

論文 / 著書情報
Article / Book Information

論題	大学入試における合格実績データに系統的な欠損がある場合の大学入試難易度序列の決定手法の提案
Title	Proposal of a method to determine the difficulty level of university entrance examinations when there are systematic defects in the results of university entrance examinations
著者	青木 亮磨, 北澤 正樹, 高橋 聡, 吉川 厚, 山村 雅幸
Authors	Ryoma Aoki, Masaki Kitazawa, Satoshi Takahashi, Atsushi Yoshikawa, Masayuki Yamamura
出典	日本科学教育学会研究会研究報告, Vol. 34, No. 3, pp. 159-164
Citation	JSSE Research Report, Vol. 34, No. 3, pp. 159-164
発行日 / Pub. date	2019, 12

大学入試における合格実績データに系統的な欠損がある場合の 大学入試難易度序列の決定手法の提案

Proposal of a method to determine the difficulty level of university entrance examinations
when there are systematic defects in the results of university entrance examinations

○青木亮磨^{*1}, 北澤正樹^{*1}, 高橋聡^{*2}, 吉川厚^{*1}, 山村雅幸^{*1}

Ryoma AOKI^{*1}, Masaki KITAZAWA^{*1}, Satoshi TAKAHASHI^{*2}, Atsushi YOSHIKAWA^{*1}, Masayuki YAMAMURA^{*1}

^{*1}東京工業大学 情報理工学院, ^{*2}関東学院大学 理工学部

^{*1}Tokyo Institute of Technology, ^{*2}Kanto Gakuin University

【要約】 教育においてデータを活用して教育施策に活用する動きが近年盛んになっており、多くの教育に関するデータが雑誌や本、WEB 上で公開されている。本研究では大学入試に関するデータを用いた大学の入試難易度序列の決定手法を提案する。各高校の公開している大学合格実績データは高校の実力と大学入試の難易度に依存する系統的な欠損である。この系統的な欠損は高校がターゲットとする大学のみ合格率が高い性質を持つ。ここで、高校からの合格率の高い大学の入試難易度が以ていると考えると、合格率の高い部分にのみ注目することで系統的な欠損のあるデータで入試難易度を比較できる。この点から、距離を用いた並び替え手法、合格率を用いた並び替え手法、レーティングを用いた序列決定手法の3つの手法を作成した。得られた序列を塾が公開する偏差値による序列との順位相関で評価した結果、各手法での順位相関は 0.83, 0.85, 0.89 という値になった。今後は精度向上を目指して、同じ順位に分類された大学の更なる序列付けをする手法を検討していく。

【キーワード】 Evidence Based Education, オープンデータ, 系統的な欠損, Colley の手法, Markov の手法, 順位相関, 大学入試難易度, レーティング手法

I. 問題の所在

教育においてデータを活用して教育政策上活かす動きが近年盛んになっている。そのうちの1つにはエビデンス・ベース・エデュケーション (以下 EBE : Evidence Based Education と呼ぶ) がある (卯月, 2018) (玉井, 藤田, 2017)。しかしながら、EBE はエビデンスという厳しい条件があるため、教育倫理に則して入手することが難しいという現状がある。その制約の中でも例えばテストデータであれば、エビデンスとして使用できるようなものがある (戸田, ほか 4 名, 2019)。しかしながら、対象となる受験者に教育効果が異なると想定できる教育手法をデータをとるために適用することが出来ないため、入手するには多大な労力が必要になる (耳塚, 浜野, 2013)。

一方、教育に関わるデータは過去から様々な形でオープンになっているものも多い。例えば、様々な受験雑誌や WEB において、中学受験対策塾の公開する合格実績や中学入試偏差値、同様に、高校受験対策塾や大学受験対策塾の公開する合格実績や高校入試偏差値/大学入試偏差値、各模試での生徒全体の問題正答率、高校の公開する大学合格実績、授業評

価アンケートなどがある。これらは生データではなく集計データであるが、膨大な量がある。これらを利用して、有意な情報を抽出したい。しかしながら、これらのオープンなデータが活用されてこなかったのも理由がある。例えば、模試の結果は受験者という母集団に依存しているため、どの程度学力が伸びたかという推定は、被験者群が異なることに加えてアンカー問題などの等化要素もないため、困難である。

