

オンライン麻雀における実力評価基準の検討

清田 達

アブストラクト

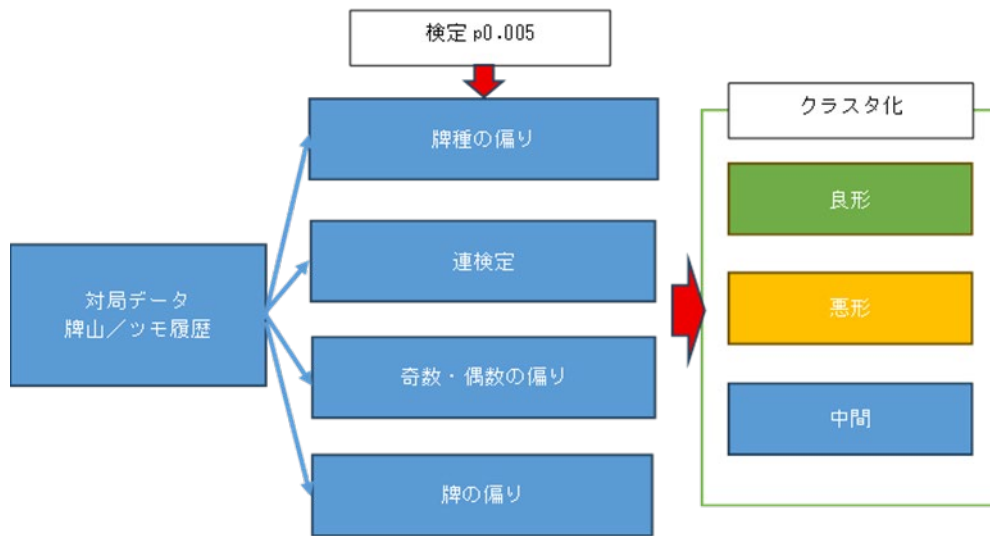
本研究は、オンライン麻雀におけるプレイヤーの実力を公平に評価するための新たな基準を提案する。オンライン麻雀の対局結果は、配牌やツモといった確率的要素に大きく影響されるため、短期的な成績が必ずしも実力を反映するとは限らない。そこで本研究では、対局データ(牌山・ツモ履歴)を複数の統計的指標(エントロピー、ラン検定など)を用いて特徴量化し、教師なしクラスタリングによって局の展開を「良形」「悪形」「中間」に分類する手法を開発した。自身の10対戦分のデータを対象としたケーススタディを通じて、特定の対局において統計的に有意な有利・不利が生じ、プレイヤー評価に影響を及ぼし得ることを実証した。本研究の貢献は、偏りの原因を問うのではなく、「その対局が実力評価の土俵として妥当か」という判定基準を提示する点にある。このアプローチは、eスポーツとしての麻雀の公平性を担保するための基盤となり、将来的には第三者認証を伴う業界標準の牌山生成アルゴリズム確立の必要性を示唆するものである。

キーワード

オンライン麻雀, 実力評価, 統計的検定, クラスタリング, eスポーツ, 公平性

1. 序論

オンライン麻雀においては、対局成績がプレイヤーの段位や勝率に直結し、ひいてはプレイヤーの評価やアイデンティティに強い影響を及ぼす。しかし、各対局には配牌やツモに代表される確率的要素が存在し、短期的にはプレイヤーの実力が正しく反映されない可能性がある。本研究では、このような偶然要因を統計的に分離し、プレイヤーの実力をより公平に評価するための検定手法を提案する。ここで重要なのは、偏りの原因が偶然であれ恣意的であれ、それが基準値から大きく逸脱している場合には、その対局結果を実力評価の根拠とすることが妥当でない、という視点である。本研究はこの「評価の土俵の妥当性」を明示する点に独自性を持つ。研究全体の流れを図1に示す。



※分析データ：某麻雀ゲームソフト中間ランク 30 戦中 10 戦

図 1 分析フロー(p 値による偏り検定と教師なしクラスタ化)

