# Decision Making
## *in Robots and Autonomous Agents*

## Behavioural Issues:
## Heuristics & Biases, etc.

Subramanian Ramamoorthy

School of Informatics

16 March, 2018

# The Rational Animal

The Greeks (Aristotle) thought that humans are rational animals.

- Many ancient philosophers agreed
- Clearly, even Aristotle was aware that people's judgments, decisions & behavior are not *always* rational.
- A more plausible revised claim: all normal humans have the *competence* to be rational.

Argument thus far in DMR has been in this spirit:

- "correct" principles of reasoning & decision making are in our minds, even if we don't always use them
- "Exceptions" were things like Bernoulli's observation that we become progressively more indifferent to larger gains

# Heuristics and Biases

- Tversky and Kahneman dramatically extended Bernoulli's account, and provided extensive psychological evidence for their model.

- What they showed is that humans depart in multiple ways from a classical picture of rational behavior.

- Many have interpreted the Kahneman & Tversky program as showing that Aristotle was wrong.

  - most people *do not* have the correct principles for reasoning & decision making

  - we get by with much simpler principles which *only sometimes* get the right answer

Amos Tversky
1937 - 1996

Daniel Kahneman

# An Intuition from Gigerenzer

*Which US city has more inhabitants, San Diego or San Antonio?*

- 62% of Americans got this correct

- 100% of Germans got this correct

The Recognition Heuristic:

> If one of two objects is recognized and the other is not, then infer that the recognized object has the higher value.

Ecological Rationality:

> Heuristic works when ignorance is systematic (non – random), i.e., when lack of recognition correlates with the criterion.

# Experimental Studies of Human Reasoning

- **The Conjunction Fallacy**
- Base Rate Neglect
- Overconfidence
- Framing

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Please rank the following statements by their probability, using 1 for the most probable and 8 for the least probable.

**(a)** Linda is a teacher in elementary school.
**(b)** Linda works in a bookstore and takes Yoga classes.
**(c)** Linda is active in the feminist movement.
**(d)** Linda is a psychiatric social worker.
**(e)** Linda is a member of the League of Women Voters.
**(f)** Linda is a bank teller.
**(g)** Linda is an insurance sales person.
**(h)** Linda is a bank teller & is active in the feminist movement.

Naïve subjects: (h) > (f) 89%

Sophisticated subjects: 85%

John is 19, wears glasses, is a little shy unless you talk to him about Star Trek or Lord of the Rings, and stays up late most nights playing video games.

Which is more likely?

1. John is a CS major who loves contemporary sculpture and golf, or

2. John loves contemporary sculpture and golf?

On problems like this, people tend to pick (1) above (2), in contradiction to probability theory.

# Experimental Studies of Human Reasoning

- The Conjunction Fallacy
- **Base Rate Neglect**
- Overconfidence
- Framing

A panel of psychologists have interviewed and administered personality tests to 30 engineers and 70 lawyers [70 engineers and 30 lawyers], all successful in their respective fields.  On the basis of this information, thumbnail descriptions of the 30 engineers and 70 lawyers [70 engineers and 30 lawyers], have been written.  You will find on your forms five descriptions, chosen at random from the 100 available descriptions.  For each description, please indicate your probability that the person described is an engineer, on a scale from 0 to 100.

Jack is a 45-year-old man. He is married and has four children. He is generally conservative, careful and ambitious. He shows no interest in political and social issues and spends most of his free time on his many hobbies which include home carpentry, sailing, and mathematical puzzles.

Dick is a 30-year-old man. He is married with no children. A man of high ability and high motivation, he promises to be quite successful in his field. He is well liked by his colleagues.

# Base Rate and Representativeness

With no personality sketch, simply asked for the probability that an unknown individual was an engineer
- subjects correctly gave the responses 0.7 and 0.3

When presented with a totally uninformative description
- the subjects gave the probability to be 0.5

Kahneman & Tversky concluded that when no specific evidence is given, prior probabilities are used properly; when worthless evidence is given, prior probabilities are ignored.

If a test to detect a disease whose prevalence is 1/1000 has a false positive rate of 5%, what is the chance that a person found to have a positive result actually has the disease, assuming that you know nothing about the person's symptoms or signs?

_____%

Harvard Medical School: "2%" = 18%; "95%" = 45%

# Why is the correct answer 2%?

- Think of a population of 10,000 people.

