

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ
РОССИЙСКОЙ ФЕДЕРАЦИИ
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ
ВЫСШЕГО ОБРАЗОВАНИЯ
«Ярославский государственный университет им. П.Г. Демидова»

Курсовая работа

Поиск приближенного равновесия в играх с неполной информацией
02.03.01. Математика и компьютерные науки

Исполнитель: Сычев Р.С.

гр. МКН-41 БО

Руководитель: к.ф.-м.н. Глазков Д.В.

Ярославль 2021

Содержание

Введение	1
1 Первая глава. Описание алгоритмов	2
1.1 Игры в развернутой форме и равновесие Нэша	2
1.2 Контрафактические сожаления и их минимизация	3
2 Вторая глава. Программная реализация алгоритма	6
2.1 Общая схема вычислений	6
2.2 Первый пример. Покер Куна	6
2.3 Второй пример. Домино	7
Заключение	9
Список использованных источников	10

Введение

В последнее время математическая теория игр с неполной информацией находит все большее применение в таких отраслях как теория операций, экономика, кибербезопасность и физическая безопасность. Не в последнюю очередь это происходит благодаря постоянному совершенствованию алгоритмов и увеличению производительности вычислительной техники. Частным случаем таких игр являются игры в развернутой форме. При такой постановке задачи можно близким к естественному способом отразить в игровой форме структуру последовательного принятия решений набором участников в конфликтной ситуации.

Существенным шагом в развитии данного направления является использование алгоритма подсчета сожалений (regret matching). Алгоритм предполагает итеративное вычисление последовательности стратегий в среднем сходящейся к оптимальному стратегическому профилю. Открытие этого метода привело к появлению ряда алгоритмов для поиска приближенного решения в играх с неполной информацией.

Данная работа посвящена рассмотрению одного из популярных в настоящее время итеративных алгоритмов - контрфактической минимизации сожалений (Counterfactual Regret Minimization) и его модификации предусматривающей использование метода монте карло (MCCFR). Данные алгоритмы появились не так давно, но на их основе уже получен ряд недостижимых до этого по сложности результатов.

Целью данной работы является описание и отработка на практике приведенных выше алгоритмов.

В соответствии с темой работы поставлены следующие задачи:

- подготовить теоритическое описание алгоритмов;
- выделить некоторые игры в развернутой форме для последующего решения;
- реализовать алгоритм и произвести расчет стратегического профиля для приведенных игр.

Данная работа может быть интересна людям желающим ознакомиться с некоторыми современными техниками решения игр с неполной информацией.

1 Первая глава. Описание алгоритмов

1.1 Игры в развернутой форме и равновесие Нэша

Игра в развернутой форме представляют компактную общую модель взаимодействий между агентами и явно отражает последовательный характер этих взаимодействий. Последовательность принятия решений игроками в такой постановке представлена деревом решения. При этом листья дерева отождествлены с терминальными состояниями, в которых игра завершается и игроки получают выплаты. Любой нетерминальный узел дерева представляет точку принятия решения. Неполнота информации выражается в том, что различные узлы игрового дерева считаются неразличимыми для игрока. Совокупность всех попарно неразличимых состояний игры называется информационными наборами. Приведем формальное определение.

Определение 1: Конечная игра в развернутой форме с неполной информацией содержит следующие компоненты:

- Конечное множество игроков N ;
- Конечное множество историй действий игроков H , такое, что $\emptyset \in H$ и любой префикс элемента из H также принадлежит H . $Z \subseteq H$ представляет множество терминальных историй (множество историй игры на являющихся префиксом). $A(h) = \{a: (h, a) \in H\}$ — доступные после нетерминальной истории $h \in H$ действия;
- Функция $P: H \setminus Z \rightarrow N \cup \{c\}$, которая сопоставляет каждой нетерминальной истории $h \in H \setminus Z$ игрока, которому предстоит принять решение, либо игрока с представляющего случайное событие;
- Функция f_c , которая сопоставляет всем $h \in H$, для которых $P(h) = c$, вероятностное распределение $f_c(\cdot|h)$ на $A(h)$. $f_c(a|h)$ представляет вероятность выбора a после истории h ;
- Для каждого игрока $i \in N$ \mathcal{I}_i обозначает разбиение $\{h \in H: P(h) = i\}$, для которого $A(h) = A(h')$ всякий раз когда h и h' принадлежат одному члену разбиения. Для $I_i \in \mathcal{I}_i$ определим $A(I_i) = A(h)$ и $P(I_i) = i$ для всех $h \in I_i$. \mathcal{I}_i называют информационным разбиением игрока i , а $I_i \in \mathcal{I}_i$ информационным набором игрока i ;
- Для каждого игрока $i \in N$ определена функция выигрыша $u_i: Z \rightarrow R$. Если для игры в развернутой форме выполняется $\forall z \in Z \sum_{i \in N} U_i(z) = 0$, то такую игру называют игрой с нулевой суммой. Определим $\Delta_{u,i} = \max_{z \in Z} u_i(z) - \min_{z \in Z} u_i(z)$ для диапазона выплат игрока.