ここで、例えば「学力を高めている学校を抽出する」というような抽象度の高い問題設定をする。この問題を解くためには、各学校における入学時と卒業時の学力レベルを得る必要がある。入学時においては入学時偏差値が仮に使用できるとしても、卒業時の学校の学力レベルはデータとしてない。そこで大学入試の合格実績データからこれを推定できれば、入学時と卒業時の学力から問題に対する一応の回答が得られる。

各高校や受験雑誌が公表しているものに、大学合格者数がある。これらは合格者数であるので、卒業者数が多い学校と少ない学校もあることや、浪人し

て合格する場合なども含まれているため、必ずしもそのまま扱えるデータでは無い。それ以上に問題になるのが、どの大学に合格した方がより学力が高いといえるかという課題である。仮に、学校の指導力がそれほど変わらないとして、浪人は次年度も同等の学力を保持していると考えても全生徒が全大学を受験するわけでは無いので、序列がそもそもつげづら。

詳しく説明すると、高校の公開する大学入試合格実績データには2つの欠損が存在する。1つは個人データである。もし個人データを扱うことが可能であれば、等価手法を用いて大学入試難易度推定や個人能力の推定を行うことが可能である。しかし、そのような個人が特定できるようなデータを手に入れることは困難である。もう1つは収集することが不可能なデータである。大学の比較をする時、大学入試難易度が低ければ全ての高校からの合格者数が多くなるのが理想である。このような合格実績になっていれば大学入試難易度序列も高校実力序列も作ることができる。しかし、実際は異なり、図1のように大学入試難易度が低くても全ての高校からの合格者数が多くなるわけではない。各高校にはターゲットとする大学が存在し、ターゲットとする大学より入試難易度の低い大学はそもそも生徒が受験を行わないため合格者数が少なくなる。このため、全ての問題を生徒が解くことを想定するテストと同じように大学入試合格実績データを扱うことはできない。大学入試合格実績データにはこのように欠損が存在するが、しかし、欠損はランダムではなく、系統的に存在するため、この性質に合わせた手法を作成する必要がある。

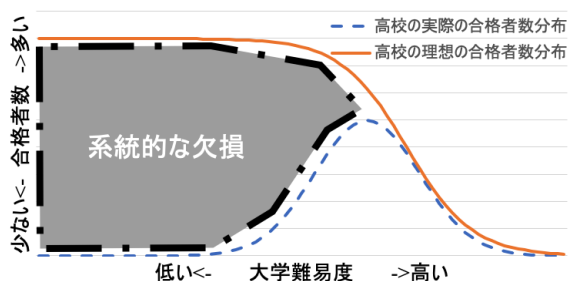


図 1: とある高校の合格者数分布

従来の序列を決定する手法では、対象が直接対決した結果に基づいてポイントを増減させて順位を決める手法(A. E. E, 1978) (M. E. G, 1999) (M. E. G,

2001)か、National Football League で用いられている Passer rating や Entertainment and Sports Programming Network の作成した Total quarterback rating のような全対象者のある項目に関して一定の基準で評価した結果を元に順位を決める手法が用いられている場合が多い。大学合格者数の比較を行おうとしたとき、高校も大学も直接対決をしておらず、また、欠損したデータであるため数値をそのまま比較することは不可能である。よって、従来の序列を決定する手法を用いることはできない。

そこで、本研究では、系統的な欠損を持つデータに対応した序列決定手法として、大学入試合格実績データを用いて大学入試難易度序列を決定する手法を作成することを目標とする。

II. 研究の方法

1. データ

大学入試に関するデータとして、各高校が公開している合格実績、各高校の卒業生数、複数の塾/大学入試対策サイトが公開している大学の偏差値データを使用する。この内、各大学への合格実績データと各高校の卒業生数のデータを大学の序列の決定に用いる。今回は地域性を考慮するため、東京都、千葉県、埼玉県、神奈川県にある大学/高校のみのデータを用いている。データを入手できた高校は886校、大学は77校となった。

2. 手法の提案と評価手法

高校の規模に差があるため、大学合格者数をそのまま比較するのではなく卒業生数で割ることで高校の規模の正規化を行う。そして、大学も規模に差があるため、卒業生数で正規化した卒業生数をさらに正規化したものを大学合格率とする。各高校からの各大学への合格率 $P = \{p_{1,1}, p_{1,2}, \dots, p_{m,n}\}$ は以下で定義する。

$$p_{ij} = \frac{q_{ij}}{g_i} \div \sum_k \frac{q_{kj}}{g_k} \quad (1)$$

ここで、 m は高校の数、 n は大学の数、 p_{ij} は高校 i から大学 j への合格率、 q_{ij} は高校 i から大学 j への合格者数、 g_i は高校 i の卒業生数を表す。

