

# Multiarmed bandits

---

## Thompson Sampling team

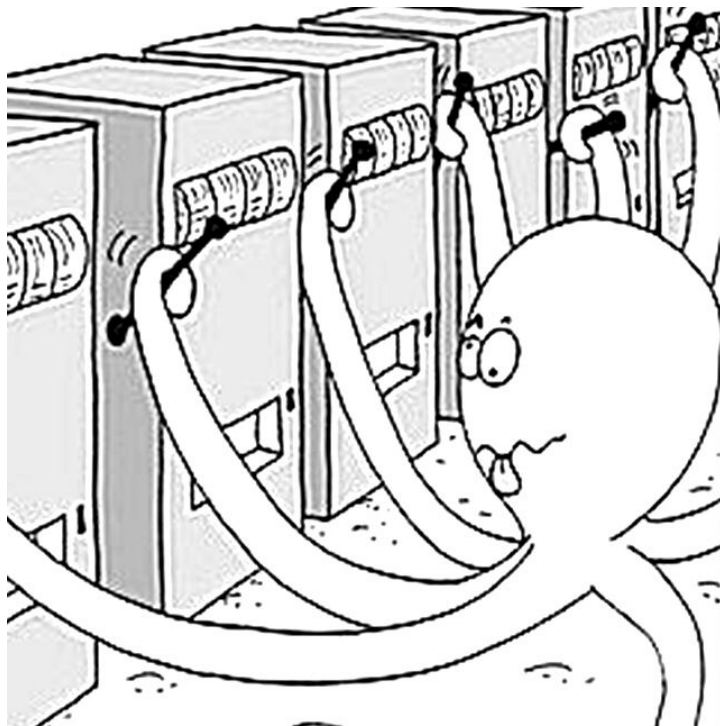
# Общая постановка задачи многоруких бандитов

Агент взаимодействует с окружением, выбирая действия (руки) из множества доступных

Каждое действие приносит случайную награду с неизвестным распределением

Цель агента – максимизировать суммарную награду за определённое количество шагов, балансируя между исследованием (сбор информации о распределениях) и эксплуатацией (использование текущих знаний для максимизации награды)

Основная метрика – regret (разница между наградой оптимального действия и выбранного действия)



# Неконтекстные бандиты

## Особенности:

- > Нет дополнительной информации (контекста) перед выбором действия
- > Каждое действие имеет фиксированное, но неизвестное распределение наград

## Примеры алгоритмов:

- >  $\epsilon$ -greedy
- > UCB
- > Thompson Sampling

## Формально:

- > На каждом шаге  $t$  агент выбирает действие  $a_t$  из  $N$  вариантов
- > Награда  $r_t$  выбирается из распределения с математическим ожиданием  $\mu_{a_t}$
- > Оптимальное действие:  $a^* = \operatorname{argmax}_a \mu_a$
- > Regret на шаге  $t$ :  $\Delta_t = \mu_{a^*} - \mu_{a_t}$
- > Цель: минимизировать  $R(T) = \sum \Delta_t$ ,  $t \in [1, T]$

# Контекстные бандиты

## Особенности:

- > Перед выбором действия агент получает контекст – вектор признаков, влияющий на награду
- > Распределение наград зависит от контекста (например, линейно)
- > Контекст может быть адаптивным (зависит от предыдущих действий агента)

## Примеры алгоритмов:

- > LinUCB
- > LinTS

## Формально:

- > На шаге  $t$ :

Агент получает контекстные векторы  $\{b_i(t)\} \in R_d$ ,  $i \in [1, N]$  для каждого действия

Ожидаемая награда действия  $i$ :  $E[r_i(t)] = b_i(t)^T \mu$ , где  $\mu \in R_d$  – неизвестный параметр

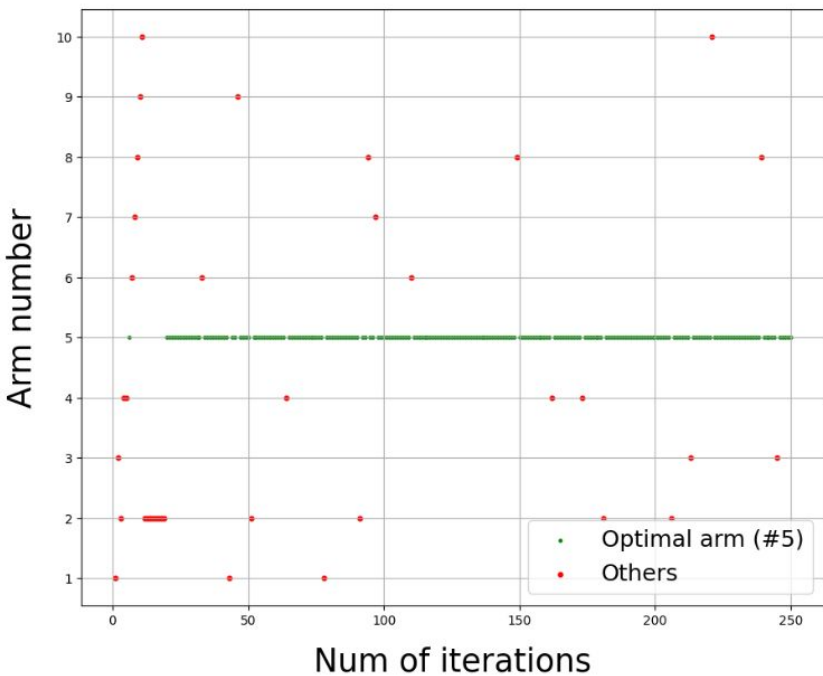
Агент выбирает действие  $a_t$ , получает награду  $r_{a_t}(t)$

- > Regret:  $\Delta_t = \max_i(b_i(t)^T \mu) - b_{a_t}(t)^T \mu$
- > Цель: минимизировать  $R(T) = \sum \Delta_t$ ,  $t \in [1, T]$

