# Adversarial Risk Analysis: Games and Auctions

David Banks

Duke University

# 1. Introduction

Classical game theory has focused upon situations in which outcomes are known. When uncertainty is addressed, it makes unreasonable assumptions about common knowledge (cf. Harsanyi, 1967/68a,b). Also, game theory makes unreasonable assumptions about human decision-making (Camerer, 2003).

Classical risk analysis has focused upon situations in which the hazards arise at random. This is appropriate for accident and life insurance, but it does not apply when hazards result from the actions of an intelligent adversary.

Corporate competition, federal regulation, and counterterrorism all entail game-theoretic problems with uncertain outcomes and partial information about the goals and actions of the opponents. This talk describes a Bayesian approach to adversarial risk analysis. It extends the decision analysis of Kadane and Larkey (1982) through the use of a **mirroring argument**.

Myerson (1991, p. 114) points up this problem clearly:

> "A fundamental difficulty may make the decision-analytic approach impossible to implement, however. To assess his subjective probability distribution over the other players' strategies, player $i$ may feel that he should try to imagine himself in their situations. When he does so, he may realize that the other players cannot determine their optimal strategies until they have assessed their subjective probability distributions over $i$'s possible strategies. Thus, player $i$ may realize that he cannot predict his opponents' behavior until he understands what an intelligent person would rationally expect him to do, which is, of course, the problem that he started with. This difficulty would force $i$ to abandon the decision analytic approach and instead undertake a game-theoretic approach, in which he tries to solve all players' decision problems simultaneously."

However, instead of following Myerson in defaulting back to game theory, we use the mirroring method. It may be viewed as a Bayesian version of Level-$k$ thinking (Stahl and Wilson, 1995).

# 2. Auctions

Suppose Apollo is bidding for a first edition of the Theory of Games and Economic Behavior. He is the only bidder, but the owner has set a secret reservation price $v^*$ below which the book will not be sold. Apollo does not know $v^*$, and expresses his uncertainty as a subjective Bayesian distribution $F(v)$.

Apollo's utility function is linear in money and his personal valuation of the book is $a^*$. If money is infinitely divisible, his choice set is $\mathcal{A} = \mathbb{R}^+$. so his expected utility from a bid of $a$ is $(a^* - a)\mathbb{P}[a > V^*]$. Thus Apollo should maximize his expected utility by bidding

$$a_0 = \mathbf{argmax}_{a \in \mathbb{R}^+}(a^* - a)F(a).$$

This is a standard approach in Bayesian auction theory (cf. Raiffa, 2002).

Now suppose that Apollo and Daphne are bidding against each other to own the first edition. Apollo needs to perform a game-theoretic calculation to find his subjective distribution $F$ over Daphne's bid $D_0$. Then Apollo can maximize his expected utility by bidding $a_0 = \operatorname{argmax}_{a \in \mathbb{R}^+} (a^* - a) F(a)$.

In order to find $F$, Apollo uses the fact that Daphne must make the symmetric calculation. This is the mirroring argument.

Specifically, suppose Daphne values the book at $d^*$ and has distribution $G$ on Apollo's bid $a_0$. Then Daphne would solve $d_0 = \operatorname{argmax}_{d \in \mathbb{R}^+} (d^* - d) G(d)$; and symmetrically, to obtain $G(d)$, Daphne would need to mirror Apollo's calculation.

But Apollo cannot duplicate Daphne's calculation since he does not know her value for the book, nor the value she thinks Apollo puts on the book, nor the value she thinks Apollo believes is her value for the book. As a Bayesian, Apollo must express his uncertainty on all three quantities through distributions.

The notation becomes complicated; the following key is helpful:

- $a^*$ is Apollo's value for the book

- $D^*$ is Daphne's value for the book; since it is unknown to Apollo, he assigns it the distribution $H_D$

- $A^*$ is the random variable that Apollo thinks Daphne uses to represent Apollo's value for the book; it has distribution $H_A$

- $F$ is Apollo's belief about the distribution of Daphne's bid.