- We would expect just 10 people in this population to have the disease (1/1000 x 10,000 = 10)

- If you test everybody in the population then false positive rate means that, in addition to the 10 people who do have the disease, another 500 (5%) will be wrongly diagnosed as having it.

- In other words only about 2% of the people diagnosed positive (10/510) actually have the disease.

- When people give a high answer like 95% they are ignoring the very low probability (i.e. rarity) of having the disease. In comparison, probability of a false positive test is relatively high.

# Experimental Studies of Human Reasoning

- The Conjunction Fallacy
- Base Rate Neglect
- **Overconfidence**
- Framing

In each of the following pairs, which city has more inhabitants?

- (a) Las Vegas              (b) Miami

  How confident are you that your answer is correct?

  50%  60%  70%  80%  90%  100%

- (a) Sydney                 (b) Melbourne

  How confident are you that your answer is correct?

  50%  60%  70%  80%  90%  100%

- (a) Hyderabad              (b) Islamabad

  How confident are you that your answer is correct?

  50%  60%  70%  80%  90%  100%

# Overconfidence effect

- Person's subjective *confidence* in his or her judgements is <u>reliably greater than</u> the objective *accuracy* of those judgements
  - especially when confidence is high

- Manifested in different ways:
  - *overestimation* of one's actual performance
  - *overplacement* of one's performance relative to others
  - *overprecision* in expressing unwarranted certainty in the accuracy of one's beliefs.

# Experimental Studies of Human Reasoning

- The Conjunction Fallacy
- Base Rate Neglect
- Overconfidence
- **Framing**

Imagine the U.S. is preparing for the outbreak of an unusual Asian disease which is expected to kill 600 people. Two alternative programs to combat the disease have been proposed. Assume that the exact scientific estimates of the consequences of the program are as follows:

- If program A is adopted, 200 people will be saved.

- If program B is adopted, there is a 1/3 probability that 600 people will be saved and a 2/3 probability that no people will be saved.

- If program A is adopted , 200 people will be saved.

- If program B is adopted, there is a 1/3 probability that 600 people will be saved and a 2/3 probability that no people will be saved.

- If program C is adopted, 400 people will die.

- If program D is adopted, there is a 1/3 probability that nobody will die and a 2/3 probability that 600 people will die.

– If program C is adopted, 400 people will die.

– If program D is adopted, there is a 1/3 probability that nobody will die and a 2/3 probability that 600 people will die.

- If program A is adopted , 200 people will be saved.

- If program B is adopted, there is a 1/3 probability that 600 people will be saved and a 2/3 probability that no people will be saved.

- If program C is adopted, 400 people will die.

- If program D is adopted, there is a 1/3 probability that nobody will die and a 2/3 probability that 600 people will die.

# Program A versus B

- Tversky and Kahneman (1981) found  that the majority choice (72%) for an analogue of problem 1 using an anonymous "Asian disease" (not SARS) was answer A.

- The prospect of saving 200 lives with certainty was more promising than the probability of a one-in-three chance of saving 600 lives.

# A and B have the Same Utility

- The standard way of calculating expected utility:
    - $\Sigma_i P(o_i) \times U(o_i)$
    - where each $o_i$ is a possible outcome, $P(o_i)$ is its probability and $U(o_i)$ is its utility or value.

- On this basis, option B (*a risky prospect*) would be of equal expected value to the first prospect A.