Отметим, что информационные наборы могут использоваться не только для реализации правил конкретной игры, но и могут быть использованы для того, чтобы заставить игрока забыть о предыдущих действиях. Игры в которых игроки не забывают о действиях называют играми с полной памятью. В дальнейшем мы бу-

дем рассматривать конечные игры в развернутой форме с нулевой суммой и полной памятью.

Стратегия игрока i — это функция σ_i , которая ставит в соответствие каждому информационному набору $I_i \in \mathcal{I}_i$ вероятностное распределение на $A(I_i)$. Обозначим за Σ_i множество всех стратегий игрока i . Стратегический профиль σ содержит стратегии для каждого игрока $i \in N$. При этом за σ_{-i} обозначим σ без σ_i .

Обозначим за $\pi^\sigma(h)$ вероятность того, что игроки достигнут h руководствуясь σ . Мы можем представить π^σ как $\pi^\sigma = \prod_{i \in N \cup \{c\}} \pi_i^\sigma(h)$, выделяя вклад каждого игрока. В таком случае, $\pi_i^\sigma(h)$ обозначает вероятность принятия совокупности решений игрока i , ведущих от \emptyset к h . Иными словами

$$\pi_i^\sigma(h) = \begin{cases} \prod_{h \sqsubset h' \wedge P(h')=i \wedge h \sqsubset (h',a)} \sigma(h')(a) & \{h'|h \sqsubset h' \wedge P(h')=i\} \neq \emptyset \\ 1 & \text{иначе.} \end{cases}$$

Запись $h \sqsubset h'$ означает, что h' является префиксом h . Обозначим за $\pi_{-i}^\sigma(h)$ вероятность достижения истории h всеми игроками (включая c) за исключением i . Для $I \subseteq H$ определим $\pi^\sigma(I) = \sum_{h \in I} \pi^\sigma(h)$. Аналогично, введем $\pi_i^\sigma(I)$ и $\pi_{-i}^\sigma(I)$.

Ожидаемое значение выплаты для игрока i обозначим как $u_i(\sigma) = \sum_{h \in Z} u_i(h) \pi^\sigma(h)$.

Традиционным способом решения игр в развернутой форме для двух игроков является поиск равновесного профиля стратегий σ , который удовлетворяет следующему условию

$$u_1(\sigma) \geq \max_{\sigma'_1 \in \Sigma_1} u_1(\sigma'_1, \sigma_2) \quad u_2(\sigma) \geq \max_{\sigma'_2 \in \Sigma_2} u_2(\sigma_1, \sigma'_2). \quad (1.1)$$

Такой стратегический профиль называют равновесием по Нэшу. В случае, если стратегический профили σ удовлетворяет условию

$$u_1(\sigma) + \epsilon \geq \max_{\sigma'_1 \in \Sigma_1} u_1(\sigma'_1, \sigma_2) \quad u_2(\sigma) + \epsilon \geq \max_{\sigma'_2 \in \Sigma_2} u_2(\sigma_1, \sigma'_2) \quad (1.2)$$

его называют ϵ – равновесием.

Для рассматриваемых далее алгоритмов наиболее интересен вариант игры с нулевой суммой для двух игроков. Именно для него имеется строгое математическое обоснование сходимости к равновесию нэша.