- $D_0$ is Daphne's bid

- $G$ is Apollo's inference about Daphne's distribution on Apollo's bid.

- $A_0$ is Apollo's bid from Daphne's perspective.

These probabilities are all belong to Apollo; he imputes the beliefs that Daphne holds. If he is mistaken, he diminishes his chance of maximizing his gain.

To determine his bid $a_0$, Apollo needs $F$, the distribution of Daphne's bid. He knows that Daphne's bid $D_0$ should satisfy $D_0 = \mathrm{argmax}_{d \in \mathbb{R}^+}(D^* - d)G(d)$ where $D^*$ is Daphne's value (a random variable, to Apollo) for the book and $G(d)$ is Apollo's estimate of Daphne's probability that a bid of $d$ exceeds Apollo's bid $A_0$.

And, to Daphne, $A_0 = \mathrm{argmax}_{d \in \mathbb{R}^+}(A^* - a)F(a)$ where $A^*$ is Daphne's belief about Apollo's value for the book and $F(a)$ is Apollo's estimate of Daphne's probability that a bid of $a$ exceeds her bid $D_0$. Thus $D_0 \sim F$ and $A_0 \sim G$.

Apollo must find his personal belief about $F$ by solving:

$$\mathbf{argmax}_{d \in \mathbf{R}^+}(D^* - d)G(d) \quad \sim \quad F$$
$$\mathbf{argmax}_{a \in \mathbf{R}^+}(A^* - a)F(a) \quad \sim \quad G.$$

The distributions for $D^*$ and $A^*$ are $H_D$ and $H_A$, respectively.

Once Apollo has $F$, he solves $a_0 = \mathrm{argmax}_{a \in \mathbb{R}^+}(a^* - a)F(a)$ to determine his bid.

To solve this system of equations, one iteratively alternates between the two equations until convergence:

1. Select $F_0$ and $G_0$ arbitrarily.

2. Simulate a large number of samples from $H_A$, and solve the argmax problem under $G_i$. The distribution of those solutions gives $F_{i+1}$.

3. Simulate a large number of samples from $H_D$, and solve the argmax problem under $F_{i+1}$. The distribution of those solutions gives $G_{i+1}$.

4. If some convergence threshold $\delta$ is satisfied (e.g., $\|F_i - F_{i+1}\| < \delta$ and $\|G_i - G_{i+1}\| < \delta$), then stop. Otherwise, return to step 2.

In simulation, this round-robin algorithm has always converged. But one wants a fixed-point theorem, and the key issue is to show this iteration is a contraction operator. For a finite dimensional space (roughly corresponding to bids in pennies, rather than infinitely divisible money), this can be done in terms of Gauss-Siedel systems of equations.

This framework allows Apollo to incorporate secret information.

For example, suppose Apollo alone knows that the book was owned by Sir Ronald Fisher, with annotations in his hand. In that case, his personal value $a^*$ is high, but his distribution for Daphne's value, $H_D$, will concentrate on much smaller values.