2.1 評価手法

大学入試の難易度は模擬試験などによる偏差値がある。例えば、3社の塾/大学入試対策サイトが公開している偏差値の度数分布を図2に示す。ここで、

縦軸は該当する偏差値に該当する大学の数の割合である。これを見てわかるように、各社とも偏差値が一致しているわけでは無い。偏差値は平均 50 標準偏差 10 の正規分布に従う必要があるが、大学の偏差値はその大学に合格するために必要な偏差値の目安であるため、正規分布に従っている必要はない。

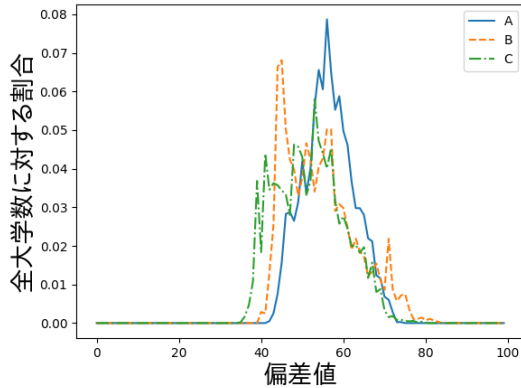


図 2: 塾の偏差値の度数分布

しかしながら、それぞれの偏差値から得られる大学序列のスピアマンの順位相関係数をとると表 1 のようになり、偏差値 A と偏差値 B、偏差値 B と偏差値 C、偏差値 C と偏差値 A の順位相関は 0.98 と非常に高い。そこで、本研究で得られる大学入試難易度序列の評価は、各偏差値から得られる序列とのスピアマンの順位相関を取るものとする。

表 1: 塾サイトの順位相関

	A	B	C
A	1.00	0.98	0.98
B	0.98	1.00	0.98
C	0.98	0.98	1.00

2.2 並び替えによる序列決定

大学の並び替えを行うことで大学入試難易度の序列を作成する手法である。本研究では、東京大学と東京医科歯科大学を大学の中で最も入試難易度の高い大学とした。

2.2.1 距離による並び替え

各高校からの合格率を各要素としたベクトルとして考えた時の大学間の距離を用いた並び替えを行う。作成した手法のアルゴリズムは以下のとおりである。

1. 大学群 U (初期値は東京大学と東京医科歯科大学) を順位 i (現在のループ回数) とする
2. 大学群 U からの距離を変数として、順位未決定の大学を X-means を用いてクラスタリ

ングする

3. 大学群 U から各クラスターへの距離の平均が最も小さいクラスターに所属する大学を新たな大学群 U とする
4. 順位未決定の大学が無くなった場合終了。順位未決定の大学が存在する場合は 1 へ

今回は距離として Canberra distance, Manhattan distance, Euclidean distance, Hamming distance, Pearson product-moment correlation coefficient を用いた。距離ごと、試行ごとに結果が異なるため、距離ごとに 4000 回のシミュレーションを行い、大学のスコア $S = \{s_1, s_2, \dots, s_n\}$ を以下の式から求める。

$$s_i = \sum_l^T \sum_j^{4000} \frac{r_{ij,l}}{\max_k(r_{i,k,l})} \quad (2)$$

ここで、 s_i を大学 i の評価値、 T を用いた 5 つの距離の集合、 $r_{i,j,k} \in \mathbb{N}$ を距離 k を用いた時のシミュレーション j における大学 i の順位とする。求めたスコアを昇順に並べることで、大学の序列が得られる。

系統的な欠損は高校がターゲットとしない大学を受験しないことから発生するものである。高校のターゲットとなる大学は合格率が高く、入試難易度が似ていると考えることができる。同じ高校からターゲットにされている、もしくは同じ高校からターゲットにされていない 2 つの大学は大学間の距離が短くなるという性質を用いることで、系統的な欠損に対応することができる。この手法は高校の並び替えは行わずに大学の並び替えのみを行うこと、使用する距離により結果が大きく異なること、クラスタリングを用いることにより複数の大学が同じ順位になることが特徴である。