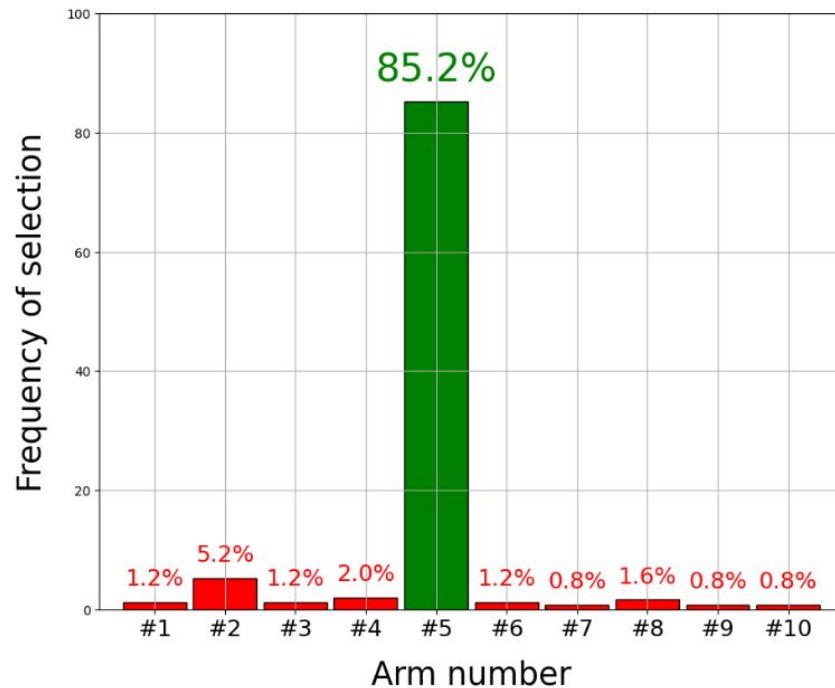
# **Сравнение алгоритмов**

# Epsilon-Greedy-Constant

Selection at each iteration

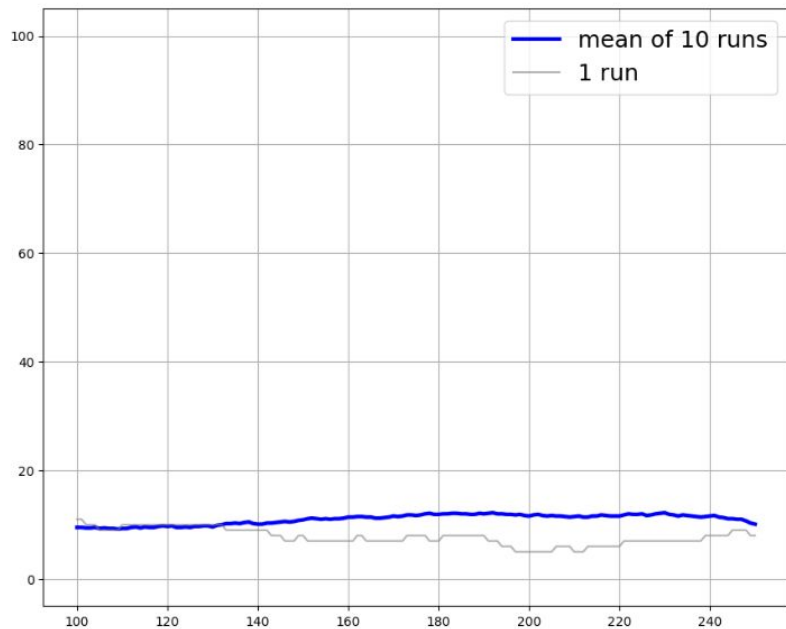


Arms selection frequency histogram at  $\epsilon = 0.13$

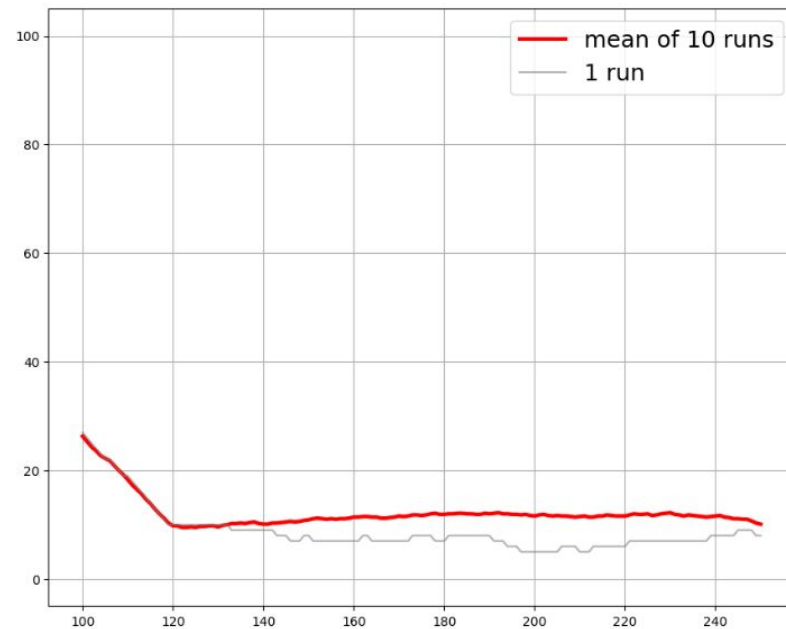


# Epsilon-Greedy-Constant

Exploration rate through past 100 iterations

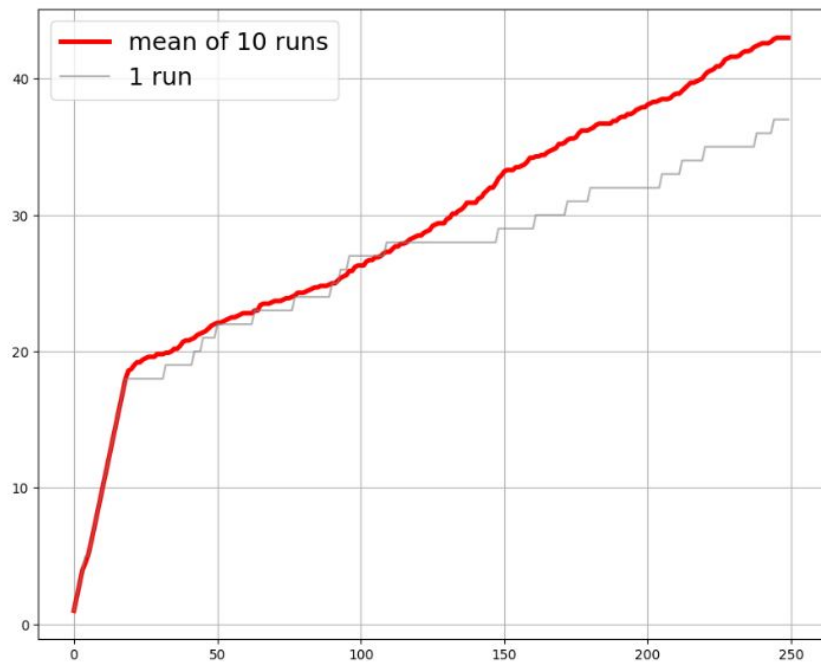


Num of unoptimal arms through past 100 iterations

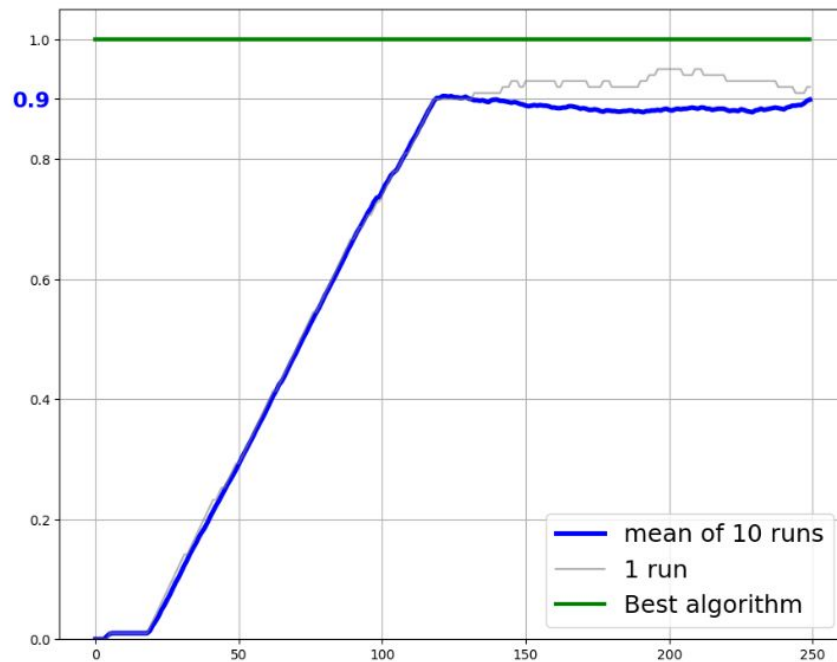


# Epsilon-Greedy-Constant

Cumulative regret



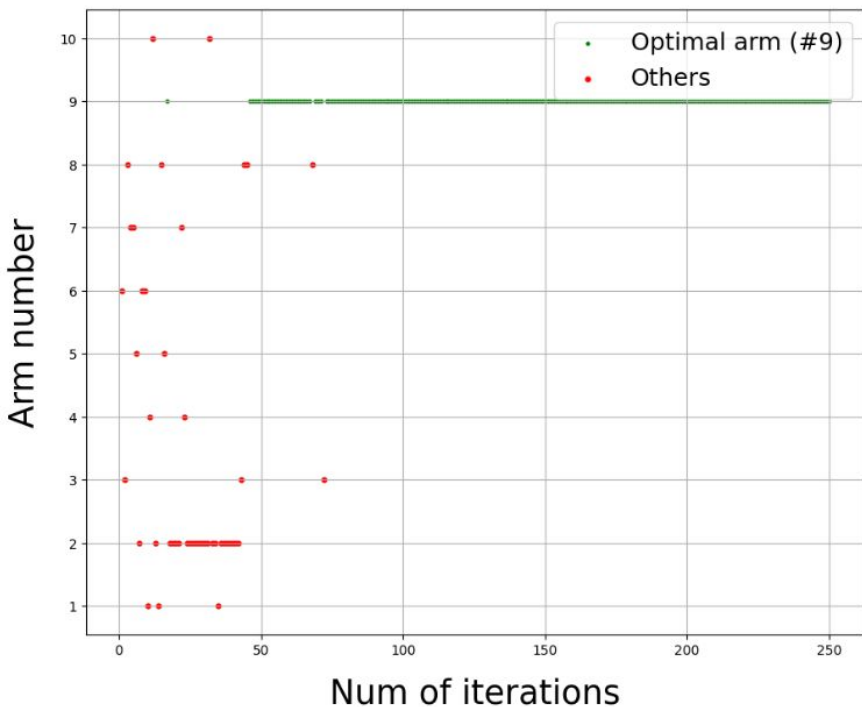
Convergence rate



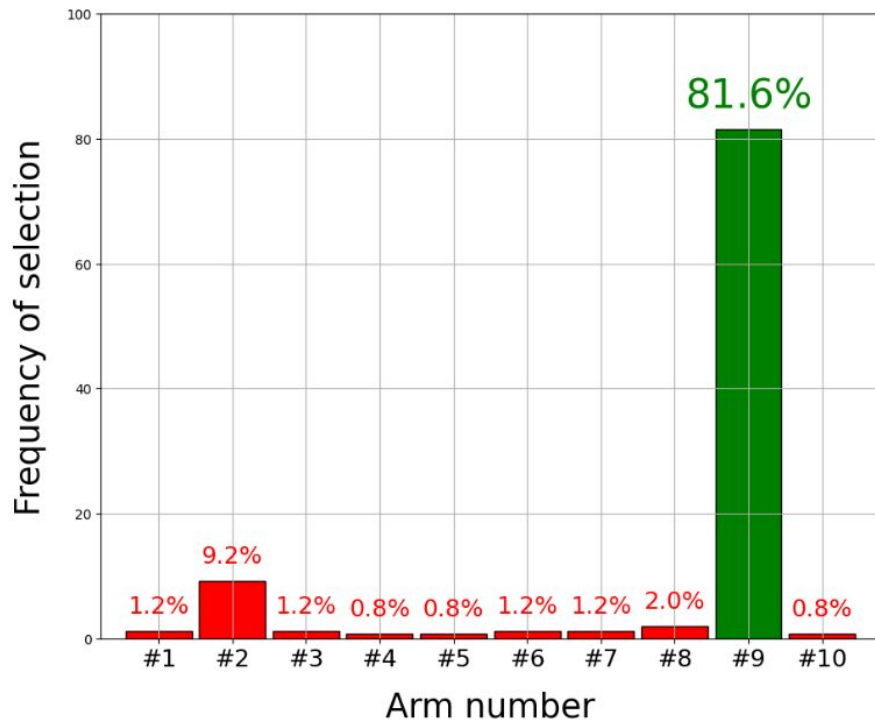


# Value-Difference-Based-Epsilon

Selection at each iteration