Similarly, he might know that Daphne knows the provenance of the book but thinks that Daphne believes, falsely, that Apollo does not. In that case $H_D$ will give concentrate on large values, but Apollo's belief about what Daphne thinks is his value for the book, $H_A$, will concentrate on small values.
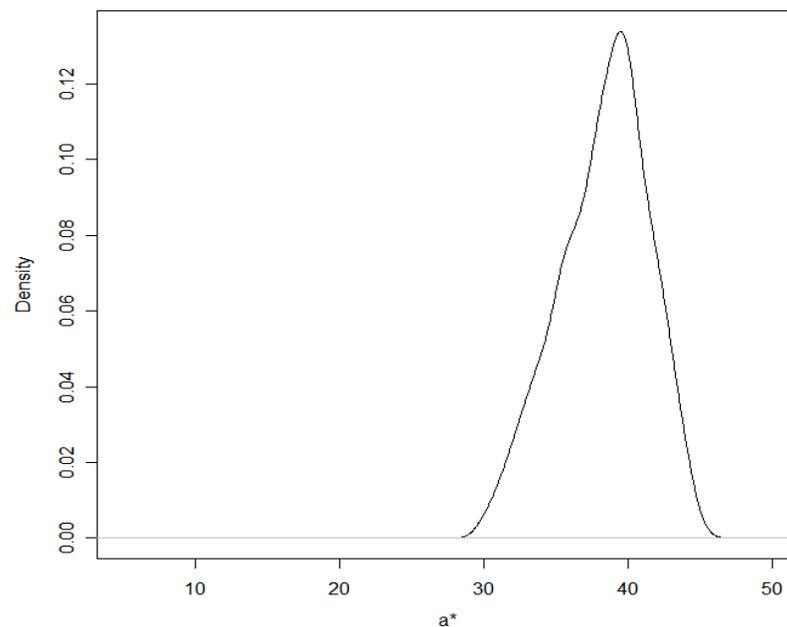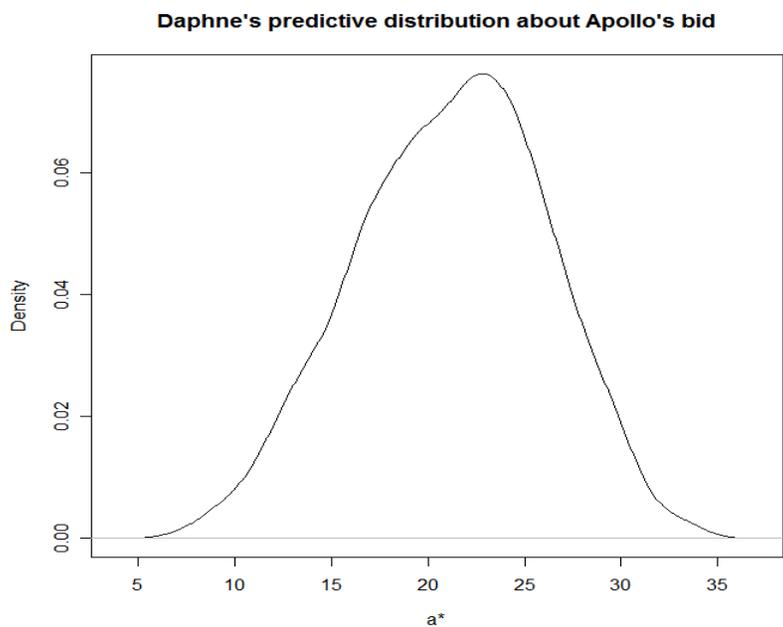
In principle, one could go into an infinite regress:

