# On the Relation Between Damage and Deception

Joel Sobel

August 29, 2023

# Three Tasks

1. "Pedagogic": Provide coherent definition(s) of deception.
2. "Applied": Relate the definition to examples faced by regulators.
3. "Technical": Propositions that relate binary relationships on beliefs to preferences of decision makers.

# Framework

- ▶ One Decision Maker.
- ▶ Beliefs $\mu \in \Delta(\Theta)$.
- ▶ Decision maker takes action $y$.
- ▶ Utility function: $U^R(\theta, y)$.
- ▶ For convenience: $\Theta$ finite.

# Standard Interpretation

Two players: Sender and Receiver.

Sender observes $\theta \in \Theta$.

Sender sends message $m \in M$. ($M$ can depend on $\theta$.)

$R$ observes $m$.

Receiver takes action $y \in Y$.

Preferences $U^i(\theta, y, m)$; $U^R(\cdot)$ independent of $m$.

Prior $P(\theta)$ (positive on $\Theta$).

# Explanation

Standard interpretation

1. ... introduces game theory.
2. ... explains where beliefs come from.
3. ... relevant for applications.

# Deception: Informal

- ▶ Deception is "inducing bad beliefs."

    I need some notion of beliefs.

    I need some notion of "bad."

- ▶ How to do this?

    Let $D(\mu)$ is the set of beliefs that are less accurate/more deceptive than $\mu$.

- ▶ So beliefs are bad compared to other possible beliefs.

Easy Case: Two states. Beliefs totally ordered (by probability placed on true state).

# Conceptual Approach: Deception

▶ Deception is partial order on beliefs. Loosely "more deceptive" means "further from truth."
($\mu$ more deceptive than $\mu'$ ...)

▶ There will be lots of definitions because there are lots of distances from truth.

▶ To evaluate the relevance of a particular definition, I relate it to when "bad" beliefs lead to "bad" utility.

# Conceptual Approach: Damage

- ▶ Damage is a(nother) partial order on beliefs. Loosely "more damaging" means "leads to lower utility."
- ▶ There will be lots of definitions because there are lots different kinds of agent/utility function.
- ▶ Goal: "Damage-Deception Result" associating a definition of deception with a class of preferences.

# Strong ($S$) Deception

### Definition (Strong deception)

The belief $\mu'$ is more **strongly deceptive** than $\mu$ given $\theta^*$ if $\mu' \neq \mu$ and there exists $p \in [0, 1)$ such that

$$\mu(\cdot) = p\mu'(\cdot) + (1 - p)I(\cdot \mid \theta^*). \tag{1}$$

$\mu$ on segment connecting point mass on $\theta^*$ to $\mu'$.

$[I(\cdot \mid \theta^*)$ is point mass on $\theta^*$.]
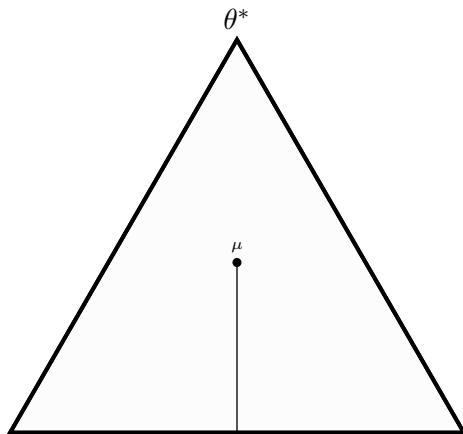
# S-Deception Illustrated



Figure: Beliefs on the line segment are more strongly deceptive than $\mu$.
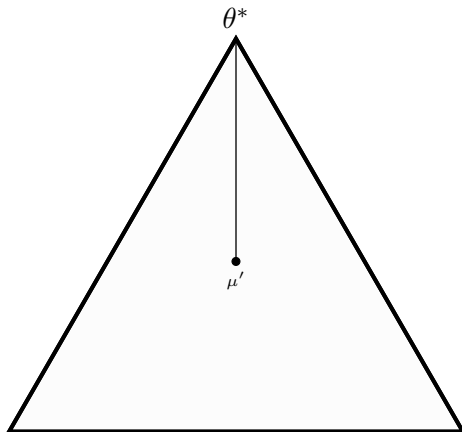
# S-Deception Illustrated Again



Figure: The belief $\mu'$ is more strongly deceptive than $\mu$ if only if $\mu$ is on line segment.

