Skill Assessment Versus Prediction in Game Play and Cheating Detection

Kenneth W. Regan University at Buffalo (SUNY)

Oxford University, 3 November 2015

うして ふゆう ふほう ふほう ふしつ

• Skill Assessment: how well people did.



・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

- Skill Assessment: how well people did.
- Prediction: how well people will do.

- Skill Assessment: how well people did.
- Prediction: how well people will do.
- Both: how unusual is how well some person did?

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Skill Assessment: how well people did.
- Prediction: how well people will do.
- Both: how unusual is how well some person did?
- Meta: Is this performance really by this person?

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Skill Assessment: how well people did.
- Prediction: how well people will do.
- Both: how unusual is how well some person did?
- Meta: Is this performance really by this person?
- Chess cheating detection needs both and more.

• E-Doping means cheating with computer assistance.



- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

• Applies to online games in much greater volume than chess.

- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.

うして ふゆう ふほう ふほう ふしつ

- Applies to online games in much greater volume than chess.
- Person X cannot cycle up that hill that fast."

- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.
- Applies to online games in much greater volume than chess.
- "Person X cannot cycle up that hill that fast."
- Person X cannot make a champion spin and jump and shoot so fast and accurately.

- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.
- Applies to online games in much greater volume than chess.
- "Person X cannot cycle up that hill that fast."
- Person X cannot make a champion spin and jump and shoot so fast and accurately. versus:

- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.
- Applies to online games in much greater volume than chess.
- "Person X cannot cycle up that hill that fast."
- Person X cannot make a champion spin and jump and shoot so fast and accurately. versus:
- "Person X has hematocrit > 50%."

- E-Doping means cheating with computer assistance.
- Jan. 2013: Lance Armstrong (cycling) and Borislav Ivanov (chess) in news at same time.
- Applies to online games in much greater volume than chess.
- "Person X cannot cycle up that hill that fast."
- Person X cannot make a champion spin and jump and shoot so fast and accurately. *versus:*
- Person X has hematocrit > 50%."
- "Person X made moves highly similar to Code Patch Y."

• Long history, worldwide competitions.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

- Long history, worldwide competitions.
- Game data readily and publicly available.

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

- Long history, worldwide competitions.
- Game data readily and publicly available.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

• Game data is precise

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ● □ ● ●

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).
- Computers play much better than best humans, which is awful!

うして ふゆう ふほう ふほう ふしつ

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).
- Computers play much better than best humans, which is great! since we can generate huge amounts of authoritative analysis data.

うして ふゆう ふほう ふほう ふしつ

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).
- Computers play much better than best humans, which is great! since we can generate huge amounts of authoritative analysis data.

うして ふゆう ふほう ふほう ふしつ

• Chess—much more than Go for instance—lends itself to robust numerical evaluation.

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).
- Computers play much better than best humans, which is great! since we can generate huge amounts of authoritative analysis data.
- Chess—much more than Go for instance—lends itself to robust numerical evaluation.
- Chess move options are *discrete*, hence closer to applications like *multiple-choice tests*.

うして ふゆう ふほう ふほう ふしつ

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).
- Computers play much better than best humans, which is great! since we can generate huge amounts of authoritative analysis data.
- Chess—much more than Go for instance—lends itself to robust numerical evaluation.
- Chess move options are *discrete*, hence closer to applications like *multiple-choice tests*.
- Both chess and online games foster notions of difficulty.

- Long history, worldwide competitions.
- Game data readily and publicly available.
- Game data is precise (except for time taken on each move?).
- Computers play much better than best humans, which is great! since we can generate huge amounts of authoritative analysis data.
- Chess—much more than Go for instance—lends itself to robust numerical evaluation.
- Chess move options are *discrete*, hence closer to applications like *multiple-choice tests*.
- Both chess and online games foster notions of difficulty.
- Chess seems better for notions of depth.

• Skill Assessment in One Number.



• Skill Assessment in One Number. "I'm a 2370."

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.

・ロト ・ 日 ・ モー・ モー・ うへぐ

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

• Difference of 200 \approx 75% expectation for higher player,

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.
- Difference of 200 \approx 75% expectation for higher player,
- Predictive content: your rating is the current best estimate of how you will perform in the next tournament.

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.
- Difference of 200 \approx 75% expectation for higher player,
- Predictive content: your rating is the current best estimate of how you will perform in the next tournament.

• **TPR**: Tournament Performance Rating.

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.
- Difference of 200 \approx 75% expectation for higher player,
- Predictive content: your rating is the current best estimate of how you will perform in the next tournament.
- TPR: Tournament Performance Rating.
- Rating and TPR based only on results of games and ratings of opponents.

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.
- Difference of 200 \approx 75% expectation for higher player,
- Predictive content: your rating is the current best estimate of how you will perform in the next tournament.
- TPR: Tournament Performance Rating.
- Rating and TPR based only on results of games and ratings of opponents.
- Indeed relatively few games: 100 in a year is a lot for pro and amateur alike.

- Skill Assessment in One Number. "I'm a 2370."
- Number has no absolute meaning—only rating differences matter.
- Difference of 200 \approx 75% expectation for higher player,
- Predictive content: your rating is the current best estimate of how you will perform in the next tournament.
- TPR: Tournament Performance Rating.
- Rating and TPR based only on results of games and ratings of opponents.
- Indeed relatively few games: 100 in a year is a lot for pro and amateur alike. Compare to 1,200 being a common need for a good election poll.

Elo Rating Examples

• Bobby Fischer hit 2800 on the US Chess Federation's Elo tabulation, 2785 on the FIDE list in July 1972.

・ロト ・ 日 ・ モー・ モー・ うへぐ

Elo Rating Examples

- Bobby Fischer hit 2800 on the US Chess Federation's Elo tabulation, 2785 on the FIDE list in July 1972.
- Current world champion Magnus Carlsen broke Garry Kasparov's record of 2851, reached 2882 a year ago. Now 2850.