2. 研究の背景

麻雀における実力評価は古くから議論されてきた。既存の研究では、得点効率、和了率、放銃率などの指標が提案されてきたが、それらは必ずしも局面ごとの有利・不利を考慮していない。実力を公平に測定するためには、局面の偏りを統計的に検出し、その影響を排除する手法が必要である。さらに従来の枠組みでは「偶然か恣意か」という原因論に焦点が当てられてきたが、本研究は「その対局結果を実力評価に用いてよいか」という妥当性の判断に着目する。

3. 手法

本研究では、対局(対戦)に使用する牌山データを以下の特徴量を用いて数値化した(大規模シミュレーション 4 万局における各クラスターの平均的特徴は表 1 参照)。牌種の偏り、連検定、偶数・奇数の偏り、牌の偏りをその牌山が持つ特徴量とした。各指標(entropy_bits, p_runs, p_empirical, markov_p, cramers_v)についての詳細は付録 1 に示す。これらを入力として教師なしクラスタリングを行い、対局を「良形」「悪形」「中間」に分類し、Isolation Forest を用いて異常局の検出も行った。シミュレーションとして約 4 万局を生成し、公平な場合の分布基準値を設定した(良形 $\approx 43\%$ 、悪形 $\approx 41\%$ 、中間 $\approx 16\%$ 、図 2 参照)。分析手法の詳細は付録 1 に示す。

表 1 大規模シミュレーションにおける各クラスターの平均的特徴

指標	クラスター0	クラスター1	クラスター2	----	----	----	----
entropy_bits	0.48(低い)	0.67(高い)	0.57(中間)	p_runs	0.08	0.08	0.4
p_empirical	0.78	0.25	0.55	markov_p	0.16	0.17	0.48
				cramers_v	0.14	0.17	0.16

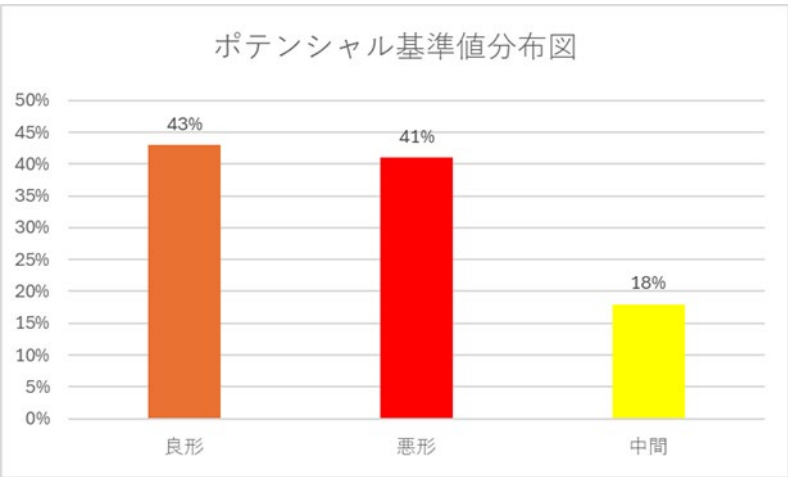


図 2 ポテンシャル基準値分布図

4. ケーススタディ

本研究では、筆者自身の 10 対戦分の対局データを対象として検定ロジックの妥当性を確認した。サンプル数としては限定的であるが、目的は大規模統計ではなく、検定手法が偏りを適切に捉えられるかをケーススタディとして示すことにある。実際に有意な偏りが観測され、短期決戦の結果が実力評価と直結しない可能性を具体的に示すことができた。もっとも、10 対戦規模のデータ数は統計的な一般化を保証するものではなく、外的妥当性を担保するためには、各プラットフォームが自社に蓄積された大規模データベースに対して本検定を適用することが不可欠である。個人研究者が数万局単位のデータを収集することは現実的に困難であるため、今後の実用化に向けては、プラットフォーム側による実装と運用が必須となる。10 対戦の特定対局における各プレイヤーの局パターン遭遇率とベースラインの比較は表 2、この対局の任意に 1 局の各指標値は表 3 に示した。分析で用いた各指標の数式やアルゴリズム設定は付録 1 を参照されたい。

表 2 特定対局における各プレイヤーの局パターン遭遇率とベースラインの比較 | プレイヤー名 | 局数 | 良型ポテンシャル | 悪形ポテンシャル | 中間 | |----| |----|