# Damage: Preliminaries

Deception: Property of Beliefs

Damage: Consequences. Property of Preferences.

Assume $R$ best replies to beliefs.

$BR(\mu)$ is $R$'s best response correspondence:

$$BR(\mu) = \arg\max_{y \in Y} \sum_{\theta \in \Theta} U^R(\theta, y)\mu(\theta).$$

# Damage: Definition

Let $\bar{u}(\theta, \mu) = U^R(\theta, BR(\mu))$.

### Definition (Damaging Behavior)

The belief $\mu'$ is more **damaging** than $\mu$ given $\theta^*$ if

$$\bar{u}(\theta^*, \mu') < \bar{u}(\theta^*, \mu).$$

Technicality:

▶ If $BR(\mu)$ is not single valued, then $\bar{u}(\theta^*, \mu')$ may not be single valued.

▶ So $\bar{u}(\theta^*, \mu') < \bar{u}(\theta^*, \mu)$ needs a definition.

▶ I need to make some assumption (otherwise ugly statements of propositions).

▶ In what follows, can rank sets using strong set order, maximum, minimum, . . . .

# Damaging Relative to a Set

Let $\mathcal{U}$ be a set of payoff functions for $R$ (real valued functions of $(\theta, y)$).

## Definition (Damaging Relative to a Set)

The belief $\mu'$ is more **damaging relative to** $\mathcal{U}$ than $\mu$ given $\theta^*$ if

1. $\bar{u}(\theta^*, \mu') \leq \bar{u}(\theta^*, \mu)$ for all $u \in \mathcal{U}$
2. $\bar{u}(\theta^*, \mu') < \bar{u}(\theta^*, \mu)$ for some $u \in \mathcal{U}$.

The same technicality about "$<$" applies.

# Where am I going?

1. $\mathcal{U}$ generates partial order on beliefs.
2. I want to compare these partial orders to definitions of deception.

# Aside: FTC

The Federal Trade Commission identifies three necessary conditions for deception.

1. Deception requires doing something that misleads the consumer.
2. FTC evaluates the impact from the perspective of a consumer who acts reasonably.
3. For a practice to be deceptive it must have a material impact on a consumer.

# FTC versus Me

- ▶ My notion of deception concentrates on "misleading" information.
- ▶ My formal results connect "material impact" (damage) to misleading information.

# Generic Proposition Statement

### Proposition

*The belief $\mu'$ is more X deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging than $\mu$ relative to the $F(X)$ family of preferences given $\theta$.*

The talk provides specific $(X, F(X))$ pairs.

Teaser:
When $X =$ strong deception, $F(X) =$ all preferences.

# Haven't I seen this before?

1. $\mu$ FOSD $\mu'$ if and only if all decision makers with increasing utility functions prefer $\mu$.
2. $\mu$ SOSD $\mu'$ if and only if all decision makers with concave utility functions prefer $\mu$.

# Contrast

My approach is interim, not ex ante.

My approach involves a decision.

# Persuasion Framework

- $R$'s decision is whether to accept or reject a proposal.
- $Y = \{0, 1\}$. $y = 0$ reject; $y = 1$ accept.
- $A^* = \{\theta : U^R(\theta, 1) \geq U^R(\theta, 0)\}$. Acceptable states for $R$.
- Persuasion preferences (with respect to $A^*$):

## Definition (Persuasion Preferences)

$R$ has **persuasion preferences** if $U^R(\theta, 0) = 0$ and

$$U^R(\theta, 1) = \begin{cases} W & \text{if } \theta \in A^* \\ -L & \text{if } \theta \notin A^* \end{cases}$$

for $W, L > 0$.

# Persuasion Result

### Definition (Binary-Action Deception)

The belief $\mu'$ is more **binary-action (BA)-deceptive** than $\mu$ given $\theta^*$ if

$$\text{if } \theta^* \in A^*, \text{ then } \mu'(A^*) < \mu(A^*)$$
$$\text{and if } \theta^* \notin A^*, \text{ then } \mu'(A^*) > \mu(A^*).$$

### Proposition

*The belief $\mu'$ is more binary-action-deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging than $\mu$ given $\theta^*$ relative to the family of persuasion preferences.*

When $X$ = binary-action deception, $F(X)$ = persuasion preferences.