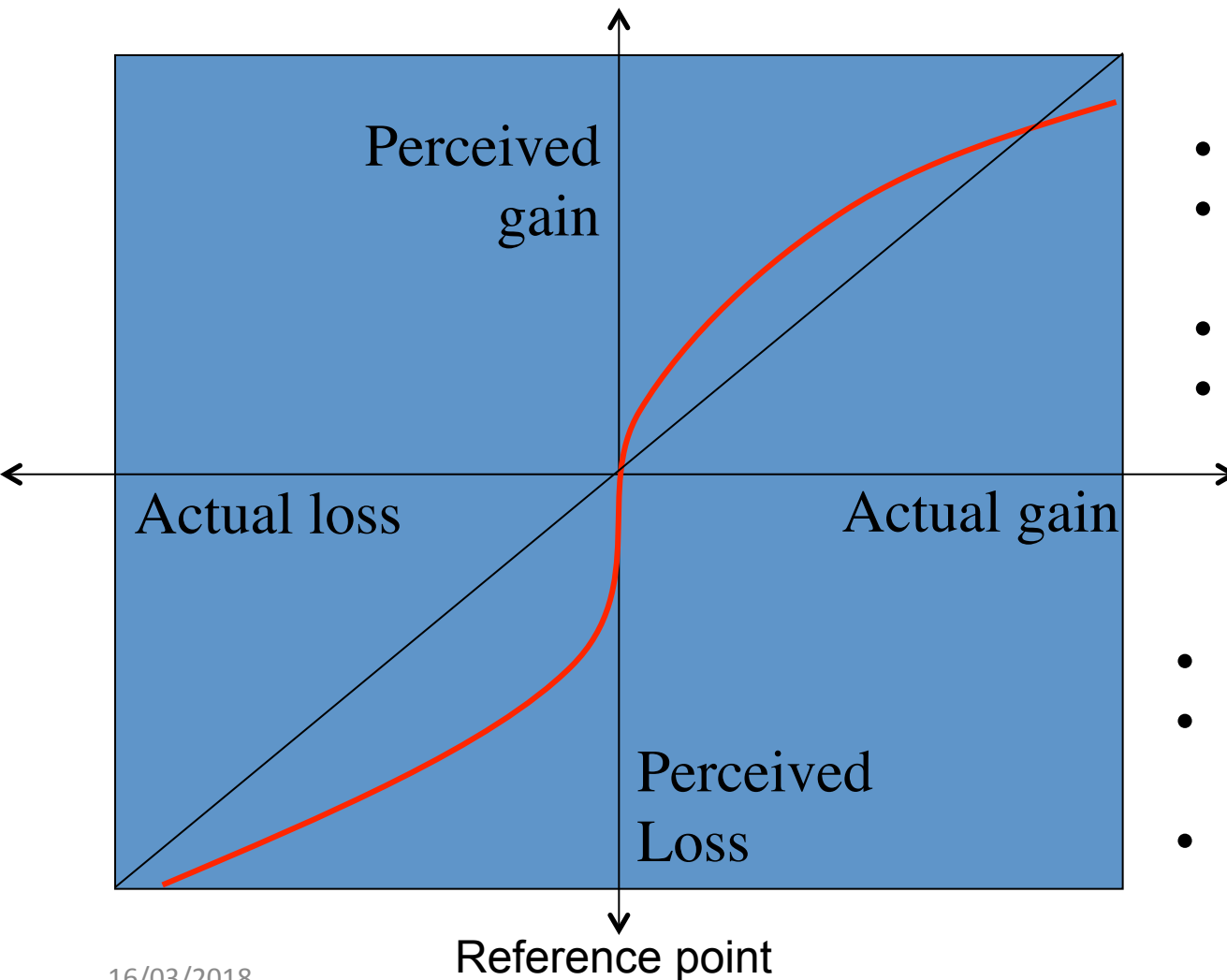- So in this case subjects are not simply using expected utility: they are **risk averse**.

# Program C versus D

- The majority of respondents in the second problem (78%) chose the riskier option.

- The certain death of 400 people is apparently less acceptable than the two-in-three chance that 600 people will die.

- Once again, there is no difference between the options (C and D) in standard expected utility.

- We say in this case that people are risk seeking.

# Framing

- Tversky and Kahneman's intention was that Problem 1 and Problem 2 are underlyingly identical, but presented in different ways.

- To the extent that they are right, we say that they are different framings of the same problem: the different responses for the two problems are what we call framing effects.

- Yet another way in which people's decisions differ from what would be predicted on the standard view of expected utility.

# Prospect Theory: S-curve of Value



- Small gain good,
- Big gain not much better,
- Small loss *terrible*,
- Big loss not much worse.

Thus we are generally:
- Loss averse,
- Risk averse for gains, but
- Risk seeking for losses.

# Estimation of Risk

Estimates of Probabilities of Death From Various Causes:

| Cause | Stanford graduate estimates | Statistical estimates |
|---|---|---|
| Heart Disease | 0.22 | 0.34 |
| Cancer | 0.18 | 0.23 |
| Other Natural Causes | 0.33 | 0.35 |
| All Natural Causes | 0.73 | 0.92 |
| | | |
| Accident | 0.32 | 0.05 |
| Homicide | 0.10 | 0.01 |
| Other Unnatural | 0.11 | 0.02 |
| All Unnatural | 0.53 | 0.08 |

(Based on a study by Amos Tversky, Thayer Watkin's report)

# Big and Small Probabilities

- Small probability events are overrated.