:--- :--- :---	A さん	8	12.5%	62.5%	25.0%		B さん	8	12.5%	12.5%	75.0%		C さん	8	37.5%	0.0%	62.5%		D さん	8	25.0%	37.5%	37.5%		ベー スライン	40,000	43.2%	41.1%	15.7%	
------------------	------	---	-------	-------	-------	--	------	---	-------	-------	-------	--	------	---	-------	------	-------	--	------	---	-------	-------	-------	--	------------	--------	-------	-------	-------	--

表 3 任意の 1 局の各指標値

	name	p_runs	p_empirical	p_value	entropy_bits	markov_p	cramers_v	リーチ順目	順目	鳴数	ツモ総数	ツモ切り数	ポイント収支
:--- :--- :--- :--- :--- :--- :--- :--- :--- :--- :---	C さん	0.000523	0.9404	0.2885	0.3585	0.0511	0.1339		14	3	17	7	13900
	D さん	0.044888	0.5657	0.4518	0.5837	0.3007	0.1561		15		14	6	-4300
	A さん	0.384978	0.6012	1	0.6252	0.2048	0.1561	12	14		13	8	-4300
	B さん	0.044888	0.4757	0.7757	0.6252	0.4273	0.1561		14		13	2	-4300

5. 関連研究との比較

近年、将棋や囲碁における人工知能研究では、局面評価関数の学習や良形・悪形に関する統計的分析が進展している(例:AlphaGo、AlphaZero など)。これらの研究は主に「最適手の探索」を目的としており、プレイヤーが直面する状況的な不利や偏りの検出を目的とするものではない。本研究は、従来の「最適性」に関する研究とは異なり、対局が実力評価として公平であったかを統計的に検証する枠組みを提示する点に新規性がある。

6. 再現性と手法選択の根拠

本研究で用いたクラスタリングおよび異常検知手法については、以下のように設定した。まず、クラスタ数 $k=3$ はシルエット係数の最大化に基づいて選定した。特徴量として指標(entropy_bits, p_runs, p_empirical, markov_p, cramers_v)を採用したのは、局面の複雑性や配牌・ツモの独立性を定量的に評価するためである。再現性については、著者が分析に用いたプログラムとデータを GitHub リポジトリに公開している[1]。これにより、他の研究者は手順に従って同一の分析を追試することが可能である。クラスタリングの例示結果を図 3 に示した。

- 「良形」(クラスターC・緑)について: クラスターC は、聴牌への近さを示唆する p_empirical で最大値を記録する一方、p_runs や p_value、entropy_bits では最小値を示します。この鋭く尖った形状は、特定の有利な状態へ収束していく、極めて特徴的な局面であることを示しており、我々はこれを『良形』と定義した。

- 「悪形」(クラスターB・赤)について: クラスターB は、 $p_empirical$ が最小値となり、手詰まり感を示す $markov_p$ が最大値を取る。これは『良形』とは正反対の特性であり、不利な局面である『悪形』と解釈するのが妥当である。
- 「中間形」(クラスターA, D・青, 黄)について: クラスターA および D は、複数の指標で 0 と 1 の中間的な値を取る、よりバランスの取れた多角形を描く。これらは『良形』『悪形』のような極端な状況ではなく、対局中に頻出する一般的な局面、『中間形』であると考えられる。

