

# Inteligência Artificial de Poker No Limit Holdem

## Pluribus

Pedro Tavares de Carvalho

31 de janeiro de 2022

**Resumo** Nesse projeto, iremos discutir a inteligência artificial de poker chamada *Pluribus* [6], seus detalhes de implementação e design, além de seus impactos na sociedade e no futuro de Inteligência Artificial (IA) em jogos. Abordaremos o poker como um jogo em teoria de jogos, seu processo de modelagem e os desafios existentes ao se tratar de um jogo como o poker em uma inteligência artificial.

**Introdução** O poker é um dos jogos mais jogados do mundo. Com regras simples de se entender, porém difíceis de se dominar, o poker atende jogadores casuais e aficionados da mesma forma.

Apesar de sua popularidade, a revolução das inteligências artificiais, que inclui damas [16], xadrez [8] e mais recentemente GO [17], atingiu diversos outros jogos antes do poker. Muito disso se deve à dificuldade de se modelar uma estratégia ótima em poker.

A história da inteligência artificial no poker é longa, incluindo diversas estratégias, dentre elas algoritmos genéticos [9] e aprendizado profundo [12]. Uma versão mais simplificada do poker, o Limit Texas Holdem, onde os jogadores possuem limites em quanto eles podem apostar, já foi fechada matematicamente [4].

Iremos explorar no que o Pluribus difere das estratégias empregadas até o momento, e o motivo dessas estratégias serem menos eficientes e menos efetivas do que as novas, modelando o poker matematicamente e observando os detalhes que fizeram o Pluribus quebrar a barreira do poker sem limites com 6 jogadores.

**Modelagem do Poker em Teoria de Jogos** O poker pode ser modelado em teoria de jogos como um jogo:

**De informação incompleta** Os agentes do jogo não possuem acesso a quais cartas os seus adversários estão na mão, ou seja, com as informações que ele tem no momento, não é possível determinar com certeza em qual estado de jogo o mesmo está.

**Adversarial** Os agentes jogam com o objetivo de ganhar, sem colaborar entre si.

**Soma zero** Quando um dos agentes joga, uma combinação dos outros perde, tornando o valor total existente no jogo fixo.

**Estocástico** Ações dentro do poker podem produzir estados diferentes, e, consequentemente, valores diferentes.

Em geral, o poker difere de jogos como xadrez e GO por alguns motivos, incluindo o fato de ele ser um jogo de informação incompleta e estocástico, e pelo fato de o mesmo ser jogado com mais de dois jogadores, o que torna a estratégia e a busca dentro da árvore de decisões mais complexa do que um minimax simples.

**Equilíbrio de Nash e Poker GTO** O poker, quando jogado por pessoas reais, é dividido em duas vertentes principais, a vertente explorativa, que tenta se aproveitar de erros na estratégia do seu adversário, e o poker Ótimo Em Teoria dos Jogos (GTO), que tenta maximizar matematicamente a estratégia, tentando atingir um Equilíbrio de Nash [13].

Em poker sem limites, não existe uma estratégia de Nash fechada, o que é feito são tentativas de aproximação dessa estratégia, no caso do poker humano, com conceitos como valor esperado, *independent chip modelling* [10], equidade, odds, entre outros que ajudam a aproximar a melhor ação em cada momento. Já em inteligências artificiais, a técnica mais usada é a de Minimização de Arrependimento Contrafactual (CRM), que atualiza a árvore de decisões tentando diminuir a quantidade de decisões erradas.

**O Pluribus por alto** O Pluribus é composto de três peças principais, abstração de informação, o croqui de estratégia construído por minimização de arrependimento, e o algoritmo de busca em tempo real que cria estratégia em nós não antes explorados.

A estratégia do Pluribus foi construída, principalmente, por jogos com versões de si mesmo. Utilizando versões antigas da sua estratégia como adversários, ele emula jogos e melhora a sua performance. Esse método é muito utilizado em inteligências artificiais, incluindo o xadrez [18] e GO [17]. Com esses jogos, o Pluribus desenvolve uma estratégia base, que é utilizada e atualizada no momento da busca em tempo real, quando o agente está realmente jogando.