ション ふゆ マ キャット キャット しょう

Elo Rating Examples

- Bobby Fischer hit 2800 on the US Chess Federation's Elo tabulation, 2785 on the FIDE list in July 1972.
- Current world champion Magnus Carlsen broke Garry Kasparov's record of 2851, reached 2882 a year ago. Now 2850.

• Current world #42 has 2702, world #100 has 2653.
Elo Rating Examples

- Bobby Fischer hit 2800 on the US Chess Federation's Elo tabulation, 2785 on the FIDE list in July 1972.
- Current world champion Magnus Carlsen broke Garry Kasparov's record of 2851, reached 2882 a year ago. Now 2850.
- Current world #42 has 2702, world #100 has 2653.
- Formal "Master" designation for USCF is 2200; "FIDE Master" is a formal *title* (IMHO) more typical of 2300.

Elo Rating Examples

- Bobby Fischer hit 2800 on the US Chess Federation's Elo tabulation, 2785 on the FIDE list in July 1972.
- Current world champion Magnus Carlsen broke Garry Kasparov's record of 2851, reached 2882 a year ago. Now 2850.
- Current world #42 has 2702, world #100 has 2653.
- Formal "Master" designation for USCF is 2200; "FIDE Master" is a formal *title* (IMHO) more typical of 2300. Likewise "International Master" ≈ 2400 , *Grandmaster* ≈ 2500 , "strong GM" ≈ 2600 .

Elo Rating Examples

- Bobby Fischer hit 2800 on the US Chess Federation's Elo tabulation, 2785 on the FIDE list in July 1972.
- Current world champion Magnus Carlsen broke Garry Kasparov's record of 2851, reached 2882 a year ago. Now 2850.
- Current world #42 has 2702, world #100 has 2653.
- Formal "Master" designation for USCF is 2200; "FIDE Master" is a formal *title* (IMHO) more typical of 2300. Likewise "International Master" ≈ 2400 , *Grandmaster* ≈ 2500 , "strong GM" ≈ 2600 .
- USCF uses 2000-2199 = "Expert," 1800-1999 = "Class A," 1600-1799 = "Class B" and so on.

• Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

• Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

• László Mérő formalized the 75%-gap as a "Class Unit."

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.
- From 600 to 2800 gives chess a human depth of 11. 8×8 checkers estimated at 10, backgammon and bridge similarly.

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.
- From 600 to 2800 gives chess a human depth of 11. 8×8 checkers estimated at 10, backgammon and bridge similarly.

• Shogi (Japanese chess) at 14, Go at least above 20, maybe 25?

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.
- From 600 to 2800 gives chess a human depth of 11. 8×8 checkers estimated at 10, backgammon and bridge similarly.
- Shogi (Japanese chess) at 14, Go at least above 20, maybe 25?
- Chess computer programs (called *engines*) on desktop PC hardware reach 3200 on all rating lists, 3380 on CCRL.

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.
- From 600 to 2800 gives chess a human depth of 11. 8×8 checkers estimated at 10, backgammon and bridge similarly.
- Shogi (Japanese chess) at 14, Go at least above 20, maybe 25?
- Chess computer programs (called *engines*) on desktop PC hardware reach 3200 on all rating lists, 3380 on CCRL.
- Computers at least even at Shogi, knocking on door at Go?

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.
- From 600 to 2800 gives chess a human depth of 11. 8×8 checkers estimated at 10, backgammon and bridge similarly.
- Shogi (Japanese chess) at 14, Go at least above 20, maybe 25?
- Chess computer programs (called *engines*) on desktop PC hardware reach 3200 on all rating lists, 3380 on CCRL.
- Computers at least even at Shogi, knocking on door at Go? "Moore's Law" of Games?

- Adult beginner typically 600, tournament/club "novice" 1200; scholastics go down below 100.
- László Mérő formalized the 75%-gap as a "Class Unit." Number of class units from beginner to champion = game's Human Depth.
- From 600 to 2800 gives chess a human depth of 11. 8×8 checkers estimated at 10, backgammon and bridge similarly.
- Shogi (Japanese chess) at 14, Go at least above 20, maybe 25?
- Chess computer programs (called *engines*) on desktop PC hardware reach 3200 on all rating lists, 3380 on CCRL.
- Computers at least even at Shogi, knocking on door at Go? "Moore's Law" of Games? Beyond chess ceiling of 3500??

• Primarily Skill Assessment; IPR for one event or series only.

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)
- From just 200–300 moves in a tournament, error bars are high, $2\sigma = \pm 200-300$ typical.

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)
- From just 200–300 moves in a tournament, error bars are high, $2\sigma = \pm 200$ –300 typical.
- Deep Blue played 2850-2900 in each of the matches against Garry Kasparov, while GK played

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)
- From just 200-300 moves in a tournament, error bars are high, $2\sigma = \pm 200-300$ typical.
- Deep Blue played 2850–2900 in each of the matches against Garry Kasparov, while GK played under 2600.

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)
- From just 200-300 moves in a tournament, error bars are high, $2\sigma = \pm 200-300$ typical.
- Deep Blue played 2850–2900 in each of the matches against Garry Kasparov, while GK played under 2600. But $\pm 225-300$.

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)
- From just 200-300 moves in a tournament, error bars are high, $2\sigma = \pm 200-300$ typical.
- Deep Blue played 2850–2900 in each of the matches against Garry Kasparov, while GK played under 2600. But $\pm 225-300$.
- Can pinpoint current quality of rapidly improving player.