- More generally, behavior departs from classical rationality substantially as regards very small and very large probabilities.

- Even when told the frequency of events, subjects tend to view the difference between 0 and 0.001 (or that between 0.999 and 1) as more significant than that between .5 and .5001.

# The "Bleak Implications" Hypothesis

- These results have "bleak implications" for the rationality of ordinary people (Nisbett et al.)

- "It appears that people lack the correct programs for many important judgmental tasks…. <span style="color:green">We have not had the opportunity to evolve an intellect capable of dealing conceptually with uncertainty.</span>" (Slovic, Fischhoff and Lichtenstein , 1976)

# The "Bleak Implications" Hypothesis

"I am particularly fond of [the Linda] example, because I know that the [conjunction] is least probable, yet a little homunculus is my head continues to jump up and down, shouting at me – "but she can't just be a bank teller; read the description." … Why do we consistently make this simple logical error?    Tversky and Kahneman argue, correctly I think, that our minds are not built (for whatever reason) to work by the rules of probability."   (Stephen J. Gould , 1992, p. 469)

# Challenge from Evolutionary Psychology

- The "frequentist" hypothesis
  - Useful information about probabilities was available to our ancestors in the form of *frequencies*
    - e.g., 3 of the last 12 hunts near rivers were successful
  - Information about the probabilities of single events was not available
  - So perhaps we evolved a mental capacity that is good at dealing with probabilistic information, but only when that information is presented in a frequency format.

# Making the Conjunction Fallacy "Disappear"…

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Please rank the following statements by their probability, using 1 for the most probable and 8 for the least probable.

**(a)** Linda is a teacher in elementary school.
**(b)** Linda works in a bookstore and takes Yoga classes.
**(c)** Linda is active in the feminist movement.
**(d)** Linda is a psychiatric social worker.
**(e)** Linda is a member of the League of Women Voters.
**(f)** Linda is a bank teller.
**(g)** Linda is an insurance sales person.
**(h)** Linda is a bank teller and is active in the feminist movement.

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

There are 100 people who fit the description above. How many of them are:
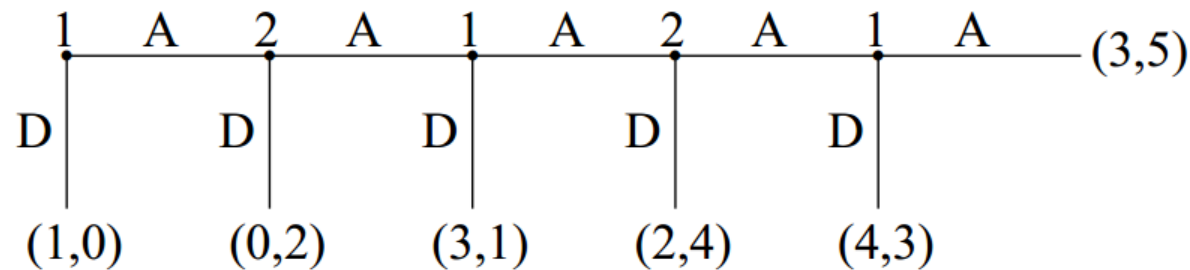
…

**(f)** bank tellers?

**(h)** bank tellers and active in the feminist movement?

...

$$\textbf{(h)} > \textbf{(f)} \quad = \quad 13\%$$

# Level-k Reasoning

The centipede game:



- In the only equilibrium involving rational players, player 1 chooses down in the first time move

- What exactly do ***real*** people do?

  – Human subjects don't choose the down option right away

  – How do we understand their behaviour?

# Hide and Seek Games

- Your opponent has hidden a prize in one of four boxes arranged in a row.

- The boxes are marked as shown below: A, B, A, A

- Your goal is, of course, to find the prize.

- His goal is that you will not find it.

- You are allowed to open only one box.

- Which box are you going to open?

