

# Проект по курсу «Продвинутые методы машинного обучения»

## Многорукие бандиты: Торговый бот

Выполнили:  
студенты 4 курса  
Клеваичук Александр,  
Михеева Екатерина,  
Молвинская Алина.

Москва, 2024

# Многорукие бандиты: торговый бот

**Многорукий бандит** — это алгоритм принятия решений, суть которого заключается в обеспечении баланса между исследованием и использованием какого-то из рассматриваемых вариантов решения с целью максимизации целевой метрики и минимизации убытков.

## Цели

- Разработка торгового бота
- Оптимизация торговой стратегии
- Подключение к реальным данным

## Задачи

- Изучение и анализ стратегий
- Сбор и обработка данных
- Проектирование и реализация торговой стратегии:
- Тестирование

# Математическая постановка

Алгоритм для задачи МАВ должен решать, какую рукоятку играть на каждом временном шаге  $t$ , исходя из результатов предыдущих  $t-1$  игр. Пусть  $i$  обозначает (неизвестное) ожидаемое вознаграждение для рукоятки  $i$ . Целью является максимизация математического ожидания общего вознаграждения за время  $T$ , т.е.

$$E \left[ \sum_{t=1}^T \mu_{i(t)} \right]$$

где  $i(t)$  — это рукоятка, сыгранная на шаге  $t$ , и ожидание берётся по случайным выборам  $i(t)$ , сделанным алгоритмом.

# Метрики

В качестве метрик в нашей задаче мы использовали:

- **total rewards** (Сумма всех полученных наград за время эксперимента)

$$r_{total,t} = \sum_{i=1}^t R_i$$

- **cumulativer regret** (Потери из-за выбора не оптимальных рук)

$$reg_t = \sum_{i=1}^t R_{opt} - R_t$$

# Датасеты и данные

Набор данных о ценах 234 российских акций. Датасет содержит исторические цены открытия, максимума, минимума, закрытия и объема акций, торгуемых на финансовых рынках Российской Биржи. Взяты почасовые данные закрытия в период с 01.06.2023 по 27.08.2024

	open	high	low	close	volume
datetime					
2023-06-01 10:00:00	2328.0	2356.0	2318.0	2343.6	143039.0
2023-06-01 11:00:00	2343.2	2348.0	2333.2	2333.2	48477.0
2023-06-01 12:00:00	2333.2	2334.2	2304.8	2311.4	86184.0
2023-06-01 13:00:00	2313.2	2336.4	2294.8	2330.0	104030.0
2023-06-01 14:00:00	2329.2	2333.6	2316.0	2326.8	33793.0
...	...	...	...	...	...
2024-08-27 19:00:00	4092.2	4096.4	4071.2	4071.2	9576.0
2024-08-27 20:00:00	4092.2	4096.4	4071.2	4071.2	9576.0
2024-08-27 21:00:00	4092.2	4096.4	4071.2	4071.2	9576.0
2024-08-27 22:00:00	4092.2	4096.4	4071.2	4071.2	9576.0
2024-08-27 23:00:00	4092.2	4096.4	4071.2	4071.2	9576.0

4437 rows × 5 columns

# Ход работы

1. Создаём класс Strategy - основа для реализации различных стратегий многоруких бандитов.
2. Thompson Sampling
3. Добавляем транзакционные издержки
4. Создаём среду, имитирующую поведение рынка акций
5. Создаём класс Bandit, который объединяет среду и стратегию.
6. Запускаем модель

# Метод Thompson Sampling

Для выбора оптимального портфеля из акций реализован класс Thompson, основанный на методе бета-распределений: у каждой акции (или "руки") есть параметры успехов и неудач (alphas и betas), которые обновляются после каждой итерации. Модель модифицирует alpha и beta параметры распределения в Томпсоне, основываясь на штрафе - комиссии брокера. В beta-распределении участвуют последние 840 часов для каждой акции (60 последних дней торгов).

Также добавлен параметр "Жадности» и скорости ее возрастания.