1.2 Контрафактические сожаления и их минимизация

Минимизация сожалений является популярным концептом, для построения итеративных алгоритмов приближенного решения игр в развернутой форме [ссылка]. Приведем связанные с ней определения. Рассмотрим дискретный отрезок времени T включающий T раундов от 1 до T . Обозначим за σ_i^t стратегию игрока i в раунде t .

Определение 1 Средним общим сожалением игрока i на момент времени T называют величину

$$R_i^T = \frac{1}{T} \max_{\sigma_i^* \in \Sigma_i} \sum_{t=1}^T (u_i(\sigma_i^*, \sigma_{-i}^t) - u_i(\sigma^t)) \quad (1.3)$$

В дополнении к этому, определим $\bar{\sigma}_i^T$ как среднюю стратегию относительно всех раундов от 1 до T . Таким образом для каждого $I \in \mathcal{I}_i$ и $a \in A(I)$ определим

$$\bar{\sigma}_i^T(I) = \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma^t(I)(a)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)}. \quad (1.4)$$

Теорема 1 Если на момент времени T средние общие сожаления игроков меньше ϵ , то σ является 2ϵ равновесием.

Говорят, что алгоритм выбора σ реализует минимизацию сожалений, если средние общие сожаления игроков стремятся к нулю при t стремящимся к бесконечности. И как результат, алгоритм минимизации сожалений может быть использован для нахождения приближенного равновесия по Нэшу, в случае игр двух игроков с нулевой суммой. Однако, стратегии сформированные для игр с большим числом игроков могут также успешно применяться на практике [1].

Понятие контрафактического сожаления служит для декомпозиции среднего общего сожаления в набор дополнительных сожалений, которые могут быть минимизированы независимо, для каждого информационного набора.

Обозначим через $u_i(\sigma, h)$ цену игры с точки зрения истории h , при условии, что h была достигнута, и игроки спользуют в дальнейшем σ .

Определение 2 Контрафактической ценой $u_i(\sigma, I)$ назовем ожидаемую цену, при условии, что информационный набор I был достигнут, когда все игроки кроме i играли в соответствии с σ . Формально

$$u_i(\sigma, I) = \frac{\sum_{h \in I, h' \in Z} \pi_{-i}^{\sigma}(h) \pi^{\sigma}(h, h') u_i(h')}{\pi_{-i}^{\sigma}(I)}, \quad (1.5)$$

где $\pi^{\sigma}(h, h')$ — вероятность перехода из h в h' .

Обозначим за $\sigma^t|_{I \rightarrow a}$ стратегический профиль идентичный σ за исключением того, что i всегда выбирает a в I .

Немедленным контрафактическим сожалением назовем

$$R_{i,imm}^T = \frac{1}{T} \max_{a \in A(I)} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{I \rightarrow a}, I) - u_i(\sigma^t, I)). \quad (1.6)$$

Интуитивно это выражение можно понимать, как аналог среднего общего сожаления в терминах контрафактической цены. Однако вместо рассмотрения всевозможных максимизирующих стратегий рассматриваются локальные модификации

стратегии. Положим $R_{i,imm}^{T,+}(I) = \max(R_{i,imm}^T(I), 0)$ Связь немедленных контрафактических сожалений и общих средних сожалений раскрывает следующая теорема.

Теорема 2 $R_i^T \leq \sum_{I \in \mathcal{I}_i} R_{i,imm}^{T,+}(I)$.

Таким образом, минимизация немедленных контрафактических сожалений минимизирует общие сожаления. В свою очередь минимизация немедленного контрафактического сожаления может происходить за счет минимизации выражений под функцией максимума. Таким образом мы приходим к понятию контрафактического сожаления

$$R_i^T(I, a) = \frac{1}{T} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{I \rightarrow a}, I) - u_i(\sigma^t, I)). \quad (1.7)$$

Контрафактическое сожаление рассматривает действие в информационном состоянии. В свою очередь, для минимизации контрафактических сожалений можно применить алгоритм приближения Блэквела, который, применимо к рассматриваемым сожалениям, приведет к следующей последовательности стратегий

$$\sigma_i^{T+1}(I)(a) = \begin{cases} \frac{R_i^{T,+}(I, a)}{\sum_{a \in A(I)} R_i^{T,+}(I, a)} & \text{если } \sum_{a \in A(I)} R_i^{T,+}(I, a) > 0, \\ \frac{1}{|A(I)|} & \text{иначе.} \end{cases} \quad (1.8)$$

Другими словами, действие выбирается в пропорции соотношения позитивных контрафактических сожалений не выбора этого действия. Обоснование сходимости полученного решения и оценку ее скорости предоставляет следующая теорема.