# Folk Wisdom

- "…in Lake Wobegon, the correct answer is usually 'c'."
    - Garrison Keillor (1997) on multiple-choice tests


- Comment on poisoning of Ukrainian's presidential candidate:
  "Any government wanting to kill an opponent
  …would not try it at a meeting with government officials."
    - Viktor Yushchenko (2004)

# Unique Equilibrium of Hide and Seek Game

- If both players randomize uniformly:

| Hider/Seeker | A | B | A | A |
|:---:|:---:|:---:|:---:|:---:|
| A | 0,1 | 1,0 | 1,0 | 1,0 |
| B | 1,0 | 0,1 | 1,0 | 1,0 |
| A | 1,0 | 1,0 | 0,1 | 1,0 |
| A | 1,0 | 1,0 | 1,0 | 0,1 |

Expected payoffs: Hider 3/4, Seeker 1/4

# Behavioural Experiment in Hide and Seek

|  | A | B | A | A |
|---|---|---|---|---|
| Hiders (624) | 0.2163 | 0.2115 | **0.3654** | 0.2067 |
| Seekers (560) | 0.1821 | 0.2054 | **0.4589** | 0.1536 |

- Central A (or 3) is most prevalent for both Hiders and Seekers
- Central A is even more prevalent for Seekers
  - As a result, Seekers do better than in equilibrium
- Shouldn't Hiders realize that Seekers will be just as tempted to look there?
- "The finding that both choosers and guessers selected the least salient alternative suggests little or no strategic thinking."

# p-Beauty Contest Games (Keynes, 1936)

"…professional investment may be likened to those newspaper competitions in which the competitors have to pick out the six prettiest faces from a hundred photographs,

the prize being awarded to the competitor whose choice most nearly corresponds to the average preferences of the competitors as a whole….

It is not a case of choosing those [faces] that, to the best of one's judgment, are really the prettiest, nor even those that average opinion genuinely thinks the prettiest.

# p-Beauty Contest Games

- We have reached the **third** degree where we devote our intelligences to…

  anticipating what average opinion expects the average opinion to be.

- And there are some, I believe, who practice the **fourth**, **fifth** and **higher** degrees."

  – Keynes, General Theory, 1936, pp. 155-56

# Level-k Reasoning

One could play at varying levels of depth of reasoning:

- Level-0: Random play

- Level-1: BR to Random play

- Level-2: BR to Level-1

- Nash: Play Nash Equilibrium (level-∞?)

- Worldly: BR to distribution of Level-0, Level-1 and Nash types

Stahl and Wilson (GEB 1995)
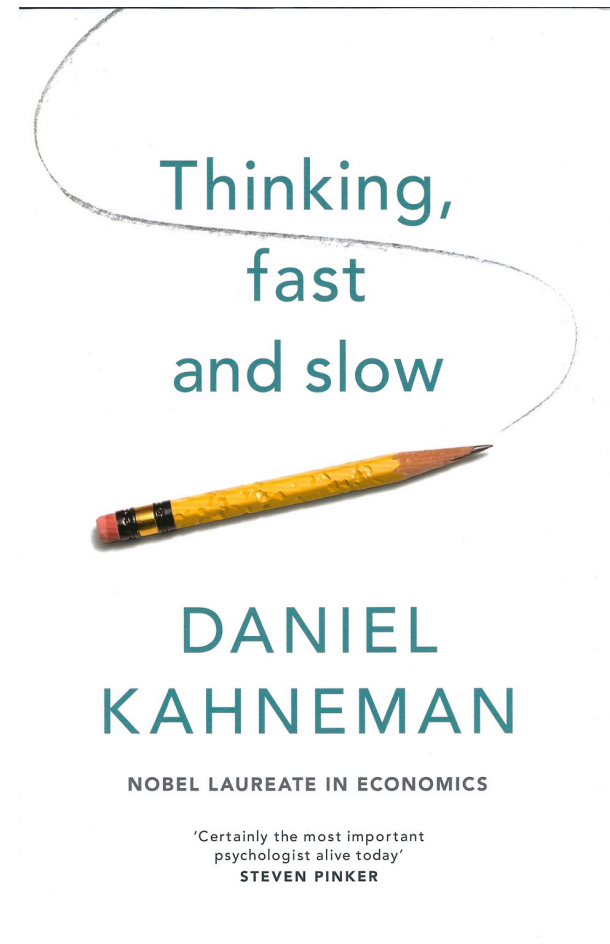
# Level-k for Hide and Seek

- Level-k: Each role is filled by Lk types: L0, L1, L2, L3, or L4 (probabilities to be estimated...)
  - Note: In Hide and Seek the types cycle after L4...