- Primarily Skill Assessment; IPR for one event or series only.
- Based only on quality of your own move decisions. Results, opponents not involved.
- Your 50-100 games will have 1,200-2,400 relevant moves. (I standardly exclude turns 1-8 and positions where one side has an overwhelming advantage.)
- From just 200-300 moves in a tournament, error bars are high, $2\sigma = \pm 200-300$ typical.
- Deep Blue played 2850–2900 in each of the matches against Garry Kasparov, while GK played under 2600. But ±225–300.
- Can pinpoint current quality of rapidly improving player.
- "Match Elo" versus "Hidden Rating" at League of Legends.

• The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.

• The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.

• Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.

(日) (日) (日) (日) (日) (日) (日) (日)

• Loud "whispers" in various circles...

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.
- Loud "whispers" in various circles...
- But my full cheating test showed only a "1.3-sigma" deviation,

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.
- Loud "whispers" in various circles...
- But my full cheating test showed only a "1.3-sigma" deviation, and his IPR was "only" 2455 also within the "2-sigma" range.

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.
- Loud "whispers" in various circles...
- But my full cheating test showed only a "1.3-sigma" deviation, and his IPR was "only" 2455 also within the "2-sigma" range.
- Was dead lost against Epishin, lucked out also in previous round,

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.
- Loud "whispers" in various circles...
- But my full cheating test showed only a "1.3-sigma" deviation, and his IPR was "only" 2455 also within the "2-sigma" range.
- Was dead lost against Epishin, lucked out also in previous round,
- World #2 Fabiano Caruana had sensational 7-win streak against the top last Sept.

- The "San Sebastian Open"—a 9-round, 8-day prize-giving Swiss—had players up to 2600, 24 above 2200, 170 players total.
- Surprise winner: 2115-rated Badr Al-Hajiri of Kuwait.
- Won last 3 games over a 2356, 2412, and GM Vl. Epishin, 2563.
- Loud "whispers" in various circles...
- But my full cheating test showed only a "1.3-sigma" deviation, and his IPR was "only" 2455 also within the "2-sigma" range.
- Was dead lost against Epishin, lucked out also in previous round,
- World #2 Fabiano Caruana had sensational 7-win streak against the top last Sept.—but his IPR was "only" 2900 while his opponents played under 2600.

• Not a crystal ball to say what move a player will make...

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.

(日) (日) (日) (日) (日) (日) (日) (日)

• How a bookie sets odds—for the *initial betting line*.

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).

(日) (日) (日) (日) (日) (日) (日) (日)

• Quantify aggregate statistics:

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).
- Quantify aggregate statistics:
 - How often did the favored horses win in a racing week?
- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).
- Quantify *aggregate statistics*:
 - How often did the favored horses win in a racing week?

(日) (日) (日) (日) (日) (日) (日) (日)

• Do basketball teams average "covering their spread"?

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).
- Quantify *aggregate statistics*:
 - How often did the favored horses win in a racing week?
 - Do basketball teams average "covering their spread"?
 - How often did Player X make the move favored by an engine?

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).
- Quantify aggregate statistics:
 - How often did the favored horses win in a racing week?
 - Do basketball teams average "covering their spread"?
 - How often did Player X make the move favored by an engine?
 - How does his/her "Average Error" compare?

- Not a crystal ball to say what move a player will make...
- Though a GM sports-analyst friend tells me there is real-time betting on chess moves in Germany.
- How a bookie sets odds—for the *initial betting line*.
- Accuracy is how well odds "even out" over hundreds of betting events (for us, moves).
- Quantify aggregate statistics:
 - How often did the favored horses win in a racing week?
 - Do basketball teams average "covering their spread"?
 - How often did Player X make the move favored by an engine?

- How does his/her "Average Error" compare?
- Also project standard deviation and confidence intervals.

・ロト ・ 日 ・ モー・ モー・ うへぐ

Obmain: A set T of decision-making situations t. Chess game turns

ション ふゆ マ キャット マックシン

- Domain: A set T of decision-making situations t. Chess game turns
- Inputs: Values v_i for every option at turn t.
 Computer values of moves m_i

ション ふゆ マ キャット マックシン

- Domain: A set T of decision-making situations t. Chess game turns
- Inputs: Values v_i for every option at turn t.
 Computer values of moves m_i
- Parameters: s, c,... denoting skills and levels. Trained correspondence to chess Elo rating E

(日) (日) (日) (日) (日) (日) (日) (日)

- Domain: A set T of decision-making situations t. Chess game turns
- Inputs: Values v_i for every option at turn t.
 Computer values of moves m_i
- Parameters: s, c,... denoting skills and levels. Trained correspondence to chess Elo rating E
- Defines fallible agent P(s, c, ...).

- Domain: A set T of decision-making situations t. Chess game turns
- Inputs: Values v_i for every option at turn t. Computer values of moves m_i
- Parameters: s, c, ... denoting skills and levels. Trained correspondence to chess Elo rating E
- Defines fallible agent P(s, c, ...).
- So Main Output: Probabilities $p_{t,i}$ for P(s, c, ...) to select option i at time t.

- Domain: A set T of decision-making situations t. Chess game turns
- Inputs: Values v_i for every option at turn t.
 Computer values of moves m_i
- Parameters: s, c,... denoting skills and levels. Trained correspondence to chess Elo rating E
- Defines fallible agent P(s, c, ...).
- So Main Output: Probabilities $p_{t,i}$ for P(s, c, ...) to select option i at time t.
- Derived Outputs:
 - Aggregate statistics: move-match MM, equal-top value EV, average scaled difference ASD, ...
 - Projected confidence intervals: Bernoulli Trials + |T|-adjustment.

シック・ 川 ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・

• IPRs similarly reflect errors from the regression.

• Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .

・ロト ・ 日 ・ モー・ モー・ うへぐ

• Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .

Given s, c,..., the model computes x_i = g_{s,c}(v₁, v_i) = the perceived inferiority of m_i by P(s, c,...).

- Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .
- Given s, c,..., the model computes x_i = g_{s,c}(v₁, v_i) = the perceived inferiority of m_i by P(s, c,...).
- Besides g, the model picks a function $h(p_i)$ on probabilities.
- Could be h(p) = p (bad), log (good enough?), $H(p_i)$, logit...

- Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .
- Given s, c,..., the model computes x_i = g_{s,c}(v₁, v_i) = the perceived inferiority of m_i by P(s, c,...).
- Besides g, the model picks a function $h(p_i)$ on probabilities.
- Could be h(p) = p (bad), log (good enough?), $H(p_i)$, logit...

(日) (日) (日) (日) (日) (日) (日) (日)

• The Main Equation:

$$\frac{h(p_i)}{h(p_1)} = 1 - x_i$$

- Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .
- Given s, c,..., the model computes x_i = g_{s,c}(v₁, v_i) = the perceived inferiority of m_i by P(s, c,...).
- Besides g, the model picks a function $h(p_i)$ on probabilities.
- Could be h(p) = p (bad), log (good enough?), $H(p_i)$, logit...
- The Main Equation:

$$rac{h(p_i)}{h(p_1)}=1-x_i=\exp(-\left(rac{\delta(v_1,v_i)}{s}
ight)^c),$$

• Here $\delta(v_1, v_i)$ scales $v_1 - v_i$ in regard to $|v_1|$.

- Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .
- Given s, c,..., the model computes x_i = g_{s,c}(v₁, v_i) = the perceived inferiority of m_i by P(s, c,...).
- Besides g, the model picks a function $h(p_i)$ on probabilities.
- Could be h(p) = p (bad), log (good enough?), $H(p_i)$, logit...
- The Main Equation:

$$rac{h(p_i)}{h(p_1)}=1-x_i=\exp(-\left(rac{\delta(v_1,v_i)}{s}
ight)^c),$$

- Here $\delta(v_1, v_i)$ scales $v_1 v_i$ in regard to $|v_1|$.
- Ratio not difference on LHS so x_i on RHS has 0-to-1 scale.

- Let v_1, v_i be values of the best move m_1 and *i*th-best move m_i .
- Given s, c,..., the model computes x_i = g_{s,c}(v₁, v_i) = the perceived inferiority of m_i by P(s, c,...).
- Besides g, the model picks a function $h(p_i)$ on probabilities.
- Could be h(p) = p (bad), log (good enough?), $H(p_i)$, logit...
- The Main Equation:

$$rac{h(p_i)}{h(p_1)} = 1 - x_i = \exp(-\left(rac{\delta(v_1,v_i)}{s}
ight)^c),$$

- Here $\delta(v_1, v_i)$ scales $v_1 v_i$ in regard to $|v_1|$.
- Ratio not difference on LHS so x_i on RHS has 0-to-1 scale.
- Given $(x_1, \ldots, x_i, \ldots, x_\ell)$, fit subject to $\sum_i p_i = 1$ to find p_1 . Other p_i follow by $p_i = h^{-1}(h(p_1)(1-x_i))$.

• Over 3 million moves of 50-PV data: > 250 GB.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

- Over 3 million moves of 50-PV data: > 250 GB.
- Over 40 million moves of Single-PV data: > 50 GB

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Over 3 million moves of 50-PV data: > 250 GB.
- Over 40 million moves of Single-PV data: > 50 GB

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

• = 150 million pages of text data at 2k/page.

- Over 3 million moves of 50-PV data: > 250 GB.
- Over 40 million moves of Single-PV data: > 50 GB
- = 150 million pages of text data at 2k/page.
- All this was taken on two quad-core home-style PC's plus a laptop.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Over 3 million moves of 50-PV data: > 250 GB.
- Over 40 million moves of Single-PV data: > 50 GB
- = 150 million pages of text data at 2k/page.
- All this was taken on two quad-core home-style PC's plus a laptop. Is this "Big Data"?

- Over 3 million moves of 50-PV data: > 250 GB.
- Over 40 million moves of Single-PV data: > 50 GB
- = 150 million pages of text data at 2k/page.
- All this was taken on two quad-core home-style PC's plus a laptop. Is this "Big Data"?



- Over 3 million moves of 50-PV data: > 250 GB.
- Over 40 million moves of Single-PV data: > 50 GB
- = 150 million pages of text data at 2k/page.
- All this was taken on two quad-core home-style PC's plus a laptop. Is this "Big Data"? New sets being taken with UB CCR cluster.



• For each Elo level E training set, find (s, c, ...) giving best fit.

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

- For each Elo level E training set, find (s, c, ...) giving best fit.
- Or do Bayesian update to infer parameter(s) that best explain data [Haworth, later work joint with me and G. Di Fatta].

- For each Elo level E training set, find (s, c, ...) giving best fit.
- Or do Bayesian update to infer parameter(s) that best explain data [Haworth, later work joint with me and G. Di Fatta].

- In frequentist view, can use many different fitting methods...
 - Can compare methods...
 - Whole separate topic...

- For each Elo level E training set, find (s, c, ...) giving best fit.
- Or do Bayesian update to infer parameter(s) that best explain data [Haworth, later work joint with me and G. Di Fatta].
- In frequentist view, can use many different fitting methods...
 - Can compare methods...
 - Whole separate topic...
 - Max-Likelihood does relatively poorly.
- Often s and c trade off markedly, but E' ~ e(s, c) condenses into one Elo.
- Strong linear fit—suggests Elo mainly influenced by error.

 Bruce Pandolfini — played by Ben Kingsley in "Searching for Bobby Fischer."