2.2.2 大学合格率による並び替え

各高校からの各大学への合格率を用いた並び替えを行う。大学と高校を交互に難易度/実力の高い順番に並べていく手法である。作成したアルゴリズムは以下のとおりである。

1. 大学群 U (初期値は東京大学と東京医科歯科大学) を順位 i (ループ回数) とする
2. 大学群 U への合格率の合計が高い高校を $|U| * m/n$ だけ選択し、これを高校群 H とする
3. 高校群 H からの合格率を用いて、順位未決定の大学を X-means を用いてクラスタリング

する

4. 高校群 H からの合格率の平均が最も高いクラスに所属する大学を新たな大学群 U とする
5. 順位未決定の大学が無くなった場合終了. 順位未決定の大学が存在する場合は 1 へ

試行ごとに結果が異なるため, 4000 回のシミュレーションを行い, 大学のスコア $S = \{s_1, s_2, \dots, s_n\}$ を以下の式から求める.

$$s_i = \sum_j \frac{r_{i,j}}{\max_k(r_{i,k})} \quad (3)$$

ここで, s_i を大学 i の評価値, $r_{i,j} \in \mathbb{N}$ をシミュレーション j における大学 i の順位とする. 求めたスコアを昇順に並べることで, 大学の序列が得られる.

大学 A をターゲットとする高校は大学 A より入試難易度の高い大学や低い大学をターゲットとしないため合格率が低く, 大学 A や大学 A に似た入試難易度の大学への合格率が高いという系統的な欠損の性質を利用し, 高校群からの合格率が高い (= 入試難易度が似ている) 大学を複数選択することで系統的な欠損に対応することができる. この手法では高校と大学にそれぞれ序列ができること, 2.2.1 と同じく複数の大学が同じ順位になることが特徴である.

2.3 レイティングを用いた序列決定

この手法でも 2.2.1 と 2.2.2 と同様に東京大学と東京医科歯科大学を大学入試難易度が最も高い大学とした. 各大学ごとの高校のレイティングは Colley の手法と Markov の手法の 2 通りの手法を用いる. また, Markov の手法を用いた場合のレイティングでは, 対戦を行った 2 校のうち合格率の低い高校が合格率の高い高校に 1 票を投票する場合と, 合格率の低い高校が合格率の高い高校に合格率の差を投票する場合の 2 通りを用いた. 作成したアルゴリズムは以下のとおりである.

1. 各大学において, 合格率をそのまま得点と考え, 各高校のレート $A = \{a_{1,1}, a_{1,2}, \dots, a_{m,n}\}$ を決定する
2. 順位決定済み大学 $D = \{d_1, d_2, \dots, d_u\}$ (初期値は東京大学と東京医科歯科大学) のレート $U = \{u_1, u_2, \dots, u_n\}$ を重みとして高校のレート $H = \{h_1, h_2, \dots, h_m\}$ を更新
3. 各高校のレート H を基準に大学のレート U

を更新

4. 順位未決定の大学のうち, 最もレートの高い大学を順位決定済み大学集合へ追加
5. 全ての大学が順位決定済みになった場合 6 へ. 順位未決定の大学が存在する場合は 2 へ
6. 順位決定済み大学集合に格納された順番を大学入試難易度順位とする

高校のレート $H = \{h_1, h_2, \dots, h_m\}$ と大学のレート $U = \{u_1, u_2, \dots, u_n\}$ は以下のように定義する.

$$h_i = \sum_j^D a_{i,j} * u_j \quad (4)$$

$$u_i = \sum_{j=1}^m h_j * p_{j,i} \quad (5)$$

ここで, $p_{i,j}$ は数式(1)で定義した合格率, $a_{i,j}$ は手順 1 で作成した大学 j における高校 i のレート, D は順位決定済み大学集合である. また, 大学のスコア $S = \{s_1, s_2, \dots, s_n\}$ を以下の式から求める.

$$s_i = \sum_j^V \frac{r_{i,j}}{\max_k(r_{i,k})} \quad (6)$$

ここで, s_i を大学 i の評価値, V を各大学毎の高校のレイティングに用いた 3 つの手法の集合, $r_{i,j} \in \mathbb{N}$ を各大学毎の高校のレイティングに用いた手法 j における大学 i の順位とする. 求めたスコアを昇順に並べることで, 大学の序列が得られる.