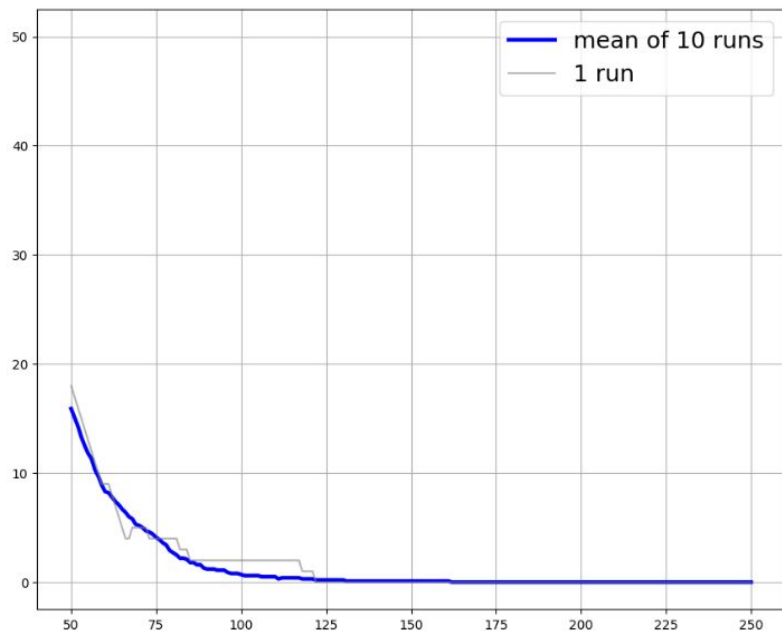


Arms selection frequency histogram

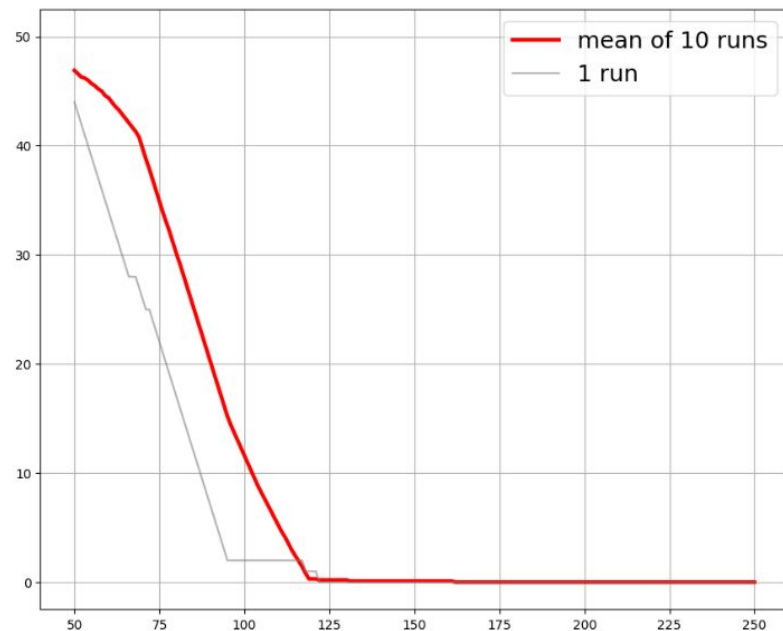


# Value-Difference-Based-Epsilon

Exploration rate through past 50 iterations

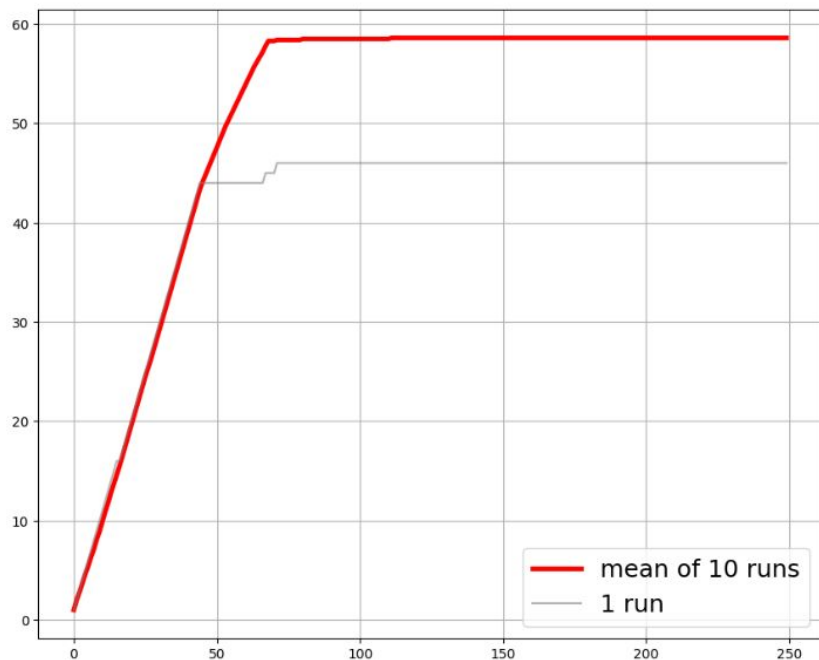


Num of unoptimal arms through past 50 iterations

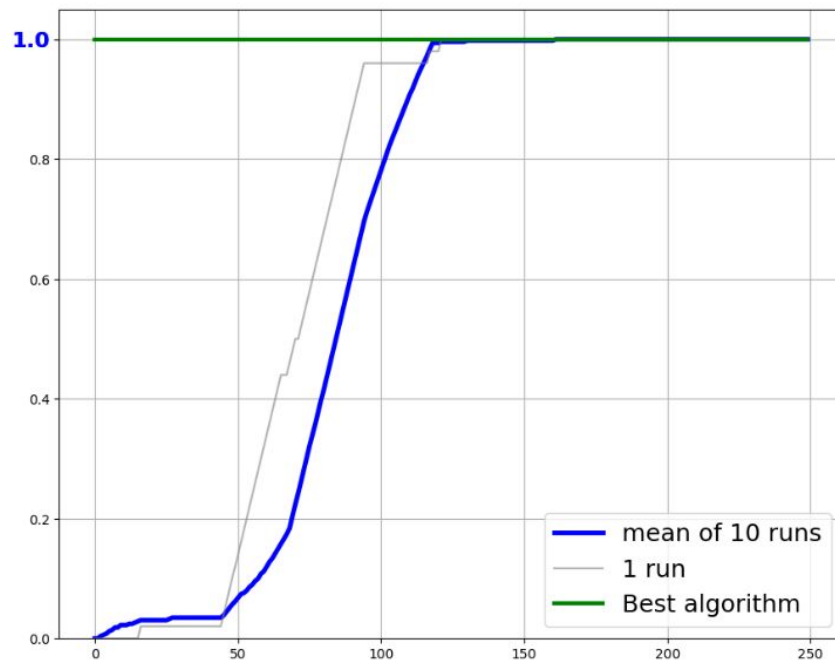


# Value-Difference-Based-Epsilon

Cumulative regret

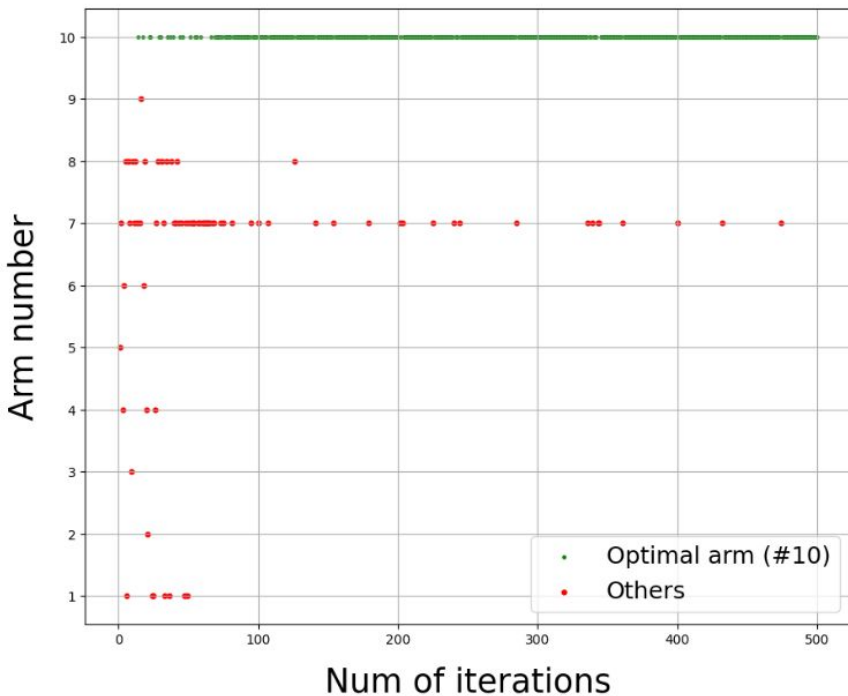


Convergence rate

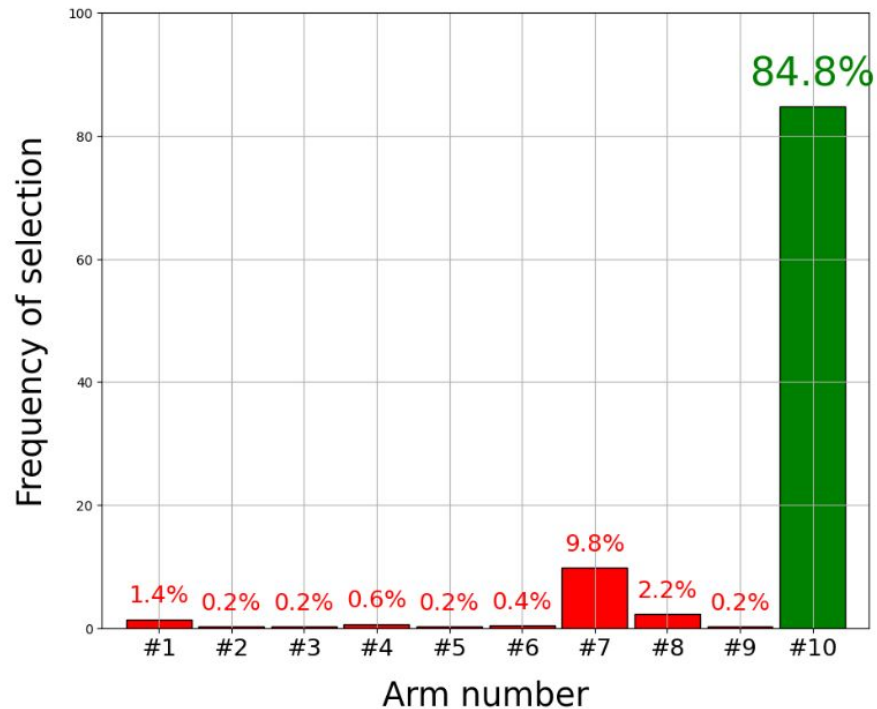


# SoftMax

Selection at each iteration

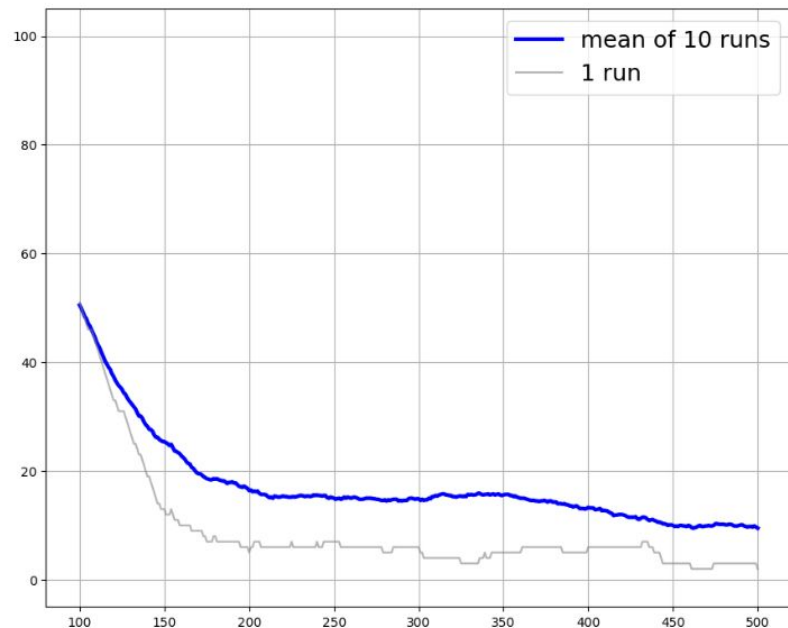


Arms selection frequency histogram

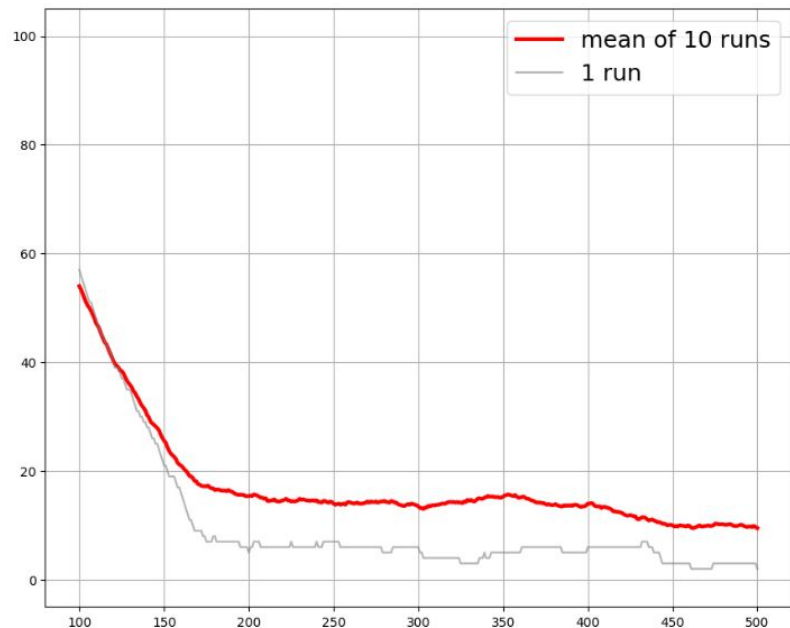