アルゴリズムの詳細設定は付録 1 にまとめてある。

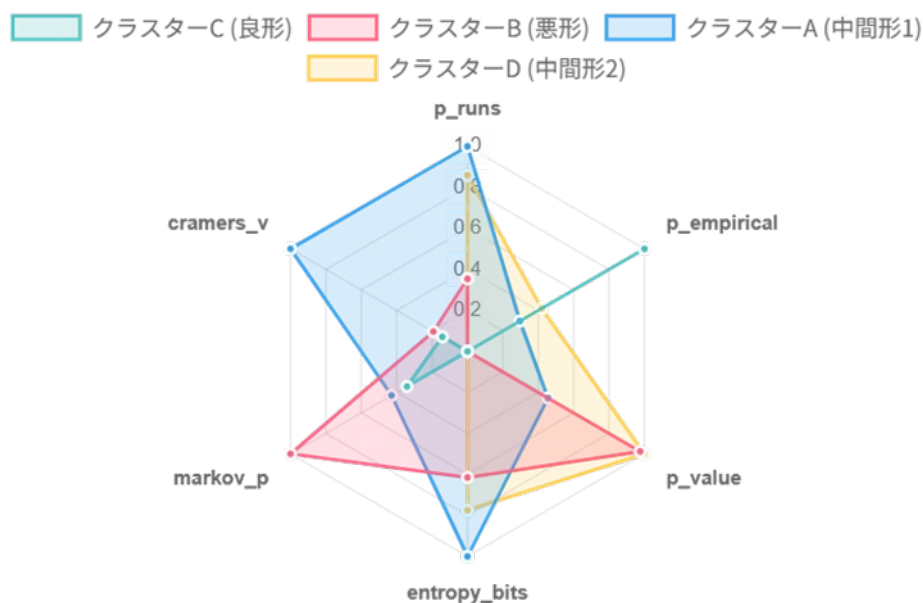


図 3 1 対局データにおける正規化されたクラスターの平均的特徴

7. 社会的意義

本研究の成果は、オンライン麻雀を単なる娯楽としてではなく、e スポーツ競技として捉えた際の公平性評価に新たな視点を提供するものである。具体的には、牌山の極端な偏りをスポーツ競技における「不正な外部要因」に例えることができる。例えば、陸上競技においてドーピングによって一部の選手が不当な優位を得た場合や、競技コース上に障害物や異常な環境条件が存在した場合、記録や順位は「実力の正当な反映」とは認められない。同様に、麻雀においても牌山の状態が統計的に大きく逸脱している場合、その対局結果を基にプレイヤーの実力を評価することは妥当ではない。本研究は「不正を検出する」ことを目的とするものではなく、むしろ「実力を正しく測定できる条件を定義する」ことを目指している。その意味で、スポーツ倫理における

「公正な競技環境」の概念を、e スポーツ麻雀に適用する枠組みを提示したと言える。この視点を踏まえれば、業界においては牌山生成アルゴリズムの第三者認証が不可欠となる。これはスポーツにおけるアンチドーピング規制や競技環境の公平性チェックに相当する役割を果たすものであり、オンライン麻雀を持続可能な競技文化として発展させるための基盤となる。

8. 限界と展望

限界 本研究で対象とした実データは著者自身の 10 対局に限定されており、サンプルサイズが小さい。さらに、検出された偏りは特定のプレイヤーや時期に依存している可能性があるため、一般化には慎重さを要する。

展望 今後は、数百局規模の実データを収集することで統計的な有意性をより強固に検証する予定である。また、他のプレイヤー集団への適用によって、個人差や戦略スタイルの影響を分離することも可能になると考えられる。さらに、麻雀以外の e スポーツやカードゲームにも手法を拡張することで、広く「不利な状況下での実力評価」を定量化できるフレームワークへと発展させられる余地がある。その際には、オンラインゲーム業界全体で共有可能な「牌山生成の業界標準化」と「第三者認証」の導入が、次のステップとして求められるだろう。

9. 結論

本研究は、牌譜データを基に実力評価の公平性を統計的に検定する枠組みを提示した。シミュレーションを参照基準とした分析により、特定の対局において統計的に有意な不利が生じ得ることを示し、これがプレイヤー評価に影響を及ぼし得ることを明らかにした。重要なのは、偏りの原因が偶然か恣意的操作かに関わらず、数値的な基準から大きく逸脱している場合には、その対局は実力を評価する土俵たり得ないという点である。例えば、10 局すべてにおいて他プレイヤーが有利な配牌やツモを得た場合、それは「運が悪い」で済む話ではなく、その局面ではそもそも実力の比較が成立していない。本研究の貢献は、原因の追及ではなく、実力評価に用いることが妥当か否かという新しい判定基準を提示する点にある。さらに本研究は、オンライン麻雀の健全化に向けた業界標準化の必要性を示唆しており、第三者認証を伴う標準的な牌山生成アルゴリズムの確立が、今後の社会的課題となる。今後の拡張研究を通じて、オンラインゲームにおける技能評価の客観性と信頼性を高める基盤の構築に寄与できると考える。

参考文献

[1] Toru1968, "mahjong-performance-analysis", (2024).