・ロト ・ 日 ・ モー・ モー・ うへぐ

• Now does "Solitaire Chess" for Chess Life magazine:

- Bruce Pandolfini played by Ben Kingsley in "Searching for Bobby Fischer."
- Now does "Solitaire Chess" for Chess Life magazine:
 - Reader covers gamescore, tries to guess each move by one side.
 - E.g. score 6 pts. if you found 15.Re1, 4 pts. for 15.h3, 1 pt. for premature 15.Ng5.
 - Add points at end: say 150=GM, 140=IM, 120=Master, 80 = 1800 player, etc.

ション ふゆ マ キャット キャット しょう

- Bruce Pandolfini played by Ben Kingsley in "Searching for Bobby Fischer."
- Now does "Solitaire Chess" for Chess Life magazine:
 - Reader covers gamescore, tries to guess each move by one side.
 - E.g. score 6 pts. if you found 15.Re1, 4 pts. for 15.h3, 1 pt. for premature 15.Ng5.
 - Add points at end: say 150=GM, 140=IM, 120=Master, 80 = 1800 player, etc.

ション ふゆ マ キャット マックタン

Is it scientific?

- Bruce Pandolfini played by Ben Kingsley in "Searching for Bobby Fischer."
- Now does "Solitaire Chess" for Chess Life magazine:
 - Reader covers gamescore, tries to guess each move by one side.
 - E.g. score 6 pts. if you found 15.Re1, 4 pts. for 15.h3, 1 pt. for premature 15.Ng5.
 - Add points at end: say 150=GM, 140=IM, 120=Master, 80 = 1800 player, etc.

うして ふゆう ふほう ふほう ふしつ

- Is it scientific?
- With my formulas, yes—using *your* games in *real* tournaments.

- Bruce Pandolfini played by Ben Kingsley in "Searching for Bobby Fischer."
- Now does "Solitaire Chess" for Chess Life magazine:
 - Reader covers gamescore, tries to guess each move by one side.
 - E.g. score 6 pts. if you found 15.Re1, 4 pts. for 15.h3, 1 pt. for premature 15.Ng5.
 - Add points at end: say 150=GM, 140=IM, 120=Master, 80 = 1800 player, etc.
- Is it scientific?
- With my formulas, yes—using your games in real tournaments.
- Goal is **natural** scoring and distribution evaluation for multiple-choice tests, especially with partial-credit answers.

- Bruce Pandolfini played by Ben Kingsley in "Searching for Bobby Fischer."
- Now does "Solitaire Chess" for Chess Life magazine:
 - Reader covers gamescore, tries to guess each move by one side.
 - E.g. score 6 pts. if you found 15.Re1, 4 pts. for 15.h3, 1 pt. for premature 15.Ng5.
 - Add points at end: say 150=GM, 140=IM, 120=Master, 80 = 1800 player, etc.
- Is it scientific?
- With my formulas, yes—using your games in real tournaments.
- Goal is **natural** scoring and distribution evaluation for multiple-choice tests, especially with partial-credit answers.
- Connect to parameters in Item-Response Theory (IRT) test-taking models.

- Bruce Pandolfini played by Ben Kingsley in "Searching for Bobby Fischer."
- Now does "Solitaire Chess" for Chess Life magazine:
 - Reader covers gamescore, tries to guess each move by one side.
 - E.g. score 6 pts. if you found 15.Re1, 4 pts. for 15.h3, 1 pt. for premature 15.Ng5.
 - Add points at end: say 150=GM, 140=IM, 120=Master, 80 = 1800 player, etc.
- Is it scientific?
- With my formulas, yes—using your games in real tournaments.
- Goal is **natural** scoring and distribution evaluation for multiple-choice tests, especially with partial-credit answers.
- Connect to parameters in Item-Response Theory (IRT) test-taking models. IRT does both skill and prediction.

Separating Skill Assessment and Prediction [BHR]

• Thus far using same formulas for both.


・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)
 - Planning tendency may show in repeated moves of same piece.

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)
 - Planning tendency may show in repeated moves of same piece.

- Prediction Idea 2: Use Player-Specific Information ("profiling").
 - Regress previous games by player.

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)
 - Planning tendency may show in repeated moves of same piece.

- Prediction Idea 2: Use Player-Specific Information ("profiling").
 - Regress previous games by player.
 - Style is more "positional"?

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)
 - Planning tendency may show in repeated moves of same piece.

- Prediction Idea 2: Use Player-Specific Information ("profiling").
 - Regress previous games by player.
 - Style is more "positional"? "tactical"?

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)
 - Planning tendency may show in repeated moves of same piece.

- Prediction Idea 2: Use Player-Specific Information ("profiling").
 - Regress previous games by player.
 - Style is more "positional"? "tactical"?
- Drawbacks: loss of neutrality and portability.

- Thus far using same formulas for both.
- Linchpin: Use best-available computer move values for assessment.
- Prediction Idea 1: Use chess-specific features.
 - Good retreating moves are harder to find(?)
 - Planning tendency may show in repeated moves of same piece.
- Prediction Idea 2: Use Player-Specific Information ("profiling").
 - Regress previous games by player.
 - Style is more "positional"? "tactical"?
- Drawbacks: loss of neutrality and portability.
- Can we find more properties in the raw numerical data?

Example of "Swing" over Increasing Depths



Move	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Nd2	103	093	087	093	027	028	000	000	056	-007	039	028	037	020	014	017	000	006	000
Bxd7	048	034	-033	-033	-013	-042	-039	-050	-025	-010	001	000	-009	-027	-018	000	000	000	000
Qg8	114	114	-037	-037	-014	-014	-022	-068	-008	-056	-042	-004	-032	000	-014	-025	-045	-045	-050
Nxd4	-056	-056	-113	-071	-071	-145	-020	-006	077	052	066	040	050	051	-181	-181	-181	-213	-213

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□▶

• In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.

・ロト ・ 日 ・ モー・ モー・ うへぐ

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.