# SoftMax

Exploration rate through past 100 iterations

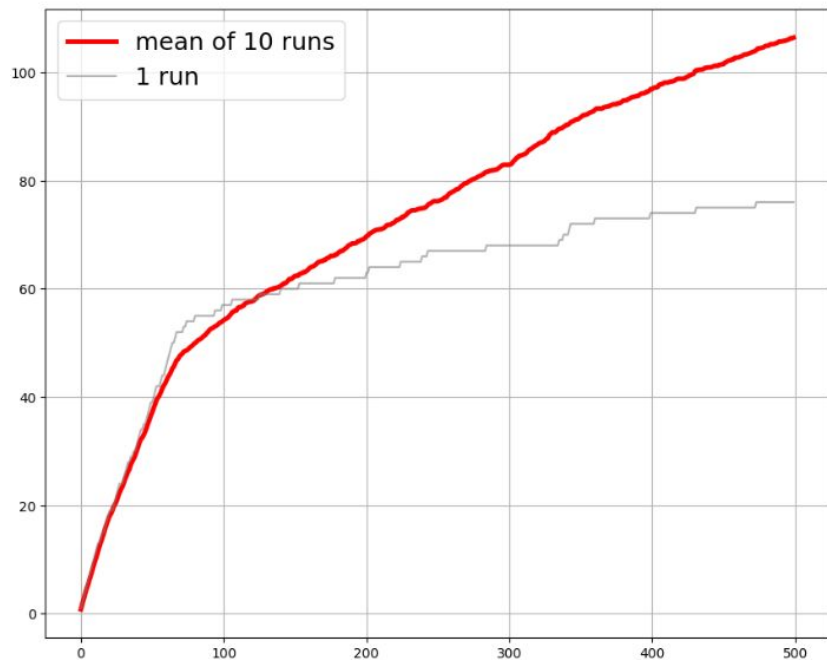


Num of unoptimal arms through past 100 iterations

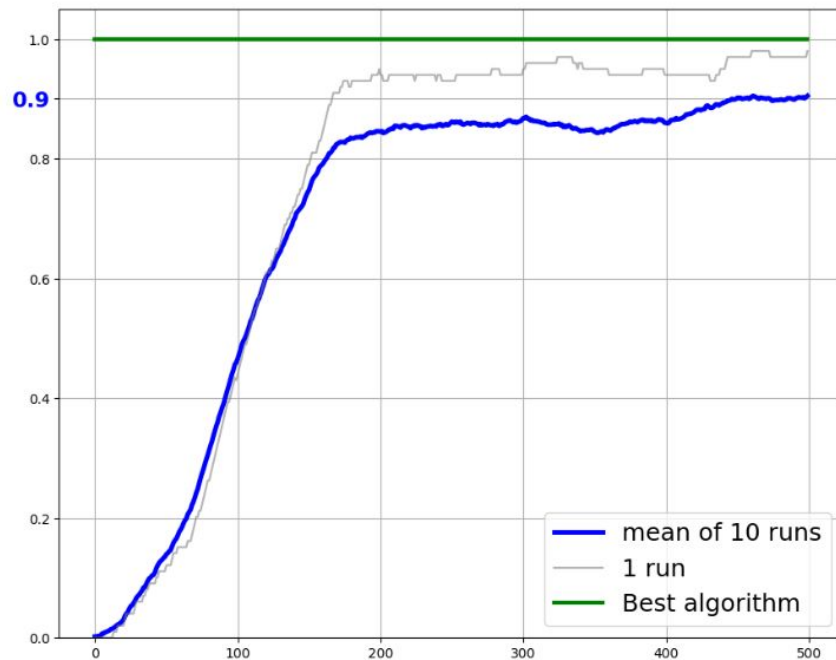


# SoftMax

Cumulative regret

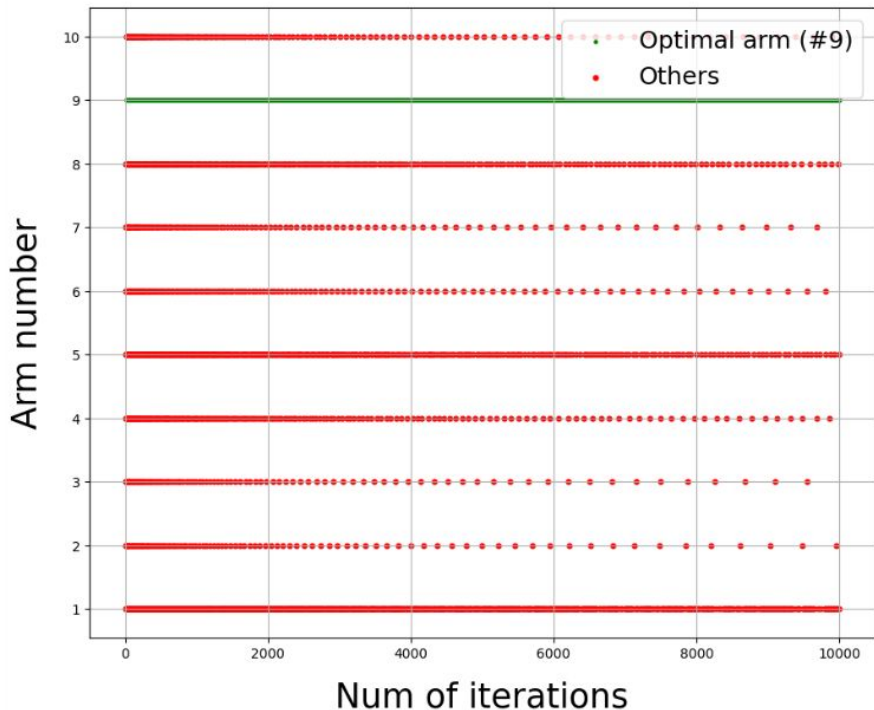


Convergence rate

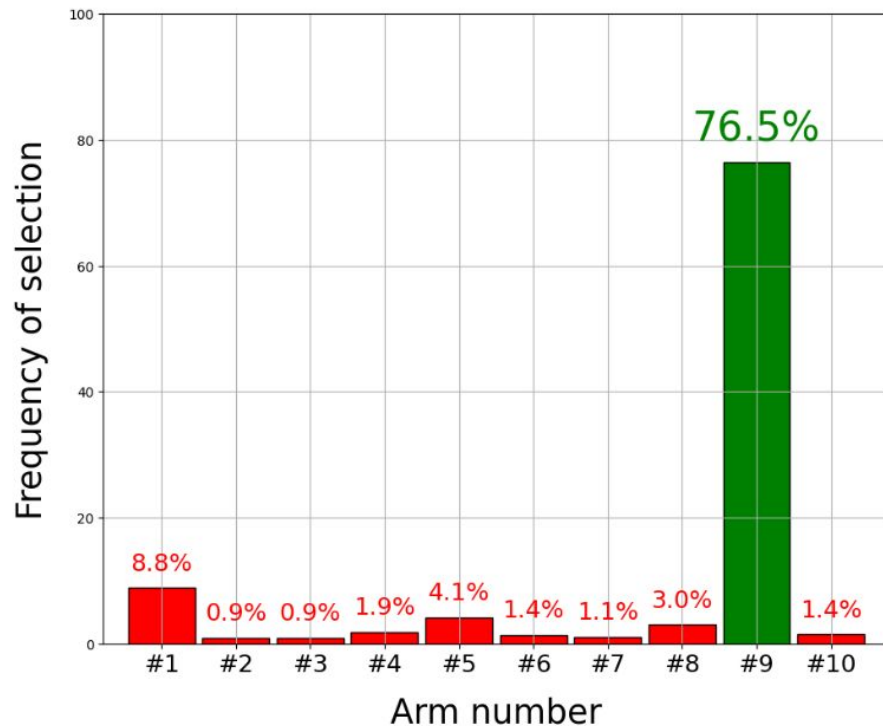


# Upper Confidence Bound

Selection at each iteration

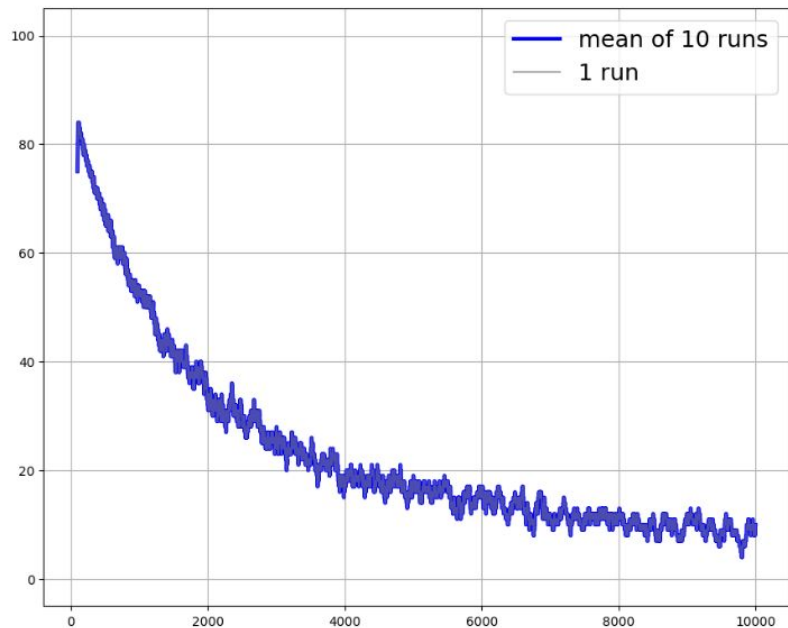


Arms selection frequency histogram

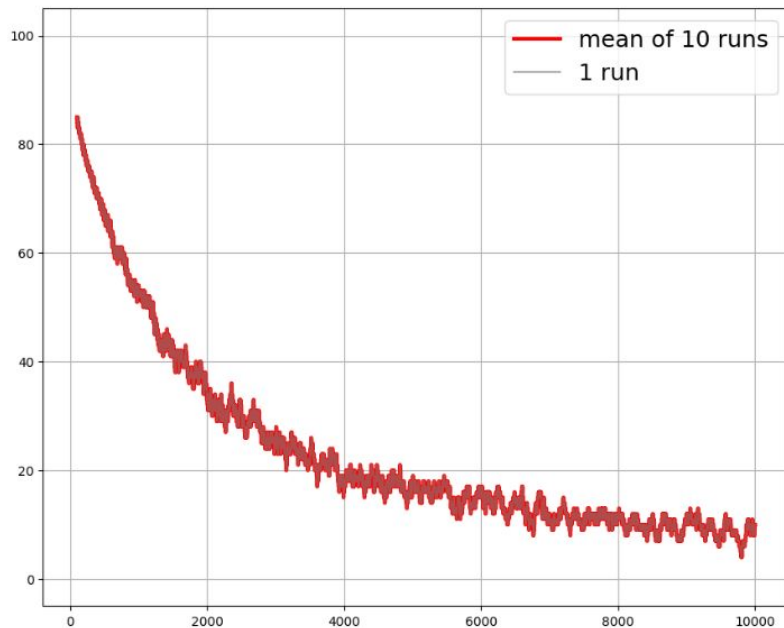


# Upper Confidence Bound

Exploration rate through past 100 iterations



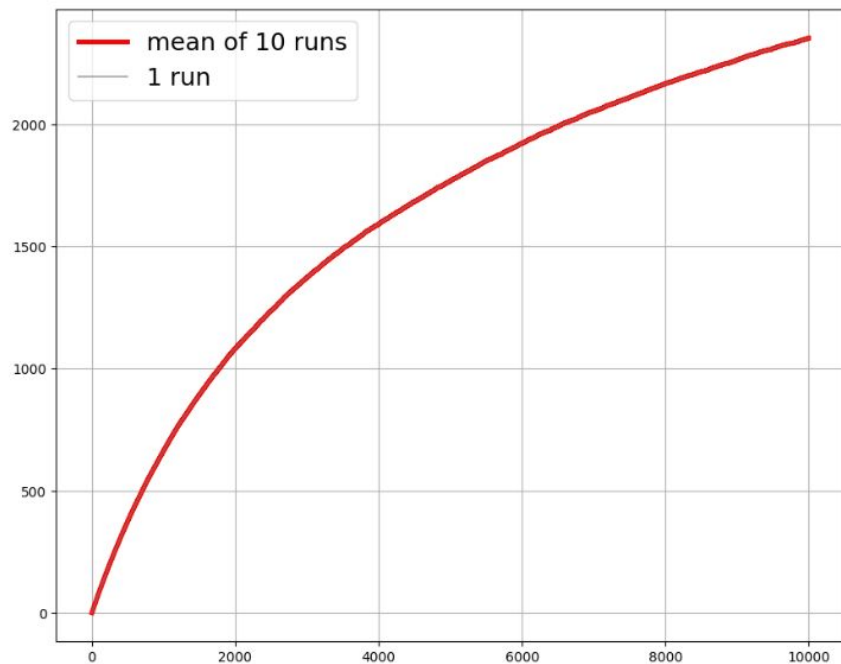
Num of unoptimal arms through past 100 iterations