<https://github.com/Toru1968/mahjong-performance-analysis> (2025 年 9 月 2 日アクセス).

付録 1. 分析手法の詳細

付録 1.1. 特徴量(指標)の定義

付録 1.1.1. entropy_bits(情報エントロピー)

- 意味: 牌山や手牌の不確実性の度合いをビット単位で表す.
- 使い方: 各牌が出る確率 p_i を計算し、シャノンエントロピーを求める. $H = -\sum_i p_i \log_2 p_i$ エラー! ファイル名が指定されていません.
- 例: 同じ牌ばかり出る偏った牌山は低エントロピー、ランダム性が高い牌山は高エントロピー.
- 検定内容: 牌山の公平性評価、プレイヤー間のランダム性の比較(不正牌山の検出).

付録 1.1.2. p_runs(ラン検定の p 値)

- 意味: 牌の連続パターンがランダムかどうかの統計的指標.
- 使い方: 同じ色の牌(萬子・筒子・索子)が連続して出る回数を観測し、ランダムな順列から期待される連続回数と比較する. p 値が小さいと「連続が偶然ではない」= 偏りありと判断する.
- 検定内容: 牌の偏り検出、偏った順番でゲーム進行していないかの確認.

付録 1.1.3. p_empirical(経験的 p 値)

- 意味: シミュレーションや実データから直接計算した p 値.
- 使い方: 実際の牌山データを n 回シミュレーションし、観測された統計量(例: 特定牌が特定位置に出る頻度)がどのくらい珍しいか評価する.
- 検定内容: 理論的な完全ランダム分布ではなく、実ゲームデータに基づく公平性評価(例: 「このプレイヤーが連続してツモに恵まれているのは偶然か?」).

付録 1.1.4. markov_p(マルコフ確率/p 値)

- 意味: 牌の順序性を考慮した確率.
- 使い方: 牌の出現が直前の牌に依存するかをマルコフモデルで評価する.

- **検定内容:** 牌山生成の偏りを検出、偏った順序が連続するとプレイヤー有利になるかの分析.

付録 1.1.5. cramers_v(クラメールの V)

- **意味:** 2 つのカテゴリ変数の関連度 (0~1).
 - **使い方:** カイ二乗検定から計算する.
-
- **検定内容:** 特定の条件(役満成立・ドラ牌)と牌の出現の依存度を定量化する. 0.5 以上なら強い依存関係、0.1 未満ならほぼ独立と見なす.

付録 1.2. 分析手法の設定

- **クラスタリング:** K-means を用い、クラスタ数 $k=3$ とした. これはシルエット係数の最大化に基づいて選定した.
- **異常検知:** Isolation Forest を利用した.

付録 1.3. 再現性

- 著者が用いた分析データおよびプログラムは GitHub 上で公開されている[1].
- 公開された手順に従うことで、同一の分析を再現可能であり、研究者が追試できる水準の透明性を確保している.

A Study on Performance Evaluation Criteria in Online Mahjong

Toru Kiyota.

Abstract

This study proposes a new framework for fairly evaluating player skill in online mahjong. Game outcomes in online mahjong are heavily influenced by probabilistic factors such as initial hands and tile draws, meaning short-term results do not always reflect true skill. To address this, we developed a method that uses several statistical indicators (e.g., entropy, runs test) to featurize game data (tile walls and draws) and classifies game situations into "Favorable," "Unfavorable," and "Neutral" forms using unsupervised clustering. A case study of our own 10 matches demonstrates that statistically significant advantages and disadvantages can occur

in specific games, potentially affecting player evaluation. The contribution of this research lies not in questioning the cause of bias, but in presenting a new criterion for judging whether a game is a “valid basis for skill assessment.” This approach provides a foundation for ensuring the fairness of mahjong as an e-sport and suggests the future necessity of establishing an industry-standard tile-shuffling algorithm with third-party certification.

Keywords

Online Mahjong, Skill Evaluation, Statistical Test, Clustering, e-Sports, Fairness