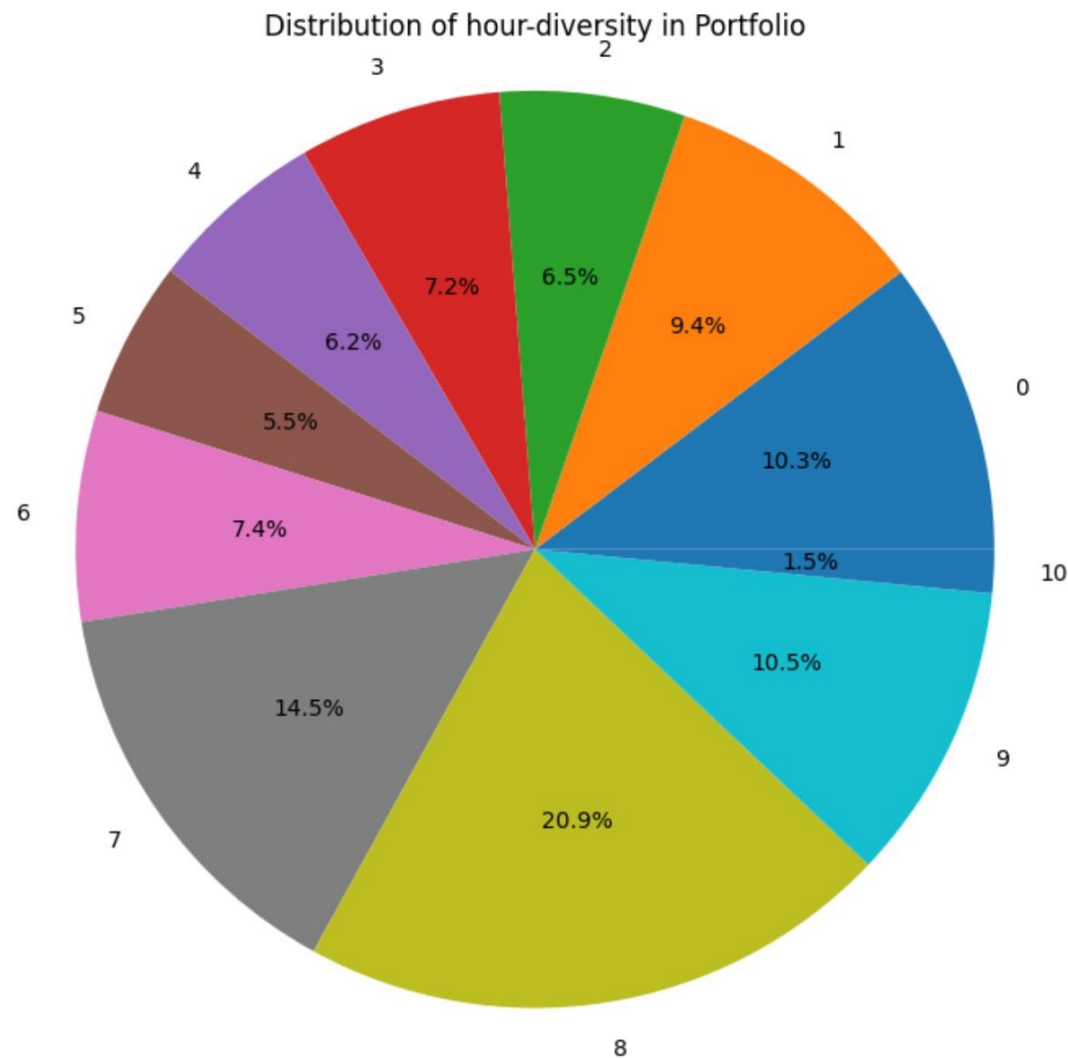
Помимо этого предусмотрена система штрафов за объём покупок - транзакционные издержки, связанные с реальными рынками.

