

Многорукие бандиты: торговый бот

18 декабря 2024 г.

Клевайчук Александр, Михеева Екатерина, Молвинская Алина

1 Введение

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.

Этот термин описывает ситуацию, в которой агент сталкивается с несколькими альтернативами (или "руками игровых автоматов"), каждая из которых имеет свою скрытую вероятность выигрыша. Метод многоруких бандитов подходит для задач, связанных с принятием решений в условиях неопределенности. В рамках этого подхода можно разработать торговую стратегию, которая будет адаптироваться и улучшаться на основе собранного опыта и данных о рынке.

Цели и задачи проекта:

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

2 Теоретическая основа

2.1 Математическая постановка

Мы рассматриваем задачу многоруких бандитов (MAB): нам дан игровой автомат с N рукоятками; на каждом временном шаге ($t = 1, 2, 3, \dots$) необходимо выбрать одну из N рукояток для игры. Каждая рукоятка i при игре приносит случайное вознаграждение, которое распределено по некоторому фиксированному (неизвестному) распределению с поддержкой на интервале $[0, 1]$. Случайные вознаграждения, получаемые при многократной игре на одной рукоятке, являются независимыми и идентично распределёнными и не зависят от игр на других рукоятках. Вознаграждение наблюдается немедленно после игры на рукоятке.

Алгоритм для задачи MAB 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.

Главная проблема задачи многоруких бандитов заключается в балансе между:

Exploration (исследование): выбор рук, чтобы собрать больше информации о распределениях награды.

Exploitation (использование): выбор наилучшей из уже известных рук для максимизации награды на текущем шаге.

Например, если стратегия будет слишком часто исследовать, она может упустить возможность получить стабильные награды от рук с высокой известной доходностью. С другой стороны, чрезмерная эксплуатация может привести к тому, что стратегия не найдет лучшую руку.

2.2 Основные алгоритмы

- ϵ -greedy;
- UCB (upper confidence bound);
- Thompson sampling

1) ϵ -greedy: с вероятностью $1 - \epsilon$ выбирает руку с наибольшей текущей наградой (использование). С вероятностью ϵ выбирает случайную руку (исследование). Значение ϵ может уменьшаться с течением времени.

2) UCB проводит своё исследование не случайно, а на основе растущей со временем неопределённости у стратегий. В начале работы алгоритм случайно задействует все стратегии, после чего рассчитывается средняя награда каждой. Далее после каждой итерации обновляются средние награды стратегий. С течением времени чем реже выбиралась та или иная стратегия, тем больше будет у неё неопределённость. Окончательный выбор стратегии — это максимальная сумма средней награды и неопределённости среди всех стратегий.

3) Thompson sampling использует байесовский подход для моделирования распределения награды каждой руки с помощью априорных и апостериорных вероятностей. На каждом шаге генерируется случайное значение из распределения для каждой руки, и выбирается рука с наибольшим значением.

2.3 Метрики

В качестве метрик в нашей задаче мы использовали:

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

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

- total regret (Потери из-за выбора не оптимальных рук)

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

Метрика total rewards была выбрана, поскольку для торгового бота основная цель — максимизация прибыли. Total rewards - метрика, которая показывает, насколько хорошо бот зарабатывает деньги.

Total regret измеряет, насколько меньше вознаграждение, полученное торговым ботом, по сравнению с тем, каким бы оно могло быть, если бы бот покупал каждую акцию по 1, не используя стратегий.. total regret оценивает, насколько хорошо алгоритм справляется с неопределенностью и принимает решения.

3 Датасет

Набор данных о ценах 249 российских акций. Датасет содержит исторические цены открытия, максимума, минимума, закрытия и объема акций, торгуемых на финансовых рынках Российской Биржи. Важной для нас переменной является цена закрытия (close). Взяты почасовые данные для 234 акций в период с 2023-06-01 по 2024-08-27 (остальные не торговались в 2023 и 2024 году).

