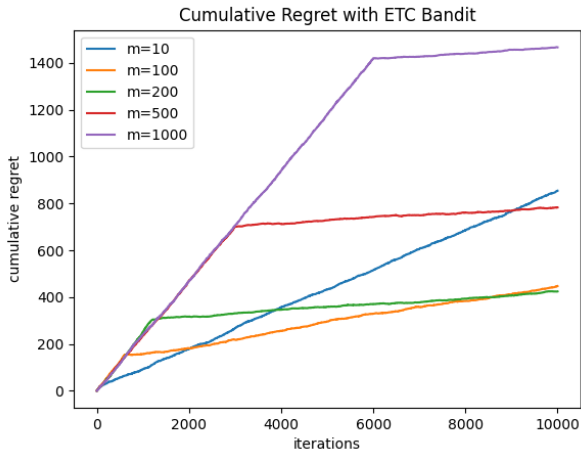
```
1 for m_val, reg in zip(m, regrets):
2     plt.plot(reg, label="m=" + str(m_val))
3
4 plt.xlabel('iterations')
5 plt.ylabel('cumulative regret')
6 plt.title('Cumulative Regret with ETC Bandit')
7 plt.legend()
8 plt.show()
```

Results



在前 m_k 轮次中寻找收益最大的 arm，在 m_k 轮次后 regret 的增加量趋于稳定 (始终选择收益最大的 arm)。

Results



在 m 增加的过程中，能够找到收益最大的 arm 的几率更高，在前 mk 轮次中所付出的代价会越高，但在 mk 轮次后所付出的代价会降低。

如何达到一个平衡？选取合适的 m ,

Choose m

$$m = \max \left\{ 1, \left\lceil \frac{4}{\Delta^2} \log \left(\frac{n\Delta^2}{4} \right) \right\rceil \right\}$$

Results