# S-Deception Result

### Proposition

*The belief $\mu'$ is more strongly deceptive given $\theta^*$ than $\mu$ if and only if it is more damaging than $\mu$ given $\theta^*$ relative to the set of all preferences.*

$X =$ strong deception
$F(X) =$ all preferences

# Proportional (P) Deception

### Definition (Proportional Deception)

The belief $\mu'$ is more **proportional (P)-deceptive** than $\mu$ given $\theta^*$ if $\mu(\theta^* \mid n) > 0$, and there exists a number $p \in [0, 1)$, and a distribution $\rho$ satisfying $\rho(\theta^*) = 0$ such that

$$\mu'(\cdot) = p\mu(\cdot) + (1 - p)\rho. \qquad (2)$$

Equivalent to

$$\frac{\mu(\theta^*)}{\mu(\theta)} \geq \frac{\mu'(\theta^*)}{\mu'(\theta)}$$

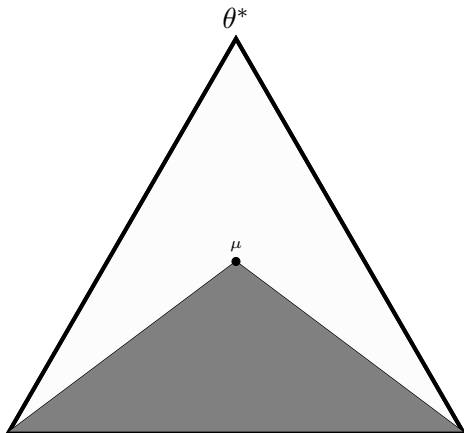($\theta^*$ relatively more likely under $n$.)

# *P*-Deception Illustrated



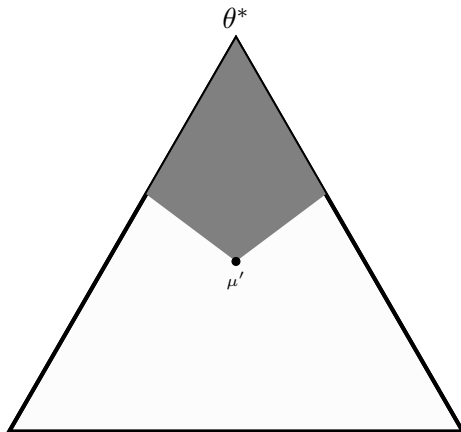Figure: Beliefs in shaded region are more proportionally deceptive than $\mu$ given $\theta^*$.

# Again . . .



Figure: The belief $\mu'$ is more proportionally deceptive than $\mu$ if and only if $\mu$ is in shaded region.

# *P*-Deception Result

### Definition (State Specific)

The Receiver's preferences are **state specific** if there is a bijection $\phi : \Theta \to Y$ and positive numbers $\alpha(\theta)$ for $\theta \in \Theta$ such that

$$U^R(\theta, y) = \begin{cases} \alpha(\theta) & \text{if } y = \phi(\theta) \\ 0 & \text{if } y \neq \phi(\theta) \end{cases}$$

Corresponds to situation in which there is a correct action for each state.

### Proposition

*The belief $\mu'$ is more proportionally deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging than $\mu$ given $\theta^*$ relative to the family of state-specific preferences.*

# Variations

1. Constrained
2. Kullback–Leibler

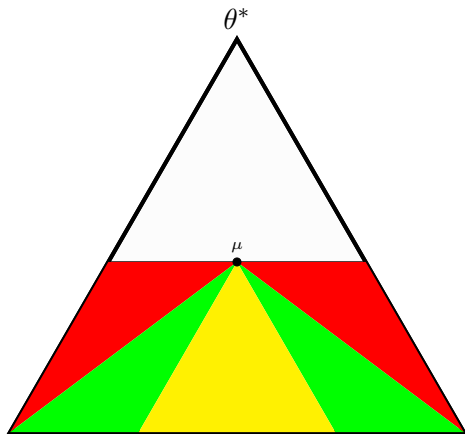Ordered states: Monotone

# Nesting Illustrated



Figure: Relative to $\mu$ and given $\theta$, beliefs in red region are *KL*-deceptive but not *P*-deceptive. Beliefs in green region are *P*-deceptive but not *C*-deceptive. Beliefs in green region are *C*-deceptive but not *S*-deceptive.

# Strategic Considerations

1. Imagine $S$ sends messages, messages determine beliefs.
2. Look for equilibria of persuasion game.
3. Define partial order on messages based on beliefs/actions induced.