	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

4 Ход работы

4.1 Strategy

Класс Strategy - основа для реализации различных стратегий многоруких бандитов. Он задаёт общую структуру, которую наследуют конкретные стратегии. Он определяет переменные для хранения состояний и действий.

```

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

```

4.2 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). Сделано для того, чтобы "Забывать" старые про-садки или подъемы, но поднять важность последних колебаний, так как мы торгуем на короткой перспективе (час). Также добавлен параметр "Жадности" и скорости ее возрастания. Его идея в том, что если акции растут в цене медленно, то уверенно уменьшать штраф и постепенно приходить к решению купить медленно-растущую акцию.

```

3 class Thompson(Strategy):
4     def __init__(self, n_arms: int, n_portfolio: int, penalty: float, greed:
      float, greed_tempo: float):
5         super().__init__(n_arms, n_portfolio, penalty, greed, greed_tempo)
6         self.alphas = {i: [1.0] for i in range(n_arms)}
7         self.betas = {i: [-1.0] for i in range(n_arms)}
8         self.greed = greed
9         self.greedness = {i: greed for i in range(n_arms)}
10        self.portfolio = []
11        self.greed_tempo = greed_tempo
12        self.penalty = penalty
13
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

```

4.3 Транзакционные издержки

В коде предусмотрена система штрафов за объём покупок - транзакционные издержки, связанные с реальными рынками. Эти издержки влияют на обновление наград для акций, которые не попадают в портфель. То есть, для акций, которые не вошли в портфель, награда уменьшается на величину штрафа, пропорционального параметру жадности. Высокие транзакционные издержки снижают привлекательность частых изменений в портфеле.

```

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

```

4.4 Создание среды

Среда StockMarketEnv имитирует поведение рынка акций. Внутри среды загружаются данные по нескольким акциям, и для каждой из них рассчитывается доходность (в виде процентного изменения цен закрытия). Если в данных присутствуют пропуски, они заполняются предыдущими значениями. Также для удобства сравнения создаётся массив ожидаемых доходностей, который обновляется на каждой итерации.

```
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
```

4.5 Механизм многорукого бандита

Создан класс Bandit, который объединяет среду и стратегию. На каждой итерации стратегия выбирает портфель из лучших акций на основе текущих параметров. Для каждой выбранной акции извлекается доходность из среды. Стратегия обновляет свои параметры на основе полученных доходностей. Параллельно с этим считается regret.

```
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)
```

4.6 Запуск модели и визуализация результатов

Модель была запущена на всём доступном временном ряде, и на каждой итерации фиксировались total rewards, total regret.

Результаты визуализированы в виде двух графиков:

total regret
total rewards

```
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

```

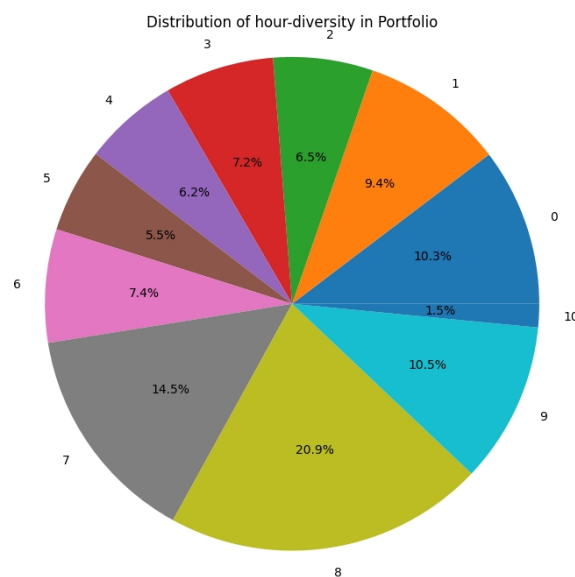
4.7 График распределения часов

Pie-chart, в котором отображено сколько акций в портфеле в этот час осталось из прошлого. Сделано это для того, чтобы отслеживать, что мы действительно заставляем алгоритм торговать реже. Без транзакционных издержек больше 50 процентов всех часов полностью обновляют свой состав акций, а с издержками около 10.

```

1 num_stocks_stayed = [sum(j in strat[i+1] for j in s) for i, s in enumerate(
2     strat[:-2])]
3 stock_counts = []
4 categories = set(num_stocks_stayed)
5 for i in set(num_stocks_stayed):
6     stock_counts.append(sum([j==i for j in num_stocks_stayed]))
7
8 plt.figure(figsize=(8, 8))
9 plt.pie(stock_counts, labels=categories, autopct='%1.1f%%', startangle=0)
10 plt.title('Distribution of hour-diversity in Portfolio')
11 plt.axis('equal')
12 plt.show()