ション ふゆ マ キャット マックシン

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.
- The values are identical to the engine: it would not matter to the quality of the output which one the engine listed first. The values give no human reason to prefer one over the other.

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.
- The values are identical to the engine: it would not matter to the quality of the output which one the engine listed first. The values give no human reason to prefer one over the other.

うして ふゆう ふほう ふほう ふしつ

• So this is a kind of ESP test.

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.
- The values are identical to the engine: it would not matter to the quality of the output which one the engine listed first. The values give no human reason to prefer one over the other.
- So this is a kind of ESP test. How well do humans perform on it?

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.
- The values are identical to the engine: it would not matter to the quality of the output which one the engine listed first. The values give no human reason to prefer one over the other.
- So this is a kind of ESP test. How well do humans perform on it?
- Dick's Dean at Princeton: PEAR—Princeton Engineering Anomalies Research.

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.
- The values are identical to the engine: it would not matter to the quality of the output which one the engine listed first. The values give no human reason to prefer one over the other.
- So this is a kind of ESP test. How well do humans perform on it?
- Dick's Dean at Princeton: PEAR—Princeton Engineering Anomalies Research.
- PEAR did 10,000s-100,000s of trials, trying to judge significance of deviations like 50.1% or even 50.01%.

- In 8%-10% of positions, engine gives the top two moves the same value. Values are discrete up to 1 centipawn.
- More often *some* pair of moves in the top 10 (say) will end up tied.
- Conditioned on one of the two moves having been played, let us invite humans to guess which move is listed first by the program.
- The values are identical to the engine: it would not matter to the quality of the output which one the engine listed first. The values give no human reason to prefer one over the other.
- So this is a kind of ESP test. How well do humans perform on it?
- Dick's Dean at Princeton: PEAR—Princeton Engineering Anomalies Research.
- PEAR did 10,000s-100,000s of trials, trying to judge significance of deviations like 50.1% or even 50.01%.

⇒ > = √Q()

• How well do real humans perform on *my* ESP test??

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

・ロト ・ 日 ・ モ ト ・ モ ・ うへぐ

(0.01, the higher move is played 53–55% of the time.

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- 0.01, the higher move is played 53–55% of the time.
- \bigcirc 0.02, the higher move is played 58–59% of the time.

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- 0.01, the higher move is played 53–55% of the time.
- \bigcirc 0.02, the higher move is played 58–59% of the time.
- \bigcirc 0.03, the higher move is played 60-61% of the time.

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- **0** 0.01, the higher move is played 53-55% of the time.
- 20.02, the higher move is played 58-59% of the time.
- 0.03, the higher move is played 60-61% of the time.
- 0.00, the higher move is played

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- **0** 0.01, the higher move is played 53-55% of the time.
- 20.02, the higher move is played 58-59% of the time.
- 0.03, the higher move is played 60-61% of the time.
- 0.00, the higher move is played 57-59% of the time.

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- **0** 0.01, the higher move is played 53-55% of the time.
- 20.02, the higher move is played 58-59% of the time.
- 0.03, the higher move is played 60-61% of the time.
- 0.00, the higher move is played 57-59% of the time.
 - Last is not a typo—see post "When is a Law Natural?"

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- **0** 0.01, the higher move is played 53-55% of the time.
- 20.02, the higher move is played 58-59% of the time.
- 0.03, the higher move is played 60-61% of the time.
- 0.00, the higher move is played 57-59% of the time.
- Last is not a typo—see post "When is a Law Natural?"
- Similar 58%-42% split seen for any pair of tied moves. What can explain it?

Conditioned on one of the top two moves being played, if their values (Rybka 3, depth 13) differ by...:

- **0** 0.01, the higher move is played 53-55% of the time.
- 20.02, the higher move is played 58-59% of the time.
- 0.03, the higher move is played 60-61% of the time.
- 0.00, the higher move is played 57-59% of the time.
 - Last is not a typo—see post "When is a Law Natural?"
 - Similar 58%-42% split seen for any pair of tied moves. What can explain it?
 - Relation to slime molds and other "semi-Brownian" systems?

• Non-Parapsychological Explanation:

• Non-Parapsychological Explanation: Stable Library Sorting.

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.
- Lead moves tend to have been higher at lower depths. Lower move "swings up."

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.
- Lead moves tend to have been higher at lower depths. Lower move "swings up."

• Formulate numerical measure of swing "up" and "down" (a trap).

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.
- Lead moves tend to have been higher at lower depths. Lower move "swings up."
- Formulate numerical measure of swing "up" and "down" (a trap).
- When best move swings up 4.0-5.0 versus 0.0-1.0, players rated 2700+ find it only 30% versus 70%.

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.
- Lead moves tend to have been higher at lower depths. Lower move "swings up."
- Formulate numerical measure of swing "up" and "down" (a trap).
- When best move swings up 4.0-5.0 versus 0.0-1.0, players rated 2700+ find it only 30% versus 70%.
- Goal is to develop a Challenge Quotient based on how much trappy play a player sets for the opponent

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.
- Lead moves tend to have been higher at lower depths. Lower move "swings up."
- Formulate numerical measure of swing "up" and "down" (a trap).
- When best move swings up 4.0-5.0 versus 0.0-1.0, players rated 2700+ find it only 30% versus 70%.
- Goal is to develop a Challenge Quotient based on how much trappy play a player sets for the opponent—and emself.
Measuring "Swing" and Complexity and Difficulty

- Non-Parapsychological Explanation: Stable Library Sorting.
- Chess engines sort moves from last depth to schedule next round of search.
- By stability, lower move can become 1st only with *strictly higher* value.
- Lead moves tend to have been higher at lower depths. Lower move "swings up."
- Formulate numerical measure of swing "up" and "down" (a trap).
- When best move swings up 4.0-5.0 versus 0.0-1.0, players rated 2700+ find it only 30% versus 70%.