# Ход работы

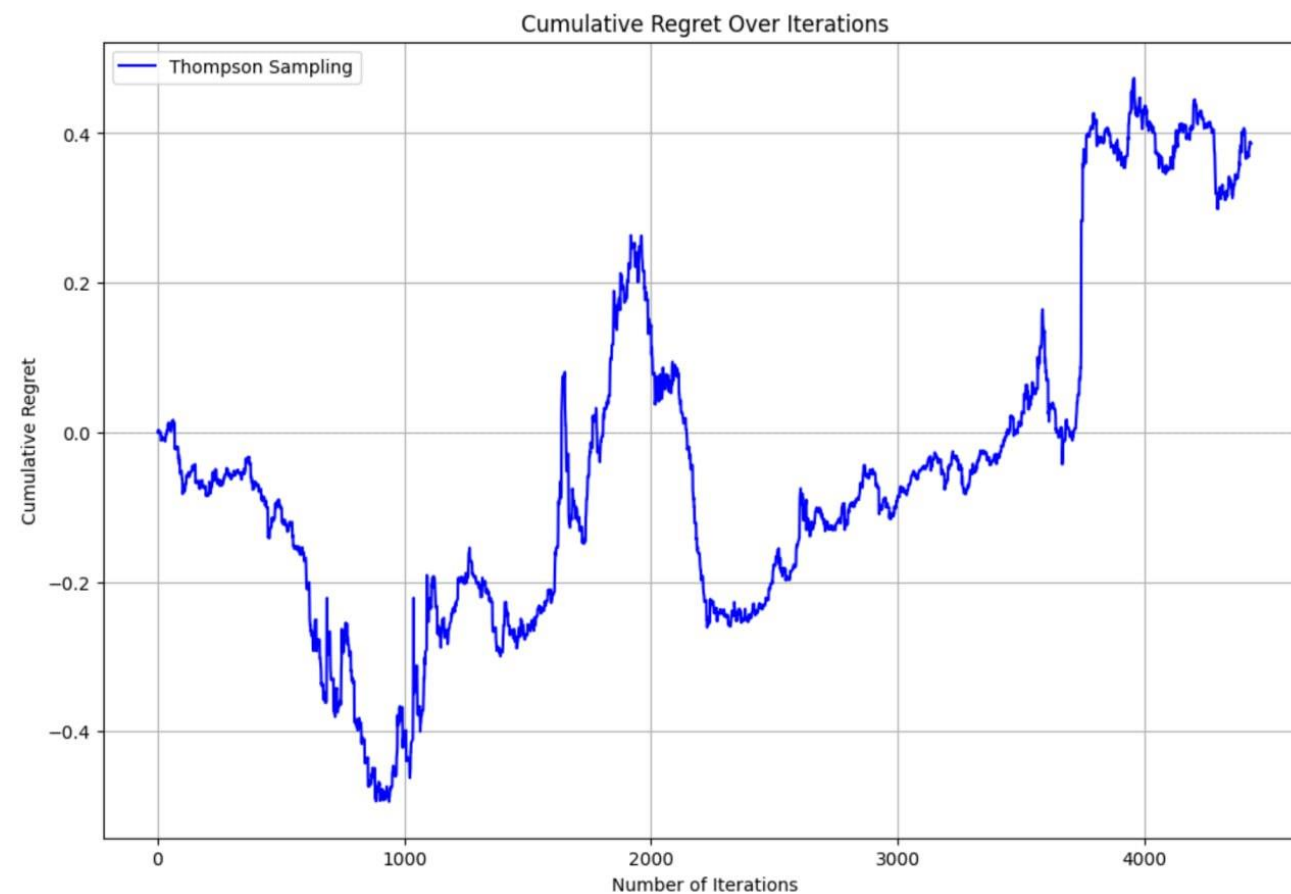
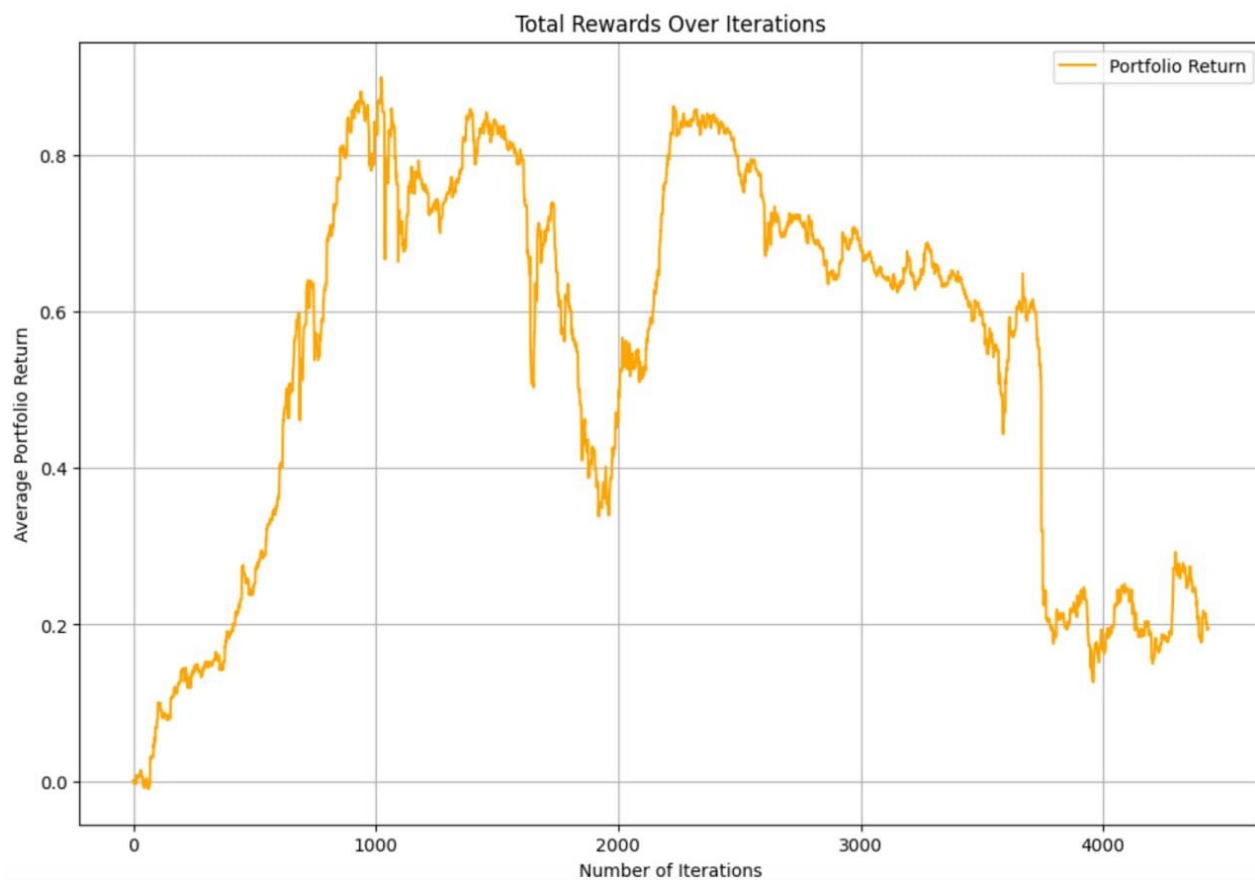
1. Создаём класс Strategy - основа для реализации различных стратегий многоруких бандитов.
2. Thompson Sampling
3. Добавляем транзакционные издержки и параметр жадности
4. Создаём среду, имитирующую поведение рынка акций
5. Создаём класс Bandit, который объединяет среду и стратегию.
6. Запускаем модель