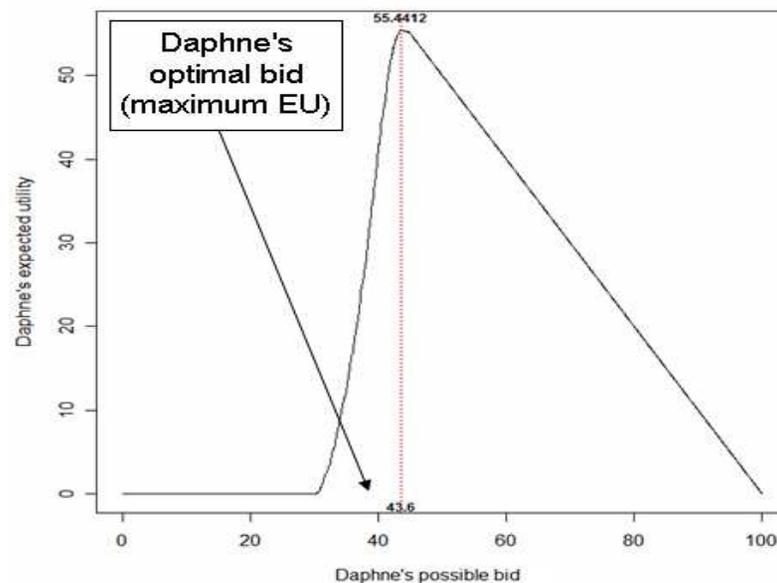<div align="center">

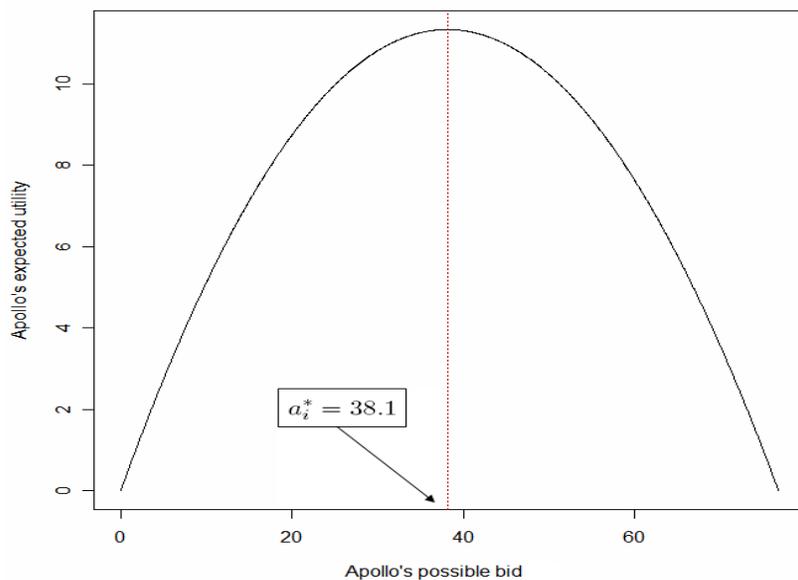## Apollo thinks that Daphne thinks that
### Apollo thinks that Daphne thinks that . . . .

</div>

But for human reasoning, it is probably quite reasonable to stop at the third step, with the distribution $H_A$ for $A^*$, as described in the mirroring analysis.

The following figures illustrate the fixed point solution. (Note that the caption reverses the roles of Apollo and Daphne.) The starting points for $H_D$ and $H_A$ were distinct triangular distributions on $[0, 100]$.

**Daphne's predictive distribution about Apollo's bid**



The left panel shows the third iterate; the right panel shows the tenth iterate.

These panels show the result of a algorithm. The left is the expected utility Daphne believes Apollo thinks he will get from a given bid. The right shows the expected utility that Daphne will receive from a given bid.

**Note:** Apollo makes the bid that maximizes his expected utility with respect to his true value for the book, not the random distribution he imputes to Daphne.

# 3. The Smallpox Decision

In 2002, U.S. policy-makers were concerned that terrorists might launch a bioterrorist attack with smallpox. They considered three scenarios:

- A major attack, involving multiple cities and weaponized virus.

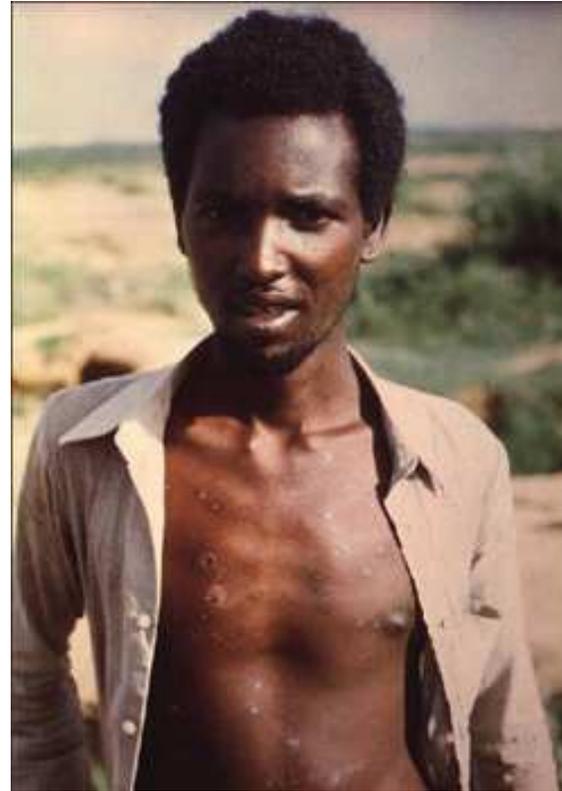- A minor attack, similar to the anthrax letters.