   $m'$ is more deceptive than $m$ if and only if $\mu(\cdot \mid m')$ is more deceptive than $\mu(\cdot \mid m)$.

   $m'$ is deceptive if there exists $m$ such that $m'$ is more deceptive than $m$.

4. Must specify $S$ preferences.

# Comments

1. $R$ can avoid being deceived. [No deception/damage if $R$'s beliefs and actions don't depend on $S$'s message.]

2. $R$ may prefer (ex ante) an eq in which he is deceived.

3. If $S$ always wants to persuade, talk is cheap, prior is favorable, and some message induces $y = 0$, then damage and deception.

4. If $S$ sometimes does not want to persuade, then some message will induce $y = 0$.

# Practical Considerations

Definition of deception depends on context:

1. Persuasion Preferences: $R$ decides whether to buy.
2. State Specific Preferences: $R$ decides which product to buy.
3. Monotone Preference: $R$ decides how much to buy.
4. Strong Deception: Conservative.

## Comments on Welfare

1. Large penalties for deception benefit $R$ because $R$ can "force" fully revealing outcome.
2. This property depends on equilibrium selection.
3. Large penalties for deception may benefit $R$ and $S$.
4. ... but not if one selects babbling when an ex ante superior equilibrium exists.
5. Restrictions on deception may harm both if it is costly to disclose (by driving firms out of business).

# Theoretical Properties

1. For which binary relations is there a damage-deception result?
2. Minimal relation?
3. Maximal relation?

# Back to the Three Tasks

1. "Pedagogic": Mission accomplished (well, too many definitions)

2. "Applied": Only suggested in this talk.

   Messages:
   - ▶ Proper definition of deception depends on the context (what is known about $R$'s preferences).
   - ▶ Deception is possible in equilibrium.
   - ▶ An $R$ who can be deceived may also be an $R$ who may benefit from communication.

3. "Technical": Dam-Dec results. Deception relation more complete corresponds to smaller families of associated preferences.
   Harming $R$ is a consequence of deception rather than part of the definition.

# Tribute

gpt

# Supplementary Material

Follows

# General Properties

- There exist (non-trivial) deception correspondences.
- The deception relations are transitive.

# More Properties

► Deception relations are not symmetric:

## Proposition

*If $D$ is a deception correspondence with respect to $\mathcal{U}$ given $\theta^*$, then $\mu' \in D(\mu'')$ implies $\mu'' \notin D(\mu')$.*

► Deception correspondences are not convex valued in general, but

## Proposition

*If $D$ is a deception correspondence with respect to $\mathcal{U}$ given $\theta^*$, $\mu \in D(\mu^*)$, and $\mu = p\mu' + (1 - p)I(\cdot \mid \theta^*)$ for $p \in (0, 1]$ then $\mu' \in D$.*

# Constrained ($C$) Deception

## Definition (Constrained deception)

The belief $\mu'$ is more **constrained ($C$)-deceptive** than $\mu$ given $\theta^*$ if $\mu(\theta^*) > 0$ and either

$$\mu'(\theta^*) = 0$$

or there exists a number $p \in [0, 1)$, a distribution $\rho$ satisfying $\rho(\theta^*) = 0$ such that

$$\mu'(\cdot) = p\mu(\cdot) + (1-p)\rho$$

and

$$\rho(\theta) \geq \mu(\theta) \text{ for all } \theta \neq \theta^*.$$

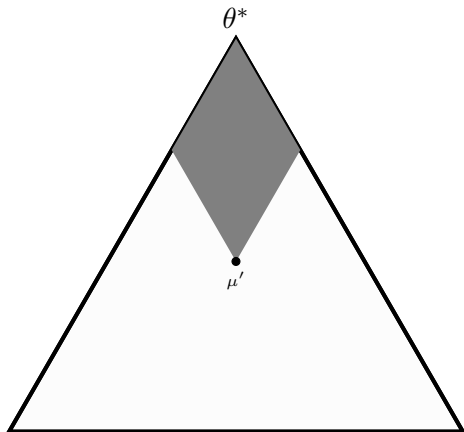$\mu'$ lowers probability of true state and raises it on others.

# C-Deception Illustrated



Figure: The belief $\mu'$ is more constrained deceptive than $\mu$ if and only if $\mu$ is in the shaded region.