**Abstrações de Informação** Para simplificar a busca e melhorar a estratégia da IA, o Pluribus aplica técnicas de simplificação de informação, tanto a nível de apostas possíveis quanto a nível de valor de mão.

**Abstração de Apostas** A simplificação de apostas age tanto no espaço de apostas consideradas pelo agente em sua jogada quanto em apostas de adversários. O processo é uma discretização das apostas, transformando valores específicos de apostas em um conjunto específico de apostas possíveis.

**Abstração de Mãos** Esse processo diminui o espaço de mãos possíveis dentro de um jogo de poker. A idéia é que mão de valores similares, como um *flush* com carta alta em 10 e

um *flush* com carta alta em 9, serão agrupadas em um mesmo conjunto, e ações similares (se não idênticas) serão tomadas em situações que surgem com essas mãos.

Isso não é uma regra, e dependendo da especificidade da situação, a abstração pode ser desconsiderada, e as mãos podem ser tratadas como diferentes pela estratégia.

**Minização de arrependimento contrafactual de Monte Carlo** Para a construção da estratégia base, é utilizada uma técnica chamada de Minimização de Arrependimento Contrafactual de Monte Carlo (MCCRM) [11, 5].

Essa técnica consiste na simulação de jogos da IA contra si mesma, e da revisão das decisões tomadas no jogo completo a fim de minimizar o "arrependimento" do agente, ou seja, aumentando a probabilidade de decisões boas e diminuindo a de decisões ruins.

Diferente da CRM normal, em que a cada geração da estratégia a árvore de decisões completa é atualizada, o MCCRM simula somente um jogo<sup>1</sup> de todos os possíveis, e atualiza as decisões que são geradas nesse jogo somente.

**Busca em jogos de informação incompleta** Diferente de jogos de informação completa, como xadrez e GO, em que você consegue ter certeza do estado em que você está e da estratégia do seu adversário, busca em ambientes de informação incompleta não conseguem se basear apenas na avaliação de um estado na estratégia de Nash, ao menos não de um estado não-folha.

Isso acontece pois a avaliação assume que o adversário não usará uma estratégia diferente no futuro, ou seja, continuará utilizando uma estratégia de Nash ótima, e portanto tomará as mesmas decisões. Isso pode ser visto na figura 1, onde uma busca em um jogo de pedra papel ou tesoura causa uma estratégia que não implica em otimalidade.

Por exemplo, se tivéssemos um algoritmo guloso, uma busca poderia sempre resultar em

<sup>1</sup>Um jogo é uma descida completa na árvore de decisões, até se chegar em um nó folha.

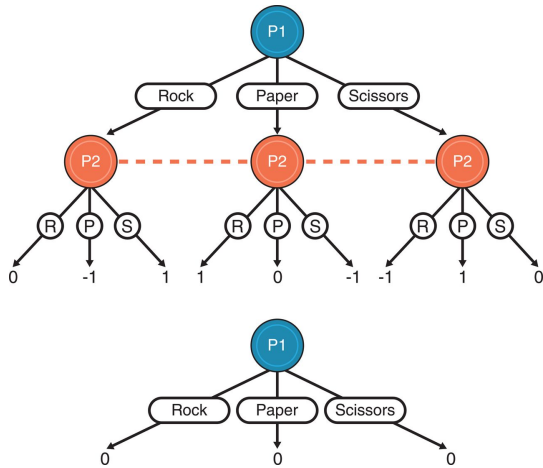


Figura 1: Avaliação e busca em um jogo de pedra papel ou tesoura. Primeiro se vê a descida das decisões e valores possíveis na árvore, depois se vê a avaliação de cada decisão em uma estratégia de Nash.