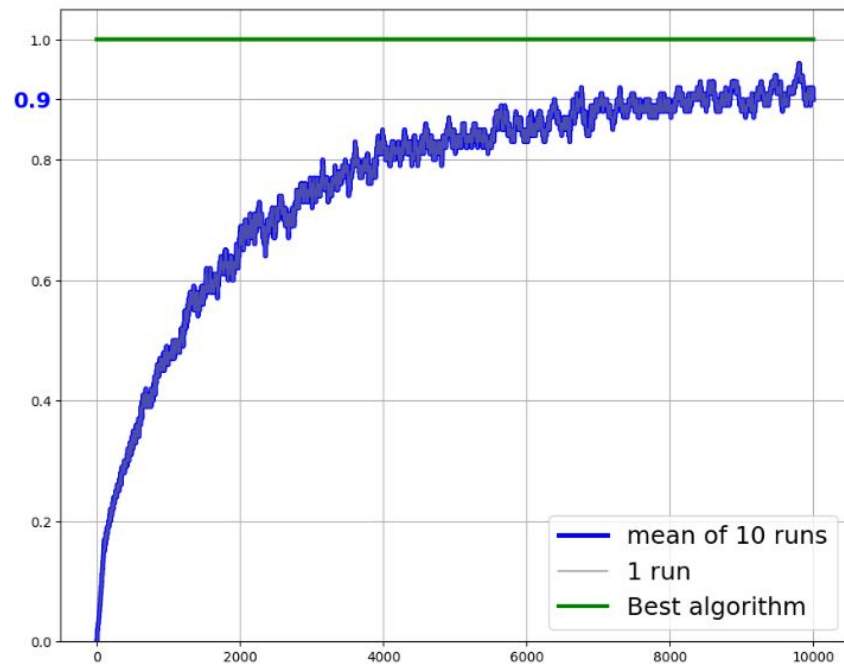


# Upper Confidence Bound

Cumulative regret

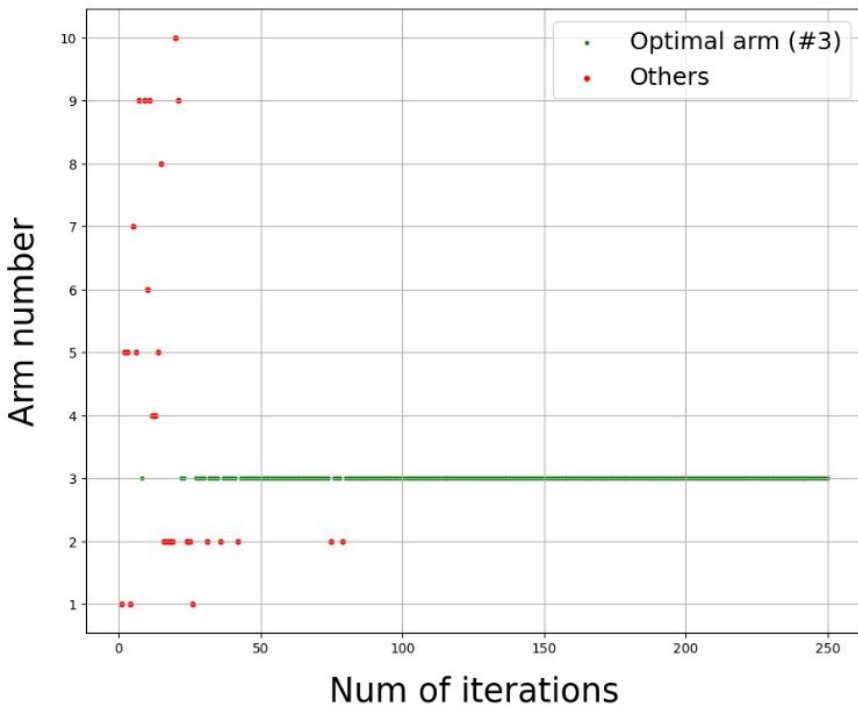


Convergence rate

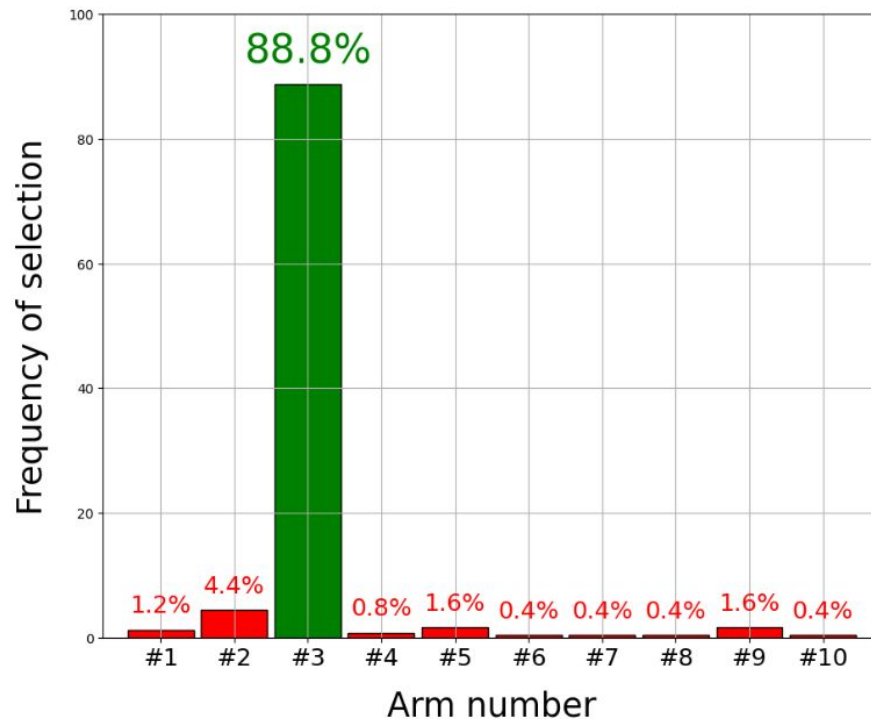


# Thompson Sampling

Selection at each iteration

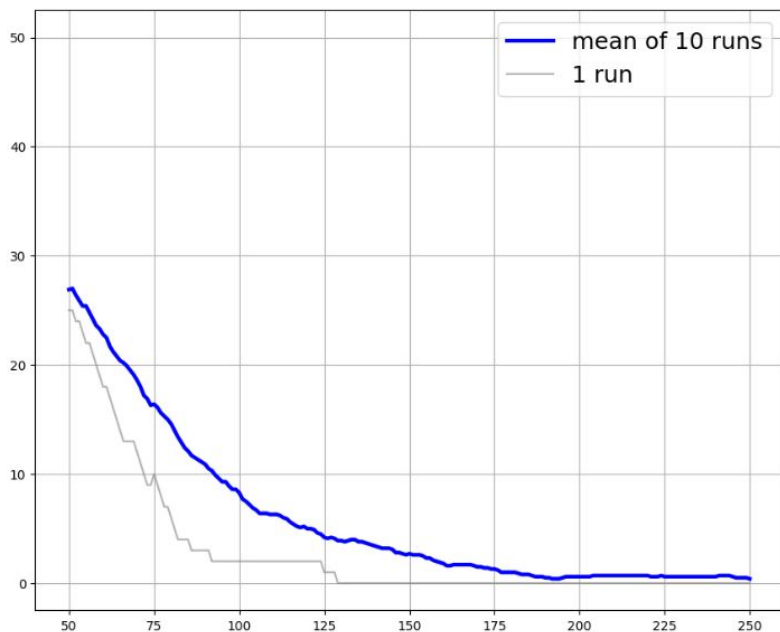


Arms selection frequency histogram

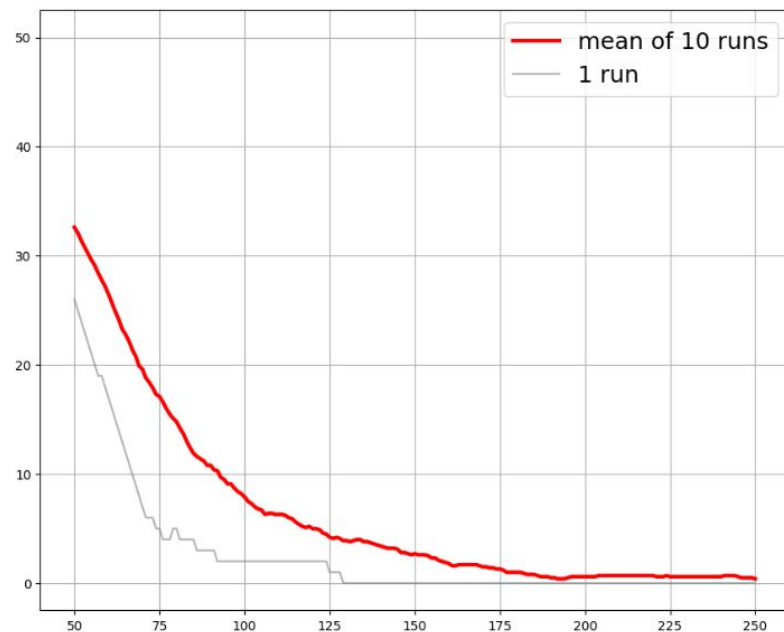


# Thompson Sampling

Exploration rate through past 50 iterations

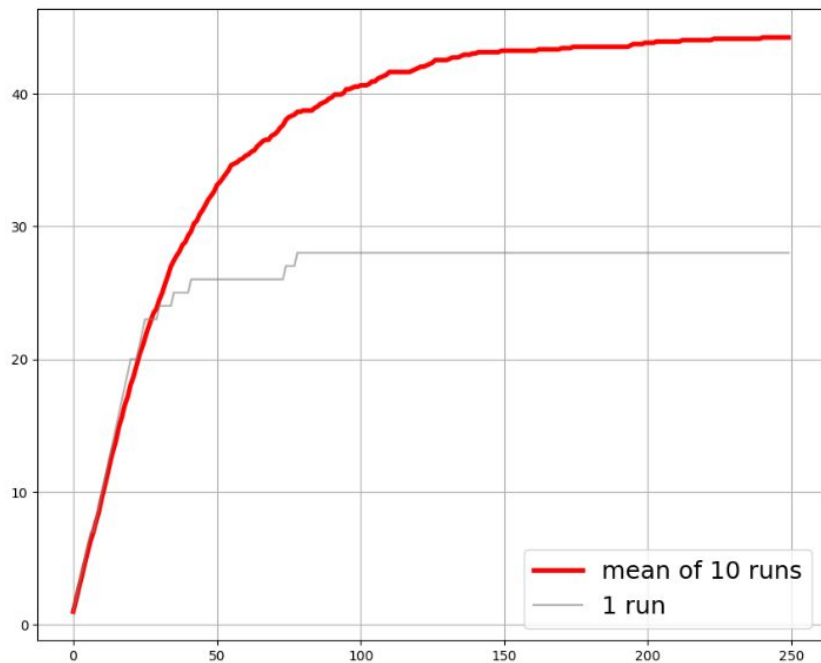


Num of unoptimal arms through past 50 iterations

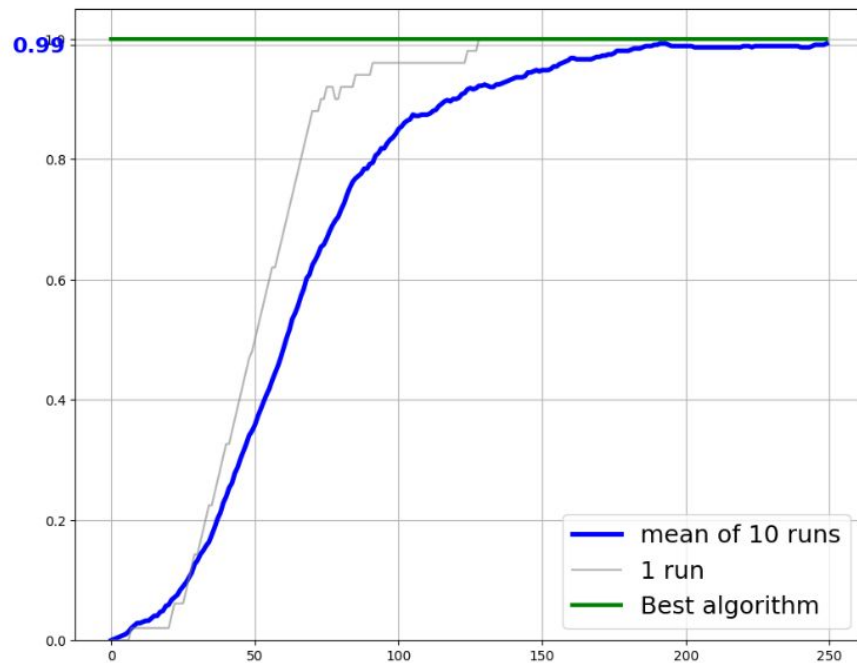


# Thompson Sampling

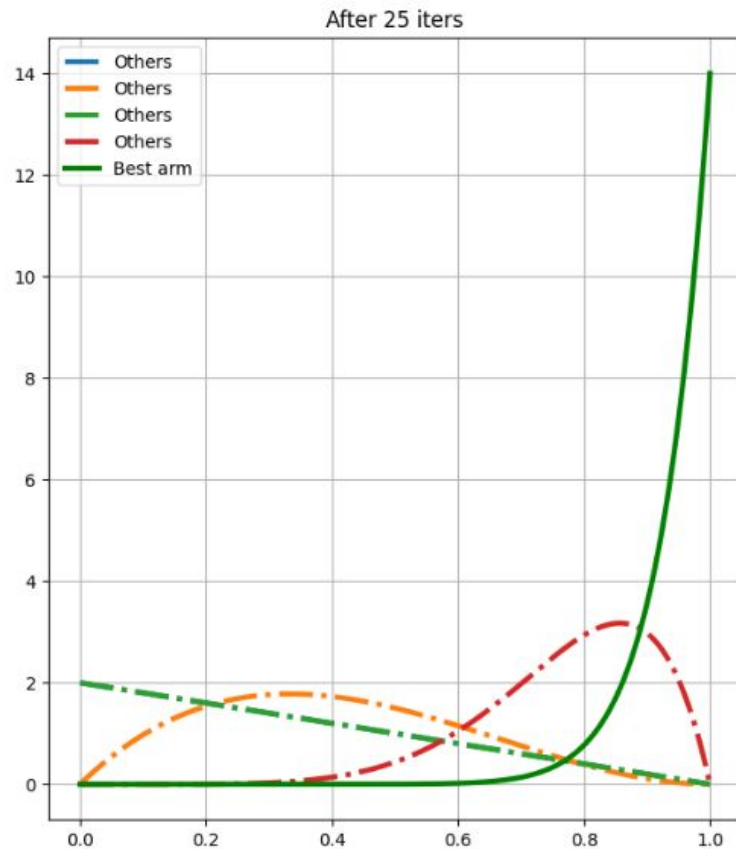
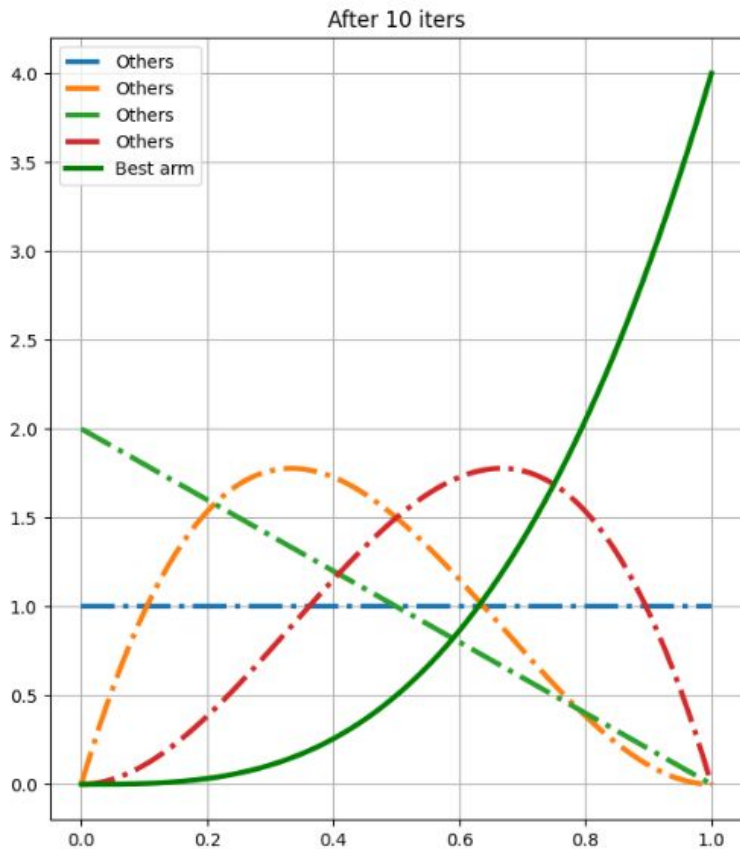
Cumulative regret



Convergence rate



# Thompson Sampling



# Модификации TS

---

**Algorithm 1** Online stochastic gradient descent with Thompson Sampling (SGD-TS)

---

**Input:**  $T, K, \tau, \alpha$ .

- 1: Randomly choose  $a_t \in [K]$  and record  $X_t, Y_t$  for  $t \in [\tau]$ .
- 2: Calculate the maximum-likelihood estimator  $\hat{\theta}_\tau$  by solving  $\sum_{t=1}^{\tau} (Y_t - \mu(X_t^T \theta)) X_t = 0$ .
- 3: Maintain convex set  $\mathcal{C} = \{\theta : \|\theta - \hat{\theta}_\tau\| \leq 2\}$ .
- 4:  $\tilde{\theta}_0 \leftarrow \hat{\theta}_\tau$ .
- 5: **for**  $t = \tau + 1$  **to**  $T$  **do**
- 6:   **if**  $t \% \tau = 1$  **then**
- 7:      $j \leftarrow \lfloor (t - 1) / \tau \rfloor$  and  $\eta_j = \frac{1}{\alpha_j}$ .
- 8:     Calculate  $\nabla l_{j,\tau}$  defined in Equation 3
- 9:     Update  $\tilde{\theta}_j \leftarrow \Pi_{\mathcal{C}} \left( \tilde{\theta}_{j-1} - \eta_j \nabla l_{j,\tau}(\tilde{\theta}_{j-1}) \right)$ .
- 10:     Compute  $\bar{\theta}_j = \frac{1}{j} \sum_{q=1}^j \tilde{\theta}_q$ .
- 11:     Compute  $A_j$  defined in Equation 5.