```



5 Результаты

На графике видны заметные пики и впадины, которые говорят о том, что ситуация на рынке довольно изменчива, что очевидно для рынка акций, так как ее стоимость меняется каждую секунду. В процессе многих итераций мы видим, что regret в целом уменьшается, из чего можно сделать вывод о том, что алгоритм лучше адаптируется к условиям задачи по мере работы.

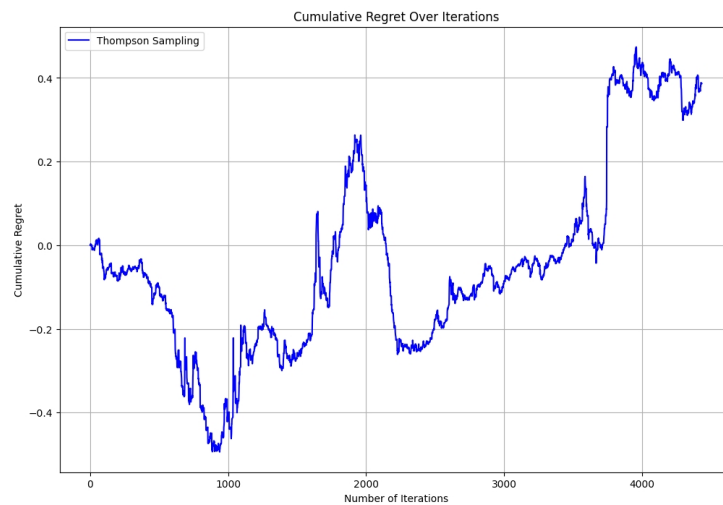
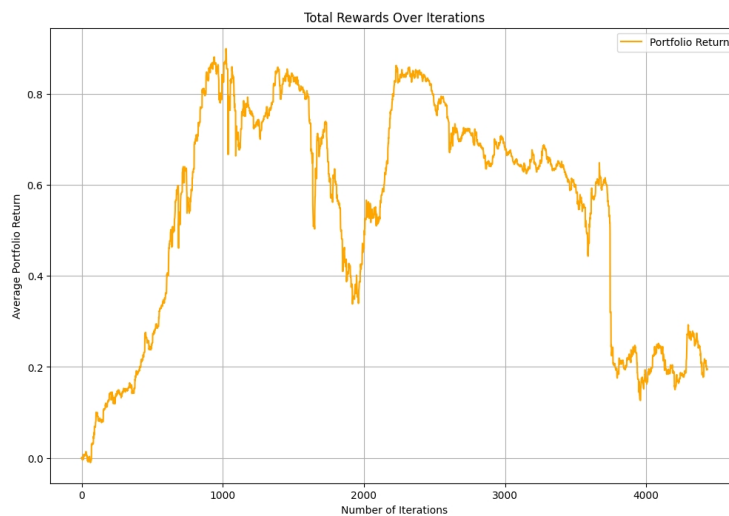


График показывает, что на протяжении всех итераций средняя доходность портфеля в целом увеличивается, это говорит о повышении эффективности алгоритма с течением времени. Колебания показывают, что наш алгоритм сталкивается с изменениями на рынке.



6 Распределение работ

6.1 Клеваичук Александр

- Расчет Thompson Sampling Strategy
- Транзакционные издержки
- Расчет среды StockMarketEnv
- Работа с датасетом
- График распределения часов

6.2 Михеева Екатерина

- Презентация
- Частичное написание теоретических основ и хода работы
- Поиск датасета
- Расчет итоговых значений в коде

6.3 Молвинская Алина

- Оформление файла с описанием работы
- Частичное написание теоретических основ и хода работы
- Поиск датасета
- Расчет класса Strategy и класса Bandit в коде