*pedra* como ação, o que faria com que se um adversário racional mudasse sua estratégia, ele pudesse explorar os nossos vieses.

Para contornar esse problema, o Pluribus aplica uma técnica de pesquisa com profundidade limitada [7] em que diferentes estratégias adversárias são emuladas e avaliadas, observando o valor em todas elas. Para o Plúribus, foram utilizadas 4 estratégias diferentes, sendo essas vieses construídos a partir da *blueprint* de estratégia gerada com o MCCRm.

**Impactos do Pluribus** O Pluribus, apesar da falta de visibilidade, representa uma nova revolução dentro do mundo de Inteligência Artificial para jogos. Nenhum projeto até o momento havia conseguido boas emulações de estratégias de Nash em jogos de informação incompleta com mais de dois jogadores.

Essa modelagem de sistema se aplica em diversas situações da vida real, incluindo negociações de vendas, interações diplomáticas e diversas outras situações em que várias pessoas disputam pelos mesmos recursos.

Dentro do poker, o Pluribus representa um possível problema para o grande cenário online de poker. Esse cenário possui diversos profissi-

onais do jogo, e possui campeonatos de milhões de dólares [15].

Alguns provedores de poker online já possuem ações contra jogadores automáticos de poker [14], porém as técnicas de detecção não são muito claras. Assim como no xadrez online [3, 2], essas técnicas de detecção se desenvolverão mais à medida que mais bots forem desenvolvidos e forem mais acessíveis.

## Glossário

**arrependimento** O quanto um agente preferiria ter tomado uma decisão à decisão real dele.

**equidade** A sua probabilidade de ganhar uma mão. Uma de suas variações é a equidade de fold, que representa a probabilidade do seu adversário foldar com a sua aposta .

**estratégia** Uma descrição de qual ação deve ser tomada em cada estado de um jogo .

**estratégia de Nash** Uma estratégia ótima em um equilíbrio de Nash. Desvios de uma estratégia de Nash, não podem melhorar a performance de um agente.

**flush** Um jogo em poker em que um jogador possui cinco cartas do mesmo naipe.

**independent chip modelling** Conceito de que suas fichas no poker possuem valores diferentes dependendo do momento do jogo em que você se encontra. Aplicado principalmente em torneios de poker.

**limit holdem** Variante do poker em que dois ou mais jogadores jogam com duas cartas na mão e cinco cartas na mesa, e que os jogadores possuem limites de aposta determinados pelo tamanho do *pot* no momento .

**no limit holdem** Variante do poker em que dois ou mais jogadores jogam com duas cartas na mão e cinco cartas na mesa, e que os jogadores não possuem limites de aposta .

**nó** Uma dos possíveis estados de jogo em uma árvore de decisão .

**nó folha** Um nó no fim de uma árvore de decisão, que não possui arestas de saída .

**odds** Quanto de equidade você precisa para, com a aposta atual na mesa, valer a pena você entrar na mão.

**pot** Quantidade de apostas que está em jogo até o momento no poker .

**valor** O valor de um nó na árvore de jogadas.

**valor esperado** Valor esperado ao se fazer uma jogada. Leva em conta a equidade das cartas em mão e no jogo.

**árvore de decisões** Representação de um jogo como uma árvore em que cada nó é um estado do jogo, e cada aresta representa uma decisão que pode ser tomada pelos jogadores.