・ロト ・ 日 ・ ・ ヨ ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

- Goal is to develop a Challenge Quotient based on how much trappy play a player sets for the opponent—and emself.
- Separates *performance* and *prediction* in the model.

Human Versus Computer Phenomena



- [show data]
- The metric correction

$$\int_{e-\delta}^e d\mu \quad ext{with} \quad d\mu = rac{c}{c+x}\,dx$$

ション ふゆ マ キャット キャット しょう

- [show data]
- The metric correction

$$\int_{e-\delta}^e d\mu \quad ext{with} \quad d\mu = rac{c}{c+x}\,dx$$

うして ふゆう ふほう ふほう ふしつ

balances evals well for Rybka, with c very near 1.0.

• A mix of three factors?

- [show data]
- The metric correction

$$\int_{e-\delta}^e d\mu \quad {
m with} \quad d\mu = rac{c}{c+x}\,dx$$

うして ふゆう ふほう ふほう ふしつ

- A mix of three factors?
- (A) Human perception of value as proportional to stakes, *per* Ariely-Kahneman-Tversky.

- [show data]
- The metric correction

$$\int_{e-\delta}^e d\mu \quad ext{with} \quad d\mu = rac{c}{c+x}\,dx$$

- A mix of three factors?
- (A) Human perception of value as proportional to stakes, *per* Ariely-Kahneman-Tversky.
- (B) Rationally playing less *catenaccio* when marginal impact of evaluation on win probability is minimal. (Leo Stedile, working under Mark Braverman)

- [show data]
- The metric correction

$$\int_{e-\delta}^e d\mu \quad ext{with} \quad d\mu = rac{c}{c+x}\,dx$$

- A mix of three factors?
- (A) Human perception of value as proportional to stakes, *per* Ariely-Kahneman-Tversky.
- (B) Rationally playing less catenaccio when marginal impact of evaluation on win probability is minimal. (Leo Stedile, working under Mark Braverman)
- (C) Greater volatility intrinsic to chess as game progresses.

A. Perception Proportional to Benefit

How strongly do you perceive a difference of 10 dollars, if:

- You are buying lunch and a drink in a pub.
- You are buying dinner in a restaurant.
- You are buying an I-pad.
- You are buying a car.

For the car, maybe you don't care. In other cases, would you be equally thrifty?

うして ふゆう ふほう ふほう ふしつ

If you spend the way you play chess, you care maybe $4 \times$ as much in the pub!

• Expectation curves according to position evaluation v are sigmoidal, indeed close to a hyperbolic tangent

$$E=rac{e^{av}-e^{-av}}{e^{av}+e^{-av}}.$$

・ロト ・ 日 ・ モ ・ ト ・ モ ・ うへぐ

• Expectation curves according to position evaluation v are sigmoidal, indeed close to a hyperbolic tangent

$$E=rac{e^{av}-e^{-av}}{e^{av}+e^{-av}}.$$

ション ふゆ マ キャット キャット しょう

• Here a gives pretty steep slope near 0, $a \approx 4.5$ for Rybka and Houdini.

• Expectation curves according to position evaluation v are sigmoidal, indeed close to a hyperbolic tangent

$$E=rac{e^{av}-e^{-av}}{e^{av}+e^{-av}}.$$

- Here a gives pretty steep slope near 0, $a \approx 4.5$ for Rybka and Houdini.
- How to test apart from cause A?

• Expectation curves according to position evaluation v are sigmoidal, indeed close to a hyperbolic tangent

$$E=rac{e^{av}-e^{-av}}{e^{av}+e^{-av}}.$$

- Here a gives pretty steep slope near 0, $a \approx 4.5$ for Rybka and Houdini.
- How to test apart from cause A?
- Expect eval-error curve to shift in games between unequally-rated players.

• Expectation curves according to position evaluation v are sigmoidal, indeed close to a hyperbolic tangent

$$E=rac{e^{av}-e^{-av}}{e^{av}+e^{-av}}.$$

- Here a gives pretty steep slope near 0, $a \approx 4.5$ for Rybka and Houdini.
- How to test apart from cause A?
- Expect eval-error curve to shift in games between unequally-rated players.
- Results so far show no shift—

Human Versus Computer Phenomena



イロト イ理ト イヨト イヨト ヨー シタの

Eval-Error Curve With Unequal Players



• Carlsen:

- 2985 at London 2011 (Kramnik 2857, Aronian 2838).
- Kasparov:
 - Was playing 2860 to Karpov's 2760 when 1984-85 match aborted.

・ロト ・ 日 ・ モー・ モー・ うへぐ

• Carlsen:

- 2985 at London 2011 (Kramnik 2857, Aronian 2838).
- Kasparov:
 - Was playing 2860 to Karpov's 2760 when 1984-85 match aborted.

• Both over 2800 in 1986, Kasparov 2905.

• Carlsen:

- 2985 at London 2011 (Kramnik 2857, Aronian 2838).
- Kasparov:
 - Was playing 2860 to Karpov's 2760 when 1984-85 match aborted.

ション ふゆ マ キャット キャット しょう

- Both over 2800 in 1986, Kasparov 2905.
- Both under 2675 in New York-Lyon match 1990.

- Carlsen:
 - 2985 at London 2011 (Kramnik 2857, Aronian 2838).
- Kasparov:
 - Was playing 2860 to Karpov's 2760 when 1984-85 match aborted.
 - Both over 2800 in 1986, Kasparov 2905.
 - Both under 2675 in New York-Lyon match 1990.
- Bobby Fischer:
 - 2920 over all 3 Candidates' Matches in 1971.
 - 2650 vs. Spassky in 1972 (Spassky 2645).
 - 2725 vs. Spassky in 1992 (Spassky 2660).
- Hou Yifan: 2970 vs. Humpy Koneru (2685) in Nov. 2011.
- Paul Morphy: 2345 in 59 most impt. games, 2125 vs. Anderssen.