- No smallpox attack.

And the kinds of defenses under consideration were:

- Stockpiling vaccine

- Stockpiling vaccine and increasing biosurveillance

- Stockpile, increase biosurveillance, inoculate all first responders

- Inoculate essentially everyone.

Smallpox facts:

- In 1967, WHO estimated that there were 15 million cases, with 2 million deaths.

- In 1979, WHO certified the eradication of smallpox (the effort cost $300 million).

- In developed nations, the most recent natural outbreak was in Yugoslavia in 1972. In the U.S., the last outbreak had been in NYC in 1947; 6.3 million people were inoculated within three weeks.

- There are two forms: Variola major (hemorrhagic) and Variola minor. After infection, the incubation period lasts about 12 days, and the disease becomes obvious and infectious on the 12-15th day. By the 28th day, the patient is dead or recovering.

- Smallpox is less infectious that measles or influenza. In developed nations, no recent epidemic has progressed beyond the second generation.

- The fatality rate from Variola major is between 3% and 35%. For Variola minor is about 1% with modern treatment.

The last natural case of hemorrhagic smallpox killed Rahima Banu, a two-year-old girl in Bangladesh in 1975. The last natural case of <u>Variola minor</u> was a Somali cook named Ali Maow Maalin, who survived.



In 1978 there was an accidental release from a Birmingham laboratory. A medical photographer died. This led to the elimination all stocks except for archival samples maintained in the U.S. and Russia.

In game-theoretic terms, the payoff matrix for this problem is shown below, where $C_{ij}$ represents the cost (or benefit) to the U.S., and $D_{ij}$ represents the corresponding cost to the terrorists.

|  | No Attack | Minor Attack | Major Attack |
|---|---|---|---|
| Stockpile | $(C_{11}, D_{11})$ | $(C_{12}, D_{12})$ | $(C_{13}, D_{13})$ |
| Biosurveillance | $(C_{21}, D_{12})$ | $(C_{22}, D_{22})$ | $(C_{23}, D_{23})$ |
| First Responders | $(C_{31}, D_{31})$ | $(C_{32}, D_{32})$ | $(C_{33}, D_{33})$ |
| Mass Inoculation | $(C_{41}, D_{41})$ | $(C_{42}, D_{42})$ | $(C_{43}, D_{43})$ |

A game theorist would replace the random payoffs by their average values, and then use the minimax theorem to find the optimal play. **But the best play for the average case is not the same as the play that is, on average, the best.**

A risk analyst would assume that terrorists attack at random, like the weather.

The ARA approach is different: it uses mirroring to model the decision processes of the terrorist, while taking account of the fact that the terrorist is modeling the defender.

As in the auction example, the Defender has two distributions. The first is the distribution that describes the Terrorist's beliefs about the random payoff matrix. The second describes what the Terrorist thinks is the Defender's beliefs about the random payoff matrix.

This leads to a pair of coupled equations, which the Defender solves in order to determine the probability distribution over the Terrorist's decision. Then the Defender makes the choice that maximizes his expected utility with respect to that probability distribution.

Obviously, neither Terrorists nor Defenders actually solve these equations. But the equations formalize the kind of strategizing that humans often perform more heuristically.

The Defender's two distributions are obtained by traditional probability elicitation from experts.