## Siglas

**CRM** Minimização de Arrependimento Contrafactual.

**GTO** Ótimo Em Teoria dos Jogos.

**IA** Inteligência Artificial.

**MCCRM** Minimização de Arrependimento Contrafactual de Monte Carlo.

## Referências

- [1] Tuomas Sandholm: *Poker and Game Theory*. <https://www.youtube.com/watch?v=b7bStIQovcY>.
- [2] Barnes, D. J. e J. Hernandez-Castro: *On the limits of engine analysis for cheating detection in chess*. Computers & Security, 48:58–73, 2015.
- [3] Bilen, E. e A. Matros: *Online cheating amid COVID-19*. Journal of Economic Behavior & Organization, 182:196–211, 2021.
- [4] Bowling, M., N. Burch, M. Johanson e O. Tammelin: *Heads-up limit hold'em poker is solved*. Science, 347(6218):145–149, 2015.
- [5] Brown, N. e T. Sandholm: *Solving imperfect-information games via discounted regret minimization*. Em *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 1829–1836, 2019.
- [6] Brown, N. e T. Sandholm: *Superhuman AI for multiplayer poker*. Science, 365(6456):885–890, 2019.
- [7] Brown, N., T. Sandholm e B. Amos: *Depth-limited solving for imperfect-information games*. arXiv preprint arXiv:1805.08195, 2018.
- [8] Campbell, M., A. Hoane e F. Hsiung Hsu: *Deep Blue*. Artificial Intelligence, 134(1):57–83, 2002, ISSN 0004-3702. <https://www.sciencedirect.com/science/article/pii/S0004370201001291>.
- [9] Carter, R. G. e J. Levine: *An Investigation into Tournament Poker Strategy using Evolutionary Algorithms*. Em *2007 IEEE Symposium on Computational Intelligence and Games*, pp. 117–124, 2007.
- [10] Gilbert, G. T.: *The independent chip model and risk aversion*. arXiv preprint arXiv:0911.3100, 2009.
- [11] Lanctot, M., K. Waugh, M. Zinkevich e M. Bowling: *Monte Carlo Sampling for Regret Minimization in Extensive Games*. Em Bengio, Y., D. Schuurmans, J. Lafferty, C. Williams e A. Culotta (eds.): *Advances in Neural Information Processing Systems*, vol. 22. Curran Associates, Inc., 2009. <https://proceedings.neurips.cc/paper/2009/file/00411460f7c92d2124a67ea0f4cb5f85-Paper.pdf>.
- [12] Moravčík, M., M. Schmid, N. Burch, V. Lisý, D. Morrill, N. Bard, T. Davis, K. Waugh, M. Johanson e M. H. Bowling: *DeepStack: Expert-Level Artificial Intelligence in No-Limit Poker*. CoRR, abs/1701.01724, 2017. <http://arxiv.org/abs/1701.01724>.
- [13] Nash, J. F. et al.: *Equilibrium points in n-person games*. Proceedings of the national academy of sciences, 36(1):48–49, 1950.
- [14] PokerStars: *PokerStars Prohibited Softwares*. <https://www.pokerstars.com/poker/room/prohibited>.
- [15] PokerStars: *SCOOP*. <https://www.pokerstars.com/poker/tournaments/scoop>.
- [16] Samuel, A. L.: *Some Studies in Machine Learning Using the Game of Checkers*. IBM Journal of Research and Development, 3(3):210–229, 1959.

- [17] Silver, D., A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel e D. Hassabis: *Mastering the game of Go with deep neural networks and tree search*. Nature, 529(7587):484–489, Jan 2016, ISSN 1476-4687. <https://doi.org/10.1038/nature16961>.
- [18] Silver, D., T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T.P. Lillicrap, K. Simonyan e D. Hassabis: *Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm*. CoRR, abs/1712.01815, 2017. <http://arxiv.org/abs/1712.01815>.



## Anexo 1: Mãos de poker

	<b>ROYAL STRAIGHT FLUSH</b> SEQUÊNCIA REAL
	<b>STRAIGHT FLUSH</b> SEQUÊNCIA DO MESMO NAIPE
	<b>FOUR OF A KIND</b> QUADRA
	<b>FULL HOUSE</b>
	<b>FLUSH</b>
	<b>STRAIGHT</b> SEQUÊNCIA
	<b>THREE OF A KIND</b> TRINCA
	<b>TWO PAIRS</b> DOIS PARES
	<b>ONE PAIR</b> UM PAR
	<b>HIGH CARD</b> CARTA ALTA



Figura 2: Ranking de mãos de poker