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト ・ らくぐ

- Capablanca: 2935 at New York 1927.
- Alekhine: 2810 in 1927 WC match over Capa (2730).

Computer and Freestyle IPRs

Analyzed Ratings of Computer Engine Grand Tournament (on commodity PCs) and PAL/CSS Freestyle in 2007–08, plus the Thoresen Chess Engines Competition (16-core) Nov–Dec. 2013.

Event	Rating	2σ range	#gm	#moves
CEGT g1,50	3009	2962-3056	42	4,212
CEGT g25,26	2963	2921-3006	42	5,277
PAL/CSS 5ch	3102	3051–3153	45	3,352
PAL/CSS 6ch	3086	3038–3134	45	3,065
PAL/CSS 8ch	3128	3083–3174	39	3,057
TCEC 2013	3083	3062-3105	90	11,024

Computer and Freestyle IPRs—To Move 60

Computer games can go very long in dead drawn positions. TCEC uses a cutoff but CEGT did not. Human-led games tend to climax (well) before Move 60. This comparison halves the difference to CEGT, otherwise similar:

Sample set	Rating	2σ range	#gm	#moves
CEGT all	2985	2954-3016	84	9,489
PAL/CSS all	3106	3078–3133	129	9,474
TCEC 2013	3083	3062–3105	90	11,024
CEGT to60	3056	3023–3088	84	7,010
PAL/CSS to60	3112	3084–3141	129	8,744
TCEC to60	3096	3072-3120	90	8,184

Degrees of Forcing Play



Forcing Index (2500 perspective)

◆□▶ ◆□▶ ◆注▶ ◆注▶ 注 のへで

Add Human-Computer Tandems



Forcing Index (2500 perspective)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Add Human-Computer Tandems



Forcing Index (2500 perspective)

Evidently the humans called the shots.

Add Human-Computer Tandems



Forcing Index (2500 perspective)

Evidently the humans called the shots. But how did they play?

2007–08 Freestyle Performance



Forcing Index (2500 perspective)

Adding 210 Elo was significant. Forcing but good teamwork.

2014 Freestyle Tournament Performance



Tandems had marginally better W-L, but quality not clear...

Add Topalov Forcing Kramnik



Forcing Index (2500 perspective)

Last bar goes way off the chart

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Is There Room to Grow?

- In chess, alas some hints of "no."
- If (randomizing) 3200-level programs can score 10% against any strategy, then no strategy can ever exceed 3550.
- In 2010–2014 many more games between players rated under 1600 and between 2800+ became available.
- Analysis in my model shows a linear relationship between rating and my Average Scaled Difference ASD statistic clear down to 1200 level.
- The y-intercept of the line is consistently near 3370.
- But Komodo and Stockfish on 4-core PCs are rated over 3370 on CCRL. How can this be?
 - Well, CCRL uses a 40 moves in 40 minutes time control. Other lists use other times and show ratings still in the 3100s.
- Best explanation: IPR correlates 85–90% with ASD and 10–15% with move-matching—which has y-intercept near 4500.

Solution and Opportunities

- Hence my model projects a ceiling around **3500-3550**.
- Still not much room to grow... in chess that is.
- This may already explain the diminishing returns from adding humans... in chess.
- But the larger marriage of Shallow but Broad to Deep but Narrow that was theoretically driving the gains still has potential.
- Revisit trying to "humanize" chess programs?
- Complexity theory classifies chess as "Hard to Parallelize."
- Whether chess endgame tables are "Associatively Compressible" is an indicator worth pursuing.
- Model has many other applications: study human performance under distraction; design multiple choice tests to standards of difficulty; extend *intrinsic* Elo quality measures to other domains.

- Lots more potential for research and connections...
- Can use support—infrastructure, student helpers...
 - Run data with other engines Houdini, Stockfish, Komodo....
 - Run more tournaments.
 - Run to higher depths—how much does that matter?
- Spread word about general-scientific aspects, including public outreach over what isn't (and is) cheating.

- Lots more potential for research and connections...
- Can use support—infrastructure, student helpers...
 - Run data with other engines Houdini, Stockfish, Komodo....
 - Run more tournaments.
 - Run to higher depths—how much does that matter?
- Spread word about general-scientific aspects, including public outreach over what isn't (and is) cheating.

うして ふゆう ふほう ふほう ふしつ

• Detect and deter cheating too—generally.

- Lots more potential for research and connections...
- Can use support—infrastructure, student helpers...
 - Run data with other engines Houdini, Stockfish, Komodo....
 - Run more tournaments.
 - Run to higher depths—how much does that matter?
- Spread word about general-scientific aspects, including public outreach over what isn't (and is) cheating.

- Detect and deter cheating too—generally.
- Learn more about human decision making.

- Lots more potential for research and connections...
- Can use support—infrastructure, student helpers...
 - Run data with other engines Houdini, Stockfish, Komodo....
 - Run more tournaments.
 - Run to higher depths—how much does that matter?
- Spread word about general-scientific aspects, including public outreach over what isn't (and is) cheating.

- Detect and deter cheating too—generally.
- Learn more about human decision making.
- Thus the Turing Tour comes back to the human mind.

- Lots more potential for research and connections...
- Can use support—infrastructure, student helpers...
 - Run data with other engines Houdini, Stockfish, Komodo....
 - Run more tournaments.
 - Run to higher depths—how much does that matter?
- Spread word about general-scientific aspects, including public outreach over what isn't (and is) cheating.

- Detect and deter cheating too—generally.
- Learn more about human decision making.
- Thus the Turing Tour comes back to the human mind.
- Thank you very much for the invitation.