- 12:     Draw  $\theta_j^{\text{TS}} \sim \mathcal{N}(\bar{\theta}_j, A_j)$ .
- 13:   **end if**
- 14:   Pull arm  $a_t \leftarrow \arg\max_{a \in [K]} \mu(x_{t,a}^T \theta_j^{\text{TS}})$  and observe reward  $Y_t$ .
- 15: **end for**

## Online Stochastic Gradient Descent and Thompson Sampling

---

**Algorithm 1** BootstrapLinTS for partially observable delayed feedback

---

**Input:**  $n_{\text{prior}}, D_{\text{max}}, T, d, K$ .

- 1: Data  $D_0 = ()$
- 2: **for**  $n = 1, \dots, T$  **do**
- 3:   Update data  $D_n$  with observed conversions
- 4:   **for**  $j = 1, \dots, n_{\text{prior}}$  **do**
- 5:     Sample prior  $\vartheta_j$  and  $x_j$  uniformly over  $[0, 1]^d$
- 6:     Normalise sampled  $\vartheta_j$  and  $x_j$
- 7:     Sample prior reward from Bernoulli( $\vartheta_j \cdot x_j$ )
- 8:     Sample delays uniformly over  $[0, D_{\text{max}}]$
- 9:   **end for**
- 10:   Concatenate  $n_{\text{prior}}$  times and rewards with  $D_n$
- 11:   Sample with replacement  $n + n_{\text{prior}}$  data points
- 12:   Estimate  $\hat{S}(t, x)$  and  $\hat{p}_1(x)$  via EM
- 13:   Observe current contexts  $x_A, A = 1, \dots, K$
- 14:   **for**  $A = 1, \dots, K$  **do**
- 15:     Calculate probability  $(1 - \hat{S}(T, x_A)) \hat{p}_1(x_A)$
- 16:   **end for**
- 17:   Select arm  $\arg\max_i (1 - \hat{S}(T, x_A)) \hat{p}_1(x_A)$
- 18: **end for**

## Bootstrapped Thompson Sampling

# Linear Thompson Sampling (LinTS)

---

**Algorithm 1** Thompson Sampling for Contextual bandits

---

Set  $B = I_d, \hat{\mu} = 0_d, f = 0_d$ .

**for all**  $t = 1, 2, \dots$ , **do**

    Sample  $\tilde{\mu}(t)$  from distribution  $\mathcal{N}(\hat{\mu}, v^2 B^{-1})$ .

    Play arm  $a(t) := \arg \max_i b_i(t)^T \tilde{\mu}(t)$ , and observe reward  $r_t$ .

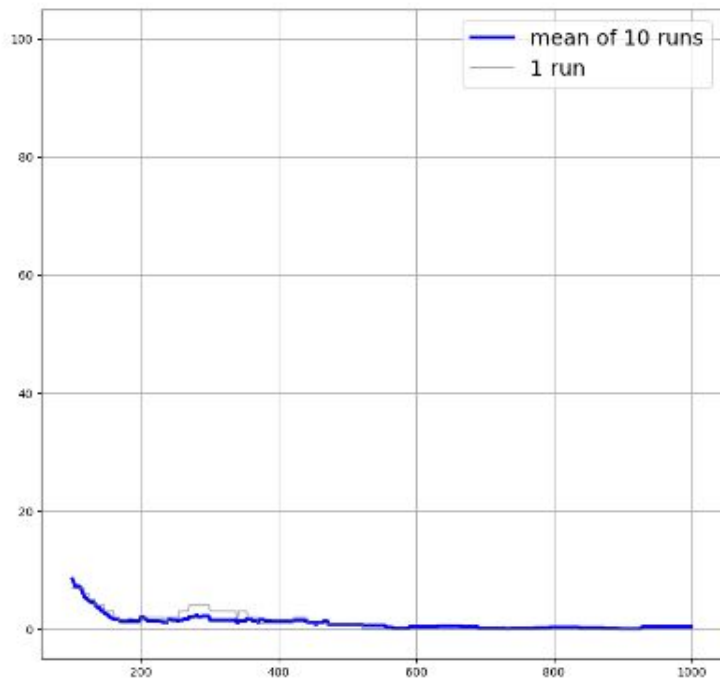
    Update  $B = B + b_{a(t)}(t)b_{a(t)}(t)^T, f = f + b_{a(t)}(t)r_t, \hat{\mu} = B^{-1}f$ .

**end for**

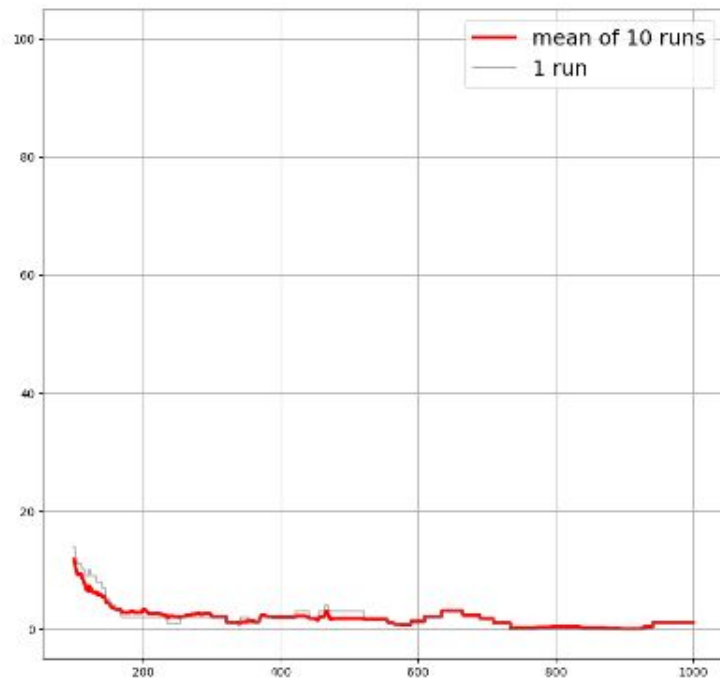
---

# LinTS

Exploration rate through past 100 iterations



Num of unoptimal arms through past 100 iterations

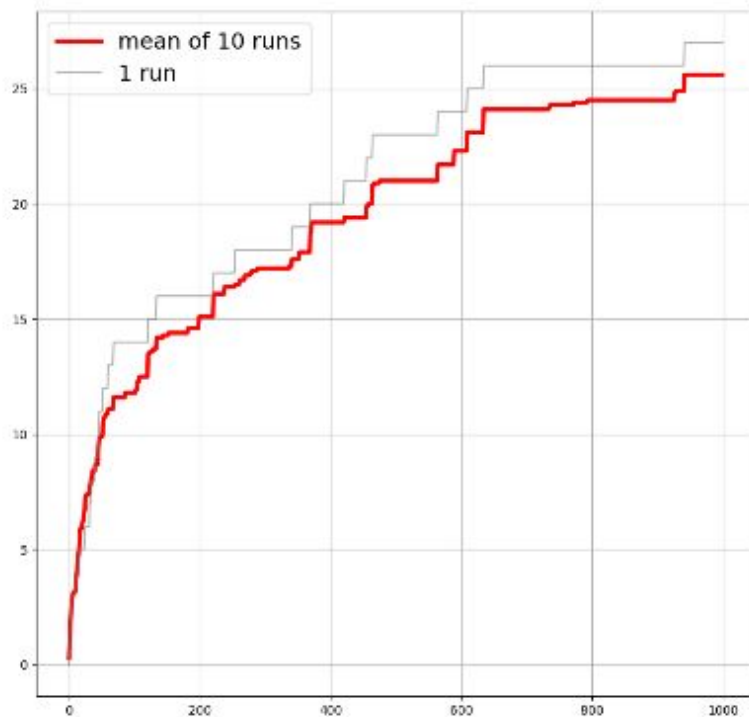


Mushroom dataset

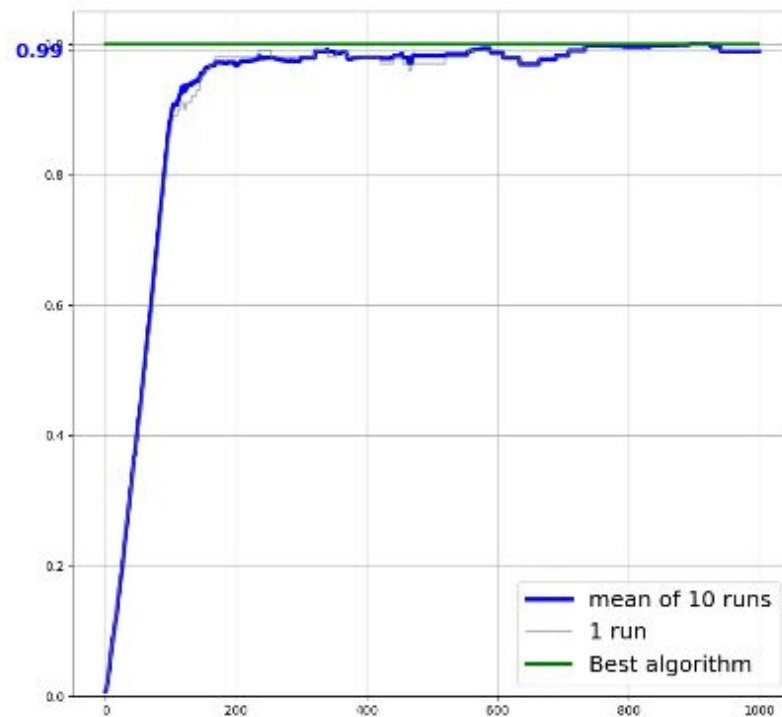


# LinTS

Cumulative regret



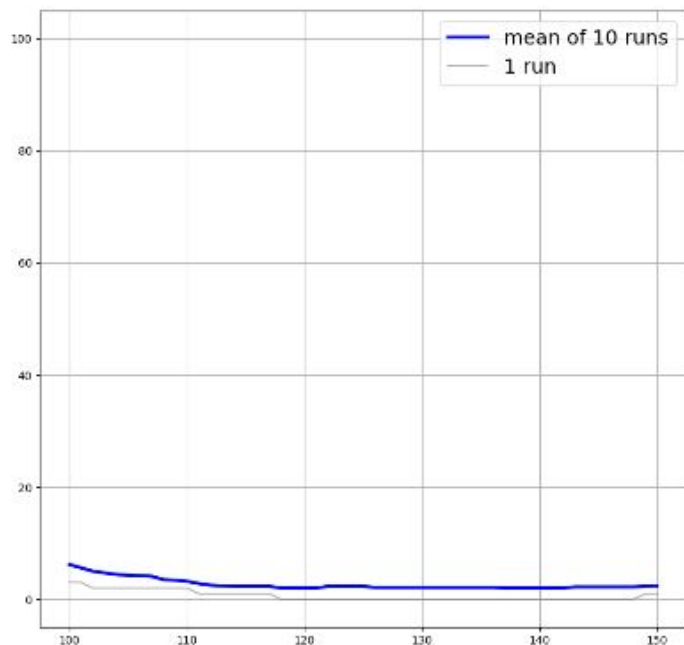
Convergence rate



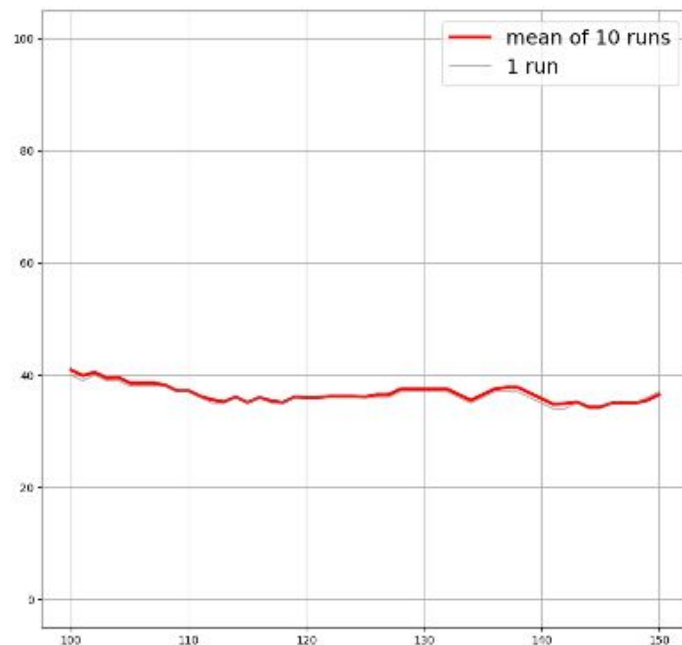
Mushroom dataset

# LinTS

Exploration rate through past 100 iterations



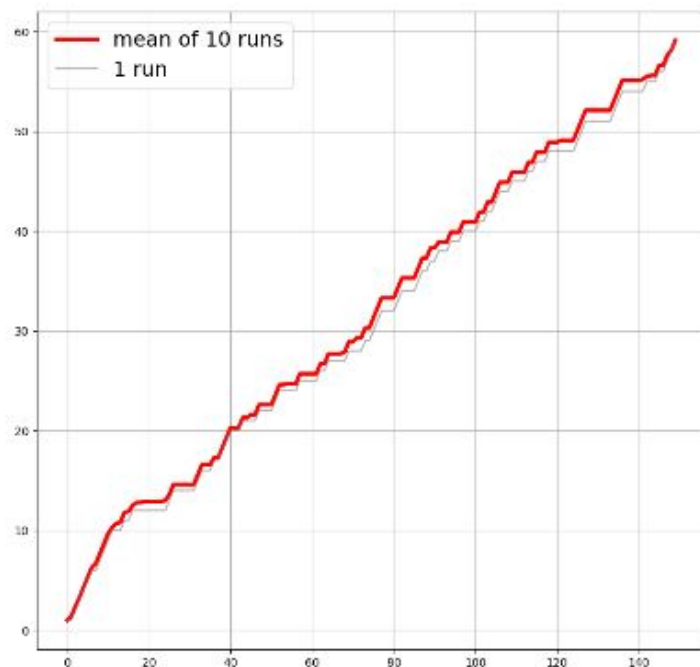
Num of unoptimal arms through past 100 iterations



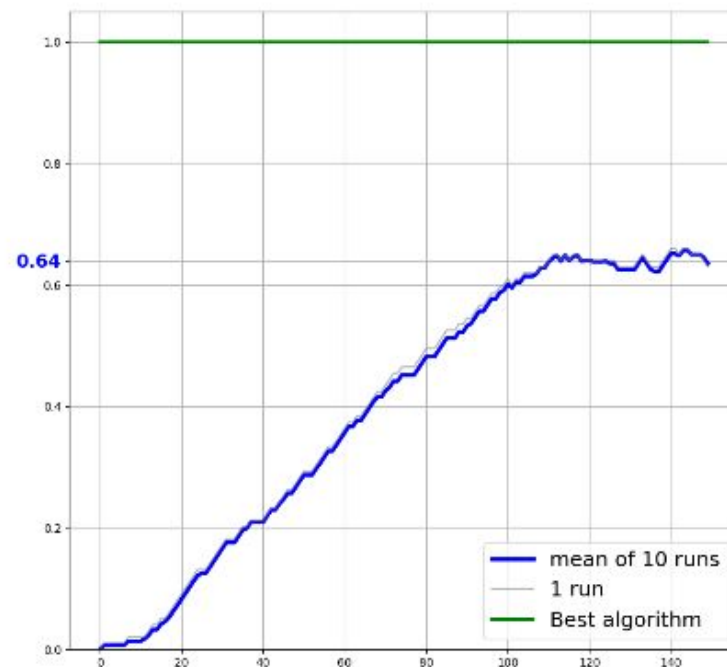
Iris dataset

# LinTS

Cumulative regret



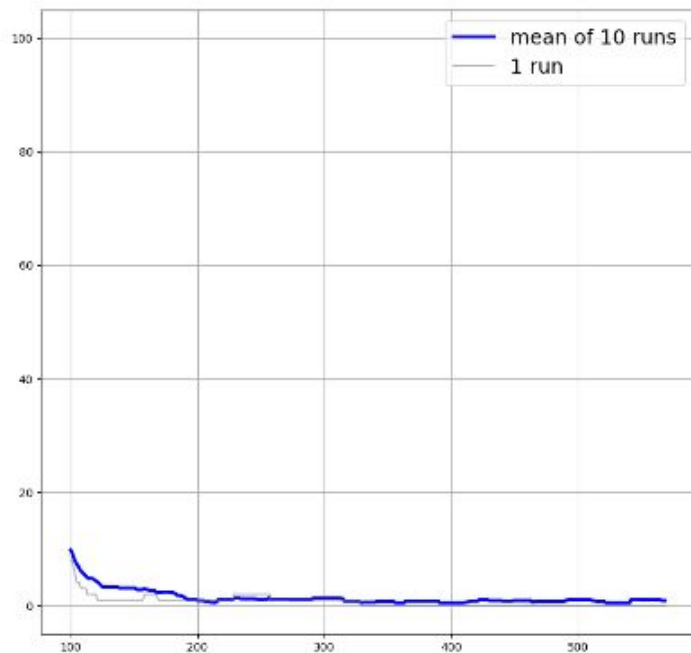
Convergence rate



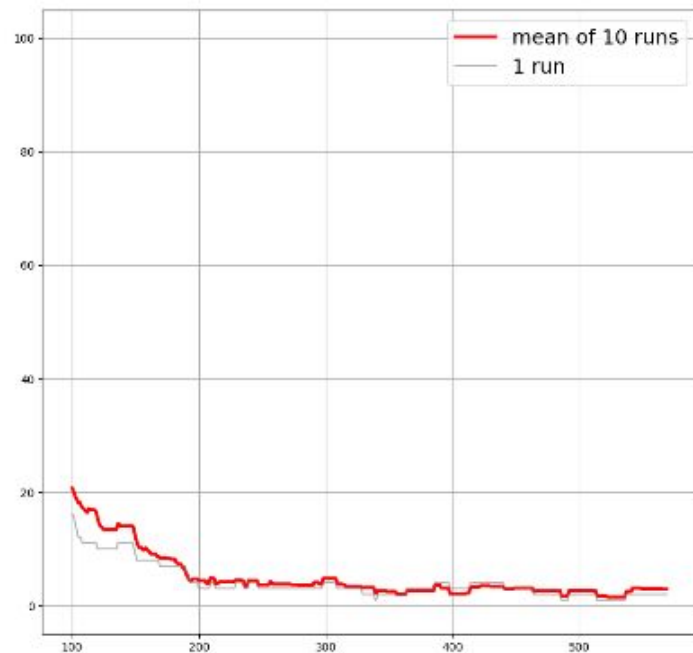
Iris dataset

# LinTS

Exploration rate through past 100 iterations



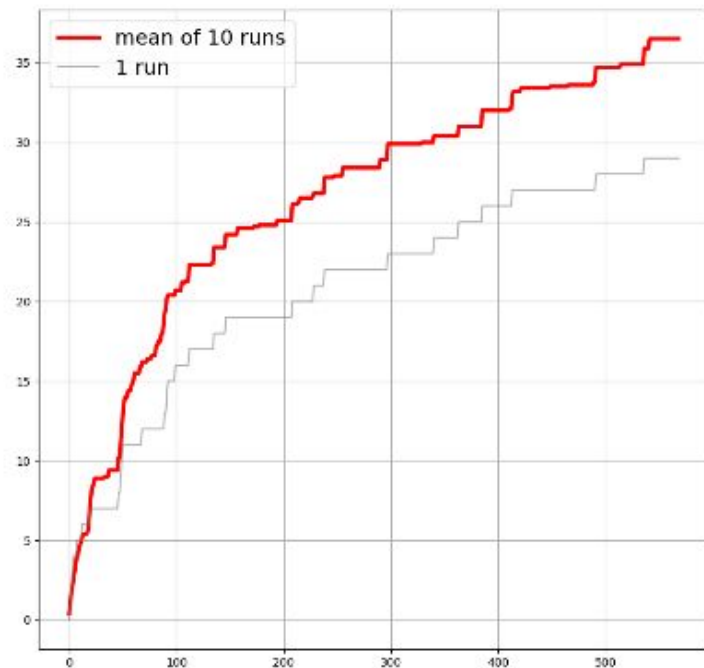
Num of unoptimal arms through past 100 iterations



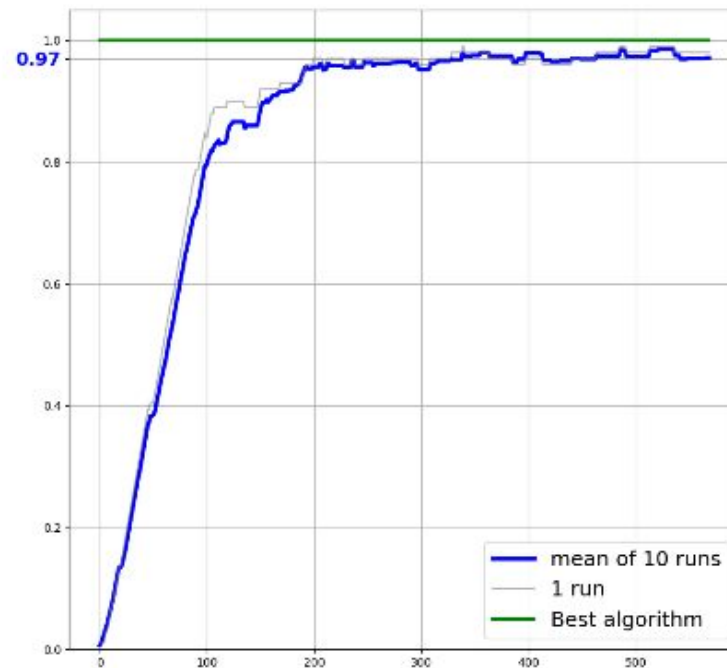
Breast cancer dataset

# LinTS

Cumulative regret

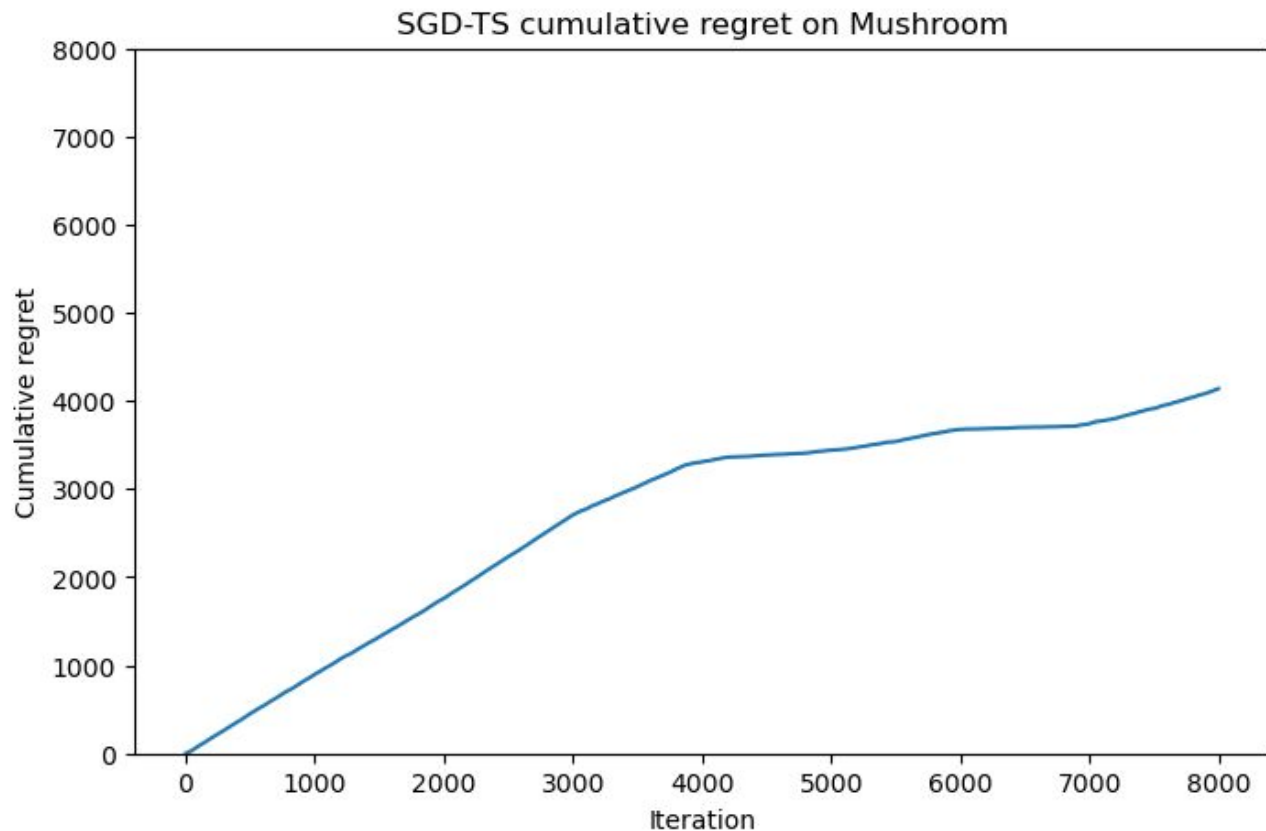


Convergence rate



Breast cancer dataset

# SGD-TS



# GTS

на всём датасете GTS сходится, но медленно (Mushroom dataset)