Теорема 3 Если игроки придерживаются стратегий, заданных выражением (1.8), то $R_{i,imm}^T(I) \leq \Delta_{u,i} \sqrt{|A_i|}/\sqrt{T}$ и следовательно $R_i^T \leq \Delta_{u,i} |\mathcal{I}_i| \sqrt{|A_i|}/\sqrt{T}$, где $|A_i| = \max_{h: P(h)=i} |A(h)|$.

2 Вторая глава. Программная реализация алгоритма

2.1 Общая схема вычислений

В рассмотренных далее примерах рассматривалась вероятностная реализация алгоритма. При использовании данного метода значения случайных событий генерируются перед началом каждой обучающей итерации. Данный подход позволяет сократить объем памяти и ускорить вычисления в некоторых случаях[2]. При реализации примеров были выделены следующие компоненты:

- настройки игры (произвольная параметризация составных частей игры);
- описание правил игры (зависит от настроек);
- модуль с реализацией алгоритма относительно определенных правил и настроек.

Настройки игры, например, по возможности могут включать число игроков, состав костяшек домино и т.п.

Правила игры включают структуру игрового дерева, механизм распределения случайных событий и функцию выплат. Игровое дерево строится с применением узлов – объектов с информацией о истории игры, о игроке и о возможных действиях.

Сам расчет итераций CFR происходит в выделенном модуле, на основе, определенных для каждого конкретного случая, правил игры. Работа алгоритма начинается с создания игрового дерева. Далее происходит расчет заданного числа итераций. После любой итерации можно получить средние стратегии игроков, которые представляют из себя приближенное коррелирующее равновесие.

2.2 Первый пример. Покер Куна

Покер Куна – это максимально упрощенная версия карточной игры покер[5].

Данная игра достаточно проста, чтобы быть решенной аналитически. Правила следующие:

- в игре участвуют 2 игрока;
- игра начинается с раздачи карт игрокам. Всего имеется 3 карты (1, 2 и 3). Каждый игрок получает одну карту. Причем каждый игрок знает свою карту и не знает карту другого игрока;
- по ходу игры игрокам доступны 2 действия «пасс» («п») и «ставка» («с»). Игру начинает первый игрок. Возможны следующие терминальные игровые истории: «пп», «сс», «сп», «псп» и «псс»;
- если терминальная игровая история заканчивается на «сс», то игрок с большей картой получает 2 очка, а игрок с меньшей их теряет;

- если терминальная игровая история заканчивается на «сп», то сделавший ставку игрок получает 1 очко, а спасовавший игрок теряет 1 очко;
- если терминальная игровая история заканчивается на «пп», то игроки получают по 0 очков.

Данная игра удобна для базовой проверки алгоритма CFR и часто служит в качестве примера той или иной реализации. Мы можем смоделировать дерево игры и информационные наборы игроков. После 10^7 итераций алгоритма удалось получить следующий профиль стратегий (Рисунок).

Основная часть кода программы представлена в приложении А.

2.3 Второй пример. Домино

В данной работе в качестве основного объекта исследования была выбрана игра «Домино».

История

Однако, спортивный вариант игры обладает значительной комбинаторной сложностью и было бы трудно хранить в памяти все дерево игры. В связи с этим в данной работе рассматривались некоторые упрощенные варианты.

Из соображений вычислительной сложности целесообразно рассматривать правила игры следующего вида:

- имеется набор из не более чем 10 костяшек домино;
- игроки имеют на руках 2, 3 (размер руки) костяшки;
- игра может происходить с 2-мя, 3-мя или 4-мя игроками;
- находящиеся не на руках костяшки раздаются по мере развития игры
- игру начинает первый игрок
- все костяшки в процессе игры выкладываются в единственную линию
- если ход возможен, то он происходит по обычным правилам;
- в случае, если ход текущего игрока невозможен, то игра на этом заканчивается, и игроки получают выплаты (победитель забирает все очки, и т.к. необходима нулевая сумма, то проигравшая сторона эти очки теряет).