- High types anchor beliefs in a naïve L0 type and adjusts with iterated best responses:
  - L1 best responds to L0 (with uniform errors)
  - L2 best responds to L1 (with uniform errors)
  - Lk best responds to Lk-1 (with uniform errors)

# Level-k for Hide and Seek

- L0 Hiders and Seekers are symmetric
  - Favor salient locations equally

  1. Favor "B": choose with probability q > 1/4
  2. Favor "end A": choose with prob. p/2>1/4
     - Choice probabilities: (p/2, q, 1-p-q, p/2)

- *Note*: Specification of Anchoring Type L0 is the key to model's explanatory power
     - Can't use uniform L0 (coincide with equilibrium)...

# Modern Perspective:
# Diversity in Modes of Thinking

- **Automatic System (Type 1)**
  - Fast, unconscious
  - Parallel, associative
  - Low energy
  - "Doer"

- **Reflective System (Type 2)**
  - Slow, conscious
  - Serial
  - Expensive
  - "Planner"

Thinking,
fast
and slow

DANIEL
KAHNEMAN

NOBEL LAUREATE IN ECONOMICS

'Certainly the most important
psychologist alive today'
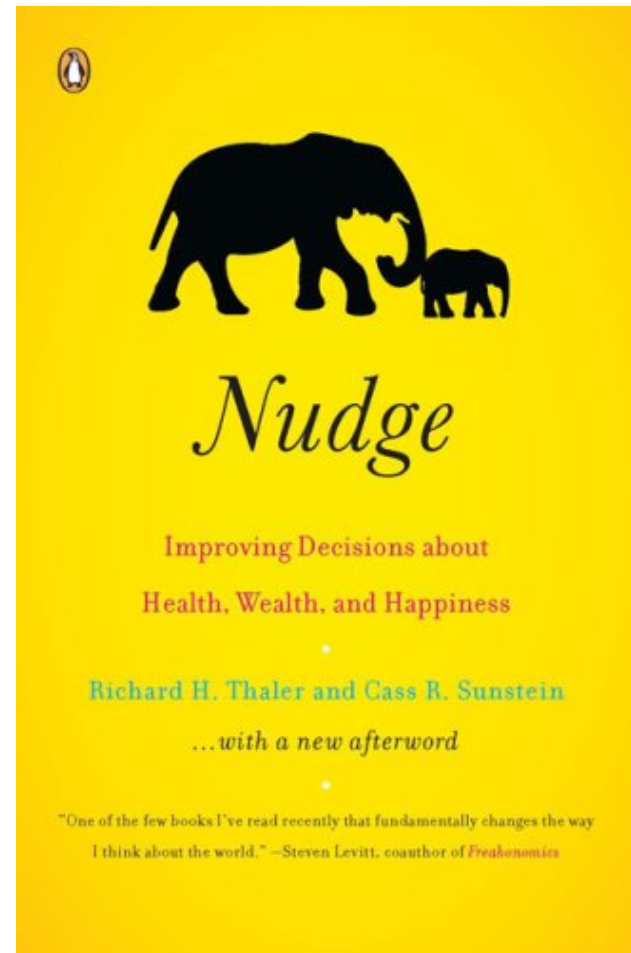**STEVEN PINKER**