大学のレートを更新するとき高校のレートにその高校からの合格率を重みとしてレートを算出し, 高校のレートを更新する時に大学のレートにその大学における高校のレートを重みとしてレートの算出する. 大学をターゲットとしている高校の影響を大きくすることで, 大学 A をターゲットとする高校は大学 A より入試難易度の高い大学や低い大学をターゲットとしないため合格率が低く, 大学 A や大学 A に似た入試難易度の大学への合格率が高いという系統的な欠損に対応することができる.

III. 結果

各手法で得られた大学入試難易度序列と各塾/大学受験対策サイトの偏差値を降順に並べた序列のスピアマンの順位相関で評価を行った. 結果を表 2 に示す. 偏差値を降順に並べた序列と距離による並び

替え手法によって作成された入試難易度序列との順位相関係数は0.82~0.84, 合格率による並び替え手法によって作成された入試難易度序列との順位相関係数は0.85と、僅かに合格率による並び替え手法によって作成された入試難易度序列の方が高いが、ほぼ同じく高い数値となった。レーティングを用いた序列決定手法によって作成された入試難易度序列と偏差値を降順に並べた序列との順位相関係数が最も高く、0.88~0.90となっている。

表 2: 各手法による大学入試難易度序列の評価結果

	距離	合格率	レーティング
A	0.84	0.85	0.90
B	0.82	0.85	0.88
C	0.84	0.85	0.89

3.1 距離関数の影響

距離による並び替え手法は採用する距離によって序列が大きく異なる。採用する距離関数それぞれの序列と各塾/大学受験対策サイトの偏差値を降順に並べた序列との順位相関を取ると表 3 となる。

表 3: 並び替え手法の各距離関数の評価結果

	Canberra	Manhattan	Euclidean	Hamming	Pearson
A	0.26	0.82	0.59	0.22	0.88
B	0.26	0.79	0.56	0.23	0.86
C	0.26	0.81	0.58	0.23	0.86

Pearson product-moment correlation coefficient を用いた序列が最も順位相関係数が高く、次いで Manhattan distance を用いた場合の順位相関係数が高い。Pearson product-moment correlation coefficient と Manhattan distance は比較する 2 つの大学への高校の合格率の傾向が似ていると短くなるような距離である。

Canberra distance を用いた序列は順位相関係数が 0.26 と低く、Hamming distance を用いた序列も順位相関が 0.22~0.23 となっており、順位相関がほぼない。Canberra distance は Manhattan distance と異なり正規化を行なっているため、比較する合格率同士が小さくても影響が大きくなる。Hamming distance は比較する 2 つの大学の合格率の差の大きさにかかわらず、2 つの合格率の同じ桁を比較し、数値が異なっていた場合に距離が一定値だけ増えるため、そもそも数値の大小を比較することに向いていない距離である。

Euclidean distance を用いた序列は順位相関が 0.56~0.59 と高くも低くもない数値となっている。このような数値になった原因として多くの大学が同じ順位になっていることが考えられる。試行回数多くの場合、東京大学と東京医科歯科大学の次に東京工業大学や一橋大学、早稲田大学、慶應義塾大学、MARCH(明治大学、青山学院大学、立教大学、中央大学、法政大学)、日東駒専(日本大学、東洋大学、駒澤大学、専修大学)など 40 以上の大学が同率順位となっており、細かい順位づけができていない。

3.2 大学合格率による並び替え

大学合格率による並び替えによって作成された序列は順位相関係数がいずれも 0.85 となっている。これは、複数の大学が同じ順位に分類されることがあるためである。今回用いている大学は 77 校であり、試行回数うち多くの場合 6 ランクの順位がついた。一番上のランクには基準としている東京大学と東京医科歯科大学を据えているため、残りの 5 ランクに 75 校が分類されている。上から 3 ランク目に MARCH を中心とした 15 校程度、上から 4 ランク目に日東駒専を中心とした 20 校程度、上から 5 ランク目には大東亜帝国(大東文化大学、東海大学、亜細亜大学、帝京大学、国士舘大学)を中心とした 30 校程度が分類されている。このように同じランクに多くの大学が分類されてしまっており、順位相関が低くなっている。