Приведенные выше правила игры позволяют на практике сформировать полное решение по методу MCCFR. Фрагмент кода приложения для расчета стратегий представлен в приложении Б.

Для проверки полученного приложения был проведен ряд тестовых запусков. Первый тест состоял в определении профиля стратегий для случая игры с полной информацией. Был взят набор из четырех костяшек, которые раздавались поровну двум игрокам.

Это крайне простая игровая ситуация, но по ней можно судить о работоспособности в целом. Ниже приведен полученный профиль стратегий (Рисунок).

Для второго теста был выбран набор из шести костяшек домино из которых каждый игрок в начале раунда получал на руки две, а остальные 2 раздавались по ходу игры. Приведем несколько фрагментов полученных стратегий (Рисунок, Рисунок, Рисунок).

Однако, данные примеры не представляют большой комбинаторной сложности. Для проверки производительности был выбран вариант игры на 10 костяшек для двух игроков. Каждый игрок получал на руки по 3 костяшки и оставшиеся 4 раздавались по ходу игры. Под эксплуатированностью стратегий подразумевается возможная выгода оппонента, если он будет менять только свою стратегию. Будем под ней понимать максимальный приближенно рассчитанный разброс выигрыша изменившего свою стратегию игрока по сравнению с оригинальным профилем. Назовем эту величину τ . Ниже представлен график расчетной эксплуатированности стратегий в зависимости от числа обучающих итераций (Рисунок).

Заключение

Реализация современных криптографических алгоритмов, безусловно, требует высокой квалификации разработчика. Это связано с тем, что подобные алгоритмы должны работать как можно более эффективно для обеспечения быстродействия обслуживаемых ими систем. Бывает весьма трудоемко оптимальное распределение переменных по регистрам процессора и уровням кэша. Также, особого рассмотрения требует использование встроенных типов. Подобные проблемы часто бывают узкоспециализированны в зависимости от архитектуры ЭВМ и выливаются в широкий спектр прикладных вопросов.

В данной работе были рассмотрены наиболее общие и надежные высокоуровневые языковые конструкции, призванные сформировать у читателя представление о процессах практического вычисления функции хэширования согласно ГОСТ Р 34.11-2012 и реализации процессов генерации и проверки ЭЦП согласно ГОСТ Р 34.10-2012.

В дополнение к изложенному, в рамках данной работы был разработан ряд приложений для операционной системы Windows с графическим интерфейсом, которые позволяют непосредственно использовать описанные механизмы.

Приложения реализуют:

- вычисление хэш-функции от файла;
- редактирование параметров и ключей схемы цифровой подписи;
- генерацию цифровой подписи;
- проверку цифровой подписи.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. *Brown, Noam*. Supplementary Materials for Superhuman AI for multiplayer poker / Noam Brown, Tuomas Sandholm. — Science First Release DOI: 10.1126/science.aay2400, 11 July 2019.
2. *Marc Lanctot Kevin Waugh, Martin Zinkevich Michael Bowling*. Monte carlo sampling for regret minimization in extensive games. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, Advances in Neural Information Processing Systems 22 / Martin Zinkevich Michael Bowling Marc Lanctot, Kevin Waugh. — MIT Press, Cambridge, 2009, pages 1078–1086.
3. *Martin Zinkevich Michael Johanson, Michael Bowling Carmelo Piccione*. Regret minimization in games with incomplete information. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, Advances in Neural Information Processing Systems 20 / Michael Bowling Carmelo Piccione Martin Zinkevich, Michael Johanson. — MIT Press, Cambridge, 2008, pages 1729–1736.
4. *Hart, Sergiu*. A simple adaptive procedure leading to correlated equilibrium / Sergiu Hart, Andreu Mas-Colell. — Econometrica, 68(5), September 2000, pages 1127–1150.
5. *W., Kuhn H.* "Simplified Two-Person Poker". In Kuhn, H. W.; Tucker, A. W. (eds.). Contributions to the Theory of Games. 1. / Kuhn H. W. — Princeton University Press, 1950, pp. 97–103.