# Training to Go Between Modes
# Example: Board Memory in Chess

[https://www.youtube.com/watch?v=rWuJqCwfjjc](https://www.youtube.com/watch?v=rWuJqCwfjjc)
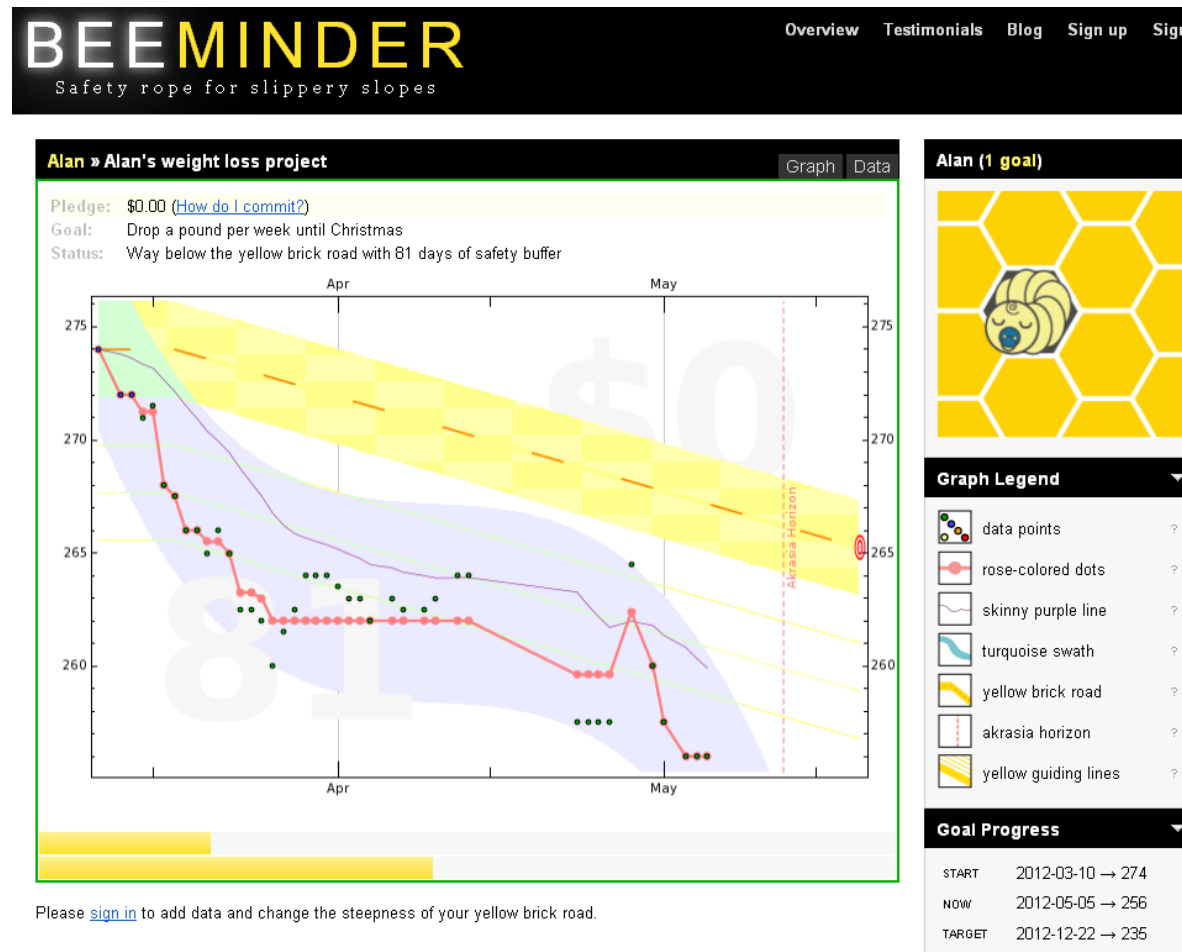
# Utilising Type I Thinking

- from behavioral science

- argues that positive reinforcement and indirect suggestions to try to achieve non-forced compliance can influence motives, incentives and decision making

- at least as effectively – if not more effectively - than direct instruction, legislation, or enforcement

# Application Example: Persuasive Technology (Beeminder)

# Application Example:
# Changing Driving Behaviour - Lottery Incentive



**Stanford | News**

Home    All News    Faculty & Staff News    For Journalists    About Us

Stanford Report, April 2, 2012

## Stanford study to try cold cash and social game to relieve rush hour traffic

Sleeping in might never feel better. To lower traffic congestion and pollution, a new program seeks to get Stanford drivers to avoid arriving and departing the campus during peak hours. Professor Balaji Prabhakar aims to deliver social benefits at low cost using people's penchant for a chance at a bigger payout over a predetermined small reward.

- Enter people in a lottery if they come to campus off-peak
- Observe that this small probability of earnings actually causes non-essential travellers to alter their behaviour

# Summary

- Much of the course progressed under the assumption of a fairly strict form of rationality
  - Nobody actually believes that this is what people always do
  - Instead, an implicit claim is that people have the capacity for this kind of reasoning

- In this lecture, we looked at how people systematically deviate from such a notion of rationality

- This has major implications for computational agents
  - In the end, you want them to work with real people
  - Modelling "real people" is challenging, but ultimately the *raison d'etre*

# Acknowledgements

Many of these  slides are adapted from Jurafsky et al.'s Symbolic Systems course at Stanford