For the distribution on what the Defender thinks is the Terrorist's payoff matrix, the Defender may use information from informants, expert opinion on terrorist psychology, and information about their resources and goals.

For example, initially the Defender might think that it would be prohibitively resource-intensive for the Terrorist to acquire smallpox. But if intelligence suggests it has been acquired, then the updated payoff matrix distribution changes to indicate that a smallpox attack has low prospective cost and high value.

The other distribution, the one that the Terrorist has on the Defender's payoff matrix, may be even more straightforward to obtain. In an open society, the Terrorist has access to essentially the same cost information that the Defender does.

The Defender might imagine that the Terrorist is modeling the distribution of costs for $C_{11}$ in the top-left cell. It represents the cost to the U.S. of stockpiling vaccine when, in fact, no smallpox attack is made.

$$C_{11} = \text{cost to test diluted Dryvax} + \text{cost to test Aventis vaccine}$$
$$+ \text{cost to make 209 x 106 doses} + \text{cost to produce VIG}$$
$$+ \text{logistic/storage/device costs.}$$

The Dryvax and Aventis vaccines had been kept in storage for years. Their residual potency was unclear; testing was needed. Two experts pooled their opinions and decided that the testing cost might be uniformly distributed between $2 million and $5 million.

The FDA said that new vaccine production was fixed by contract at $512 million. The VIG cost is modelled as N($100 million, $20 million$^2$). Logistics costs are distributed as N($940 million, $100 million$^2$).

Preliminary research with domain experts provided approximate distributions for the costs in the Defender's payoff matrix. Other terms that experts assessed included:

- The value of a human life was treated as fixed, at $750,000. (This follows the DOT human capital model; non-market methods tend to give higher values.)

- The number of key personnel to be inoculated: this was guessed to be uniform between .5 million and .6 million.

- The number of smallpox cases per attack: this was guessed to be gamma with mean 10 and sd 100.

- The cost to treat one smallpox case: normal with mean $200,000 and sd $50,000.

- The economic costs of an attack: gamma with mean $5 billion, sd $10 billion.

Expert elicitation is always problematic, but for practical reasons it has been embraced by DHS, CREATE, and other counterterrorism analysts.

**Note:** The different costs in the matrix are correlated. If the stockpiling costs turn out to be higher than expected, those higher costs should also appear as a summand in every cell in that same row. The payoff table is a random matrix with a complex correlation structure.

Using these toy models for costs, to both the Defender and the Terrorist, we generated 10,000 random bivariate payoff matrices and solved each as a non-zero-sum game. Some solutions were randomized strategies, but most were not.

Averaging these solutions gave estimates of the probability that particular actions would be taken by the Terrorist, and estimates of what we believe is the Terrorist's probabilities for our actions.

Under the plausible but inexact elicitations we used, the Terrorist will not mount a smallpox attack with probability .99; the chance of minor attack is .002, a major attack is .008.

We also estimate that the Terrorist thinks there is probability .21 that the Defender will stockpile vaccine, probability .23 that the Defender will stockpile and increase biosurveillance, probability .29 of stockpiling, surveilling, and inoculating first responders, and probability .26 of mass inoculation.

**In this scenario, our expected utility is maximized by stockpiling.**

# 4. Conclusions

Adversarial risk analysis is an important combination of tools from game theory and statistical risk analysis. In particular, Bayesian versions of ARA provide attactive alternatives to the known deficiencies of current methods.

In auctions, gambling, and counterterrorism, agents often have mental models for the decision-processes of their opponents. If those models are correct, then there is the opportunity to improve on the innate pessimism of minimax solutions.

A key component is the need to properly handle many different kinds of uncertainty, which arise in different parts of the analysis. The toy problems considered in this talk point to some of the issues, in the context of normal form and extensive form games.

For more complex applications, say in counterterrorism, work is needed to account for the effects of decentralized and tiered decision-making, constraints on the resources of the opponents, and sensitivity analysis to the Bayesian beliefs and elicitations.