3.3 レーティング手法の影響

それぞれのレーティング手法において作成した大学入試難易度序列との順位相関係数を表 4 に示す。

表 4: 各レーティング手法の評価結果

	Colley	Markov_1 票	Markov_合格率差
A	0.91	0.89	0.88
B	0.89	0.87	0.85
C	0.90	0.88	0.87

Colley の手法を用いた場合も Markov の手法を用いた場合も各高校のレートは 0 以上 1 以下の値になるが、Colley の手法によって算出されるレートは合計しても 1 になるとは限らない一方、Markov の手法によって算出されるレートは合計すると 1 になる違いがある。このため、Markov の手法によって算出されるレートは合格率の高い高校と合格率の低い高校での差が Colley の手法を用いた場合より小さくなる。

このため、Markov の手法の方が順位相関が低くなったと考えられる。

IV. 考察

本研究では、各高校における大学合格実績のデータを用いた。高校も大学もレベルに幅があり、高校それぞれにターゲットとする大学があり、ターゲットとする大学の合格率が高く、それ以外の大学の合格率が低いという系統的な欠損があった。この系統的な欠損に対応するため、とある高校から合格率の高い大学は入試難易度が似ていると考え3つのアプローチを行った。

距離による並び替え手法では合格率を要素とした距離を用いて合格率の高い部分だけに注目し、大学間の距離が短ければ似た入試難易度の大学であると考え、使用することで系統的な欠損に対応した。使用する距離関数に Pearson product-moment correlation coefficient や Manhattan distance などの比較する2つの大学への高校の合格率の傾向が似ていると短くなる距離を用いることで系統的な欠損に対応できたと考えられる。

合格率による並び替え手法では、距離による並び替え手法と同じく、合格率の高い部分だけに注目し、複数の高校からの合格率が高い大学は似た入試難易度の大学であると考え、使用することで系統的な欠損に対応した。この手法では前回のループによって決定された高校を基準に選択する大学を決定しており、その高校は前回のループによって決定された大学数に全大学数に対する全高校数の比率を掛け合わせたものだけ選ばれる。本研究のデータは大学 77 校、高校 886 校を用いており、大学 1 校につき高校が $886 / 77 \approx 11$ 校選択される。このように、序列を決定するために使用する要素が、序列を決定したい要素より多い場合は、系統的な欠損に対応しながら安定した序列が得られると考えられる。

レーティングを用いた序列決定手法では各大学において高校のレーティングを行い、合格率を用いて大学のレーティングを行うアルゴリズムにし、その大学をターゲットとしている高校と大学の影響を強め、ターゲットとしていない高校と大学の影響を弱めることで系統的な欠損に対応した。順位相関が作成した3つの手法のうち最も高く出ているが同率順位が出ない手法のため、同率順位を多く出した場合には向かないと考えられる。

V. おわりに

本研究では、各高校の大学入試合格実績データという系統的な欠損のあるデータに対応するため、とある高校から合格率の高い大学は入試難易度が似ていると考え、大学入試難易度序列として距離による並び替え手法、合格率による並び替え手法、レーティングによる序列決定手法の3つの手法を提案した。距離による並び替え手法では順位相関が 0.82~0.84、合格率による並び替え手法では順位相関が 0.85、レーティングによる序列決定手法では順位相関が 0.88~0.90 になった。

距離による並び替え手法や合格率による並び替え手法では同じランクに多くの大学が分類されているため、このランク内で更なる序列決定を行うことによって精度を高めることを今後の課題とする。

文献

- 卯月由佳(2018)：エビデンスの広がりと言われる教育政策 社会情緒的スキルの教育と調査をめぐる欧米の動向から、社会と調査, 21, 20-28
- 玉井航太, 藤田英典(2017)：エビデンスに基づく教育のための縦断データの解析方法, 国際基督教大学学報 I-A 教育研究, 59, 5-16
- 戸田綾佳, 久野弘暉, 高橋聡, 山村雅幸, 吉川厚(2019)：学習における評価の時系列変化に着目した複数科目間の評価の関係, 日本科学教育学会, 43, 421-424.
- 耳塚寛明, 浜野隆(2014)：全国学力・学習状況調査(きめ細かい調査)の結果を活用した学力に影響を与える要因分析に関する調査研究, 平成25年度「学力調査を活用した専門的な課題分析に関する調査研究」, 国立大学法人お茶の水女子大学
- Arpad E. Elo(1978)：The Rating Of Chess Players Past & Present, Arco Publishing
- Mark E. Glickman(1999)：Parameter estimation in large dynamic paired comparison experiments, Royal Statistical Society, 48, 3, 377-394
- Mark E. Glickman(2001)：Dynamic paired comparison models with stochastic variances, Journal of Applied Statistics, 28, 6, 673-689