# Результаты



# Результаты





ПРЕЗИДЕНТСКАЯ  
АКАДЕМИЯ

**СПАСИБО ЗА ВНИМАНИЕ!**

Москва, 2024

2024

РАНХиГС

```
1 class Strategy:
2     def __init__(self, n_arms: int, n_portfolio: int):
3         self.n_arms = n_arms
4         self.n_portfolio = n_portfolio
5         self.n_iters = 0
6         self.arms_states = np.zeros(n_arms)
7         self.arms_actions = np.zeros(n_arms)
8
9     def flush(self):
10        self.n_iters = 0
11        self.arms_states = np.zeros(self.n_arms)
12        self.arms_actions = np.zeros(self.n_arms)
13
14    def update_reward(self, arm: int, reward: float):
15        self.n_iters += 1
16        self.arms_states[arm] += reward
17        self.arms_actions[arm] += 1
18
19    def choose_arm(self):
20        raise NotImplementedError
```

```
3 class Thompson(Strategy):
4     def __init__(self, n_arms: int, n_portfolio: int, penalty: float, greed:
5         float, greed_tempo: float):
6         super().__init__(n_arms, n_portfolio, penalty, greed, greed_tempo)
7         self.alphas = {i: [1.0] for i in range(n_arms)}
8         self.betas = {i: [-1.0] for i in range(n_arms)}
9         self.greed = greed
10        self.greedness = {i: greed for i in range(n_arms)}
11        self.portfolio = []
12        self.greed_tempo = greed_tempo
13        self.penalty = penalty
14
15    def choose_arm(self):
16        sampled_values = {i: np.random.beta(
17                                sum(self.alphas[i][:840]),
18                                abs(sum(self.betas[i][:840]))
19                                )
20                        for i in range(self.n_arms)}
21
22        best_arms = sorted(sampled_values, key=sampled_values.get, reverse=
23                            True)[:self.n_portfolio]
24
25        self.portfolio = best_arms
```

```
1 return best_arms
2
3 def update_reward(self, arm: int, reward: float):
4
5     the received reward
6     if arm in self.portfolio:
7         if reward <= 0:
8             self.betas[arm].insert(0, reward)
9             self.alphas[arm].insert(0, 1e-20 )
10            self.greedness[arm] = self.greed
11        else:
12            self.alphas[arm].insert(0, reward)
13            self.betas[arm].insert(0, -1e-20)
14            self.greedness[arm] = self.greed
15    else:
16        reward -= self.penalty / self.greedness[arm]
17
18    if reward <= 0:
19        self.betas[arm].insert(0, reward)
20        self.alphas[arm].insert(0, 1e-20 )
21        if self.greedness[arm] > 1:
22            self.greedness[arm] -= self.greed_tempo
23    else:
24        self.alphas[arm].insert(0, reward )
25        self.betas[arm].insert(0, -1e-20 )
26        self.greedness[arm] += self.greed_tempo
```

```
1 class StockMarketEnv:
2     def __init__(self, stock_files: list):
3         self.stock_data = self.load_stock_data(stock_files)
4         self.n_arms = len(self.stock_data)
5
6     def load_stock_data(self, stock_files):
7         data = {}
8         for stock_file in stock_files:
9             df = pd.read_csv(stock_file, parse_dates=['datetime'], index_col
              = 'datetime')
10            df = df.reindex(ind)
11            df = df.fillna(method='ffill')
12            df = df.sort_index()
13            df = df.loc['2023-06-01_00:00:00':]
14            if df['close'].isna().sum() == 0:
15                df['return'] = df['close'].pct_change()
16                data[stock_file] = df['return'].dropna()
17        return data
```

```
1 class Bandit:
2
3     def __init__(self, env: StockMarketEnv, strategy: Strategy):
4         self.env = env
5         self.strategy = strategy
6
7     def action(self, inde):
8         arm = self.strategy.choose_arm()
9         reward = self.env.pull_arm(arm, inde)
10        self.strategy.update_reward(arm, reward)
```



```
1 def calculate_regret(env: StockMarketEnv, strategy: Strategy):
2     strategy.flush()
3     bandit = Bandit(env, strategy)
4     regrets = []
5     total_rewards = []
6     for i in range(len(ind)-2):
7         expected_returns = env.get_expected_returns(i)
8         optimal_return = np.mean(expected_returns)
9         selected_stocks = bandit.strategy.choose_arm()
10        strat.append(selected_stocks)
11
12        portfolio_return = np.mean([bandit.env.pull_arm(stock, i) for stock
13                                     in selected_stocks])
14
15        regrets.append(optimal_return-portfolio_return)
16
17        total_rewards.append(portfolio_return)
18
19        for stock in range(len(stock_files)):
20            reward = bandit.env.pull_arm(stock, i)
21            bandit.strategy.update_reward(stock, reward)
22
23    return regrets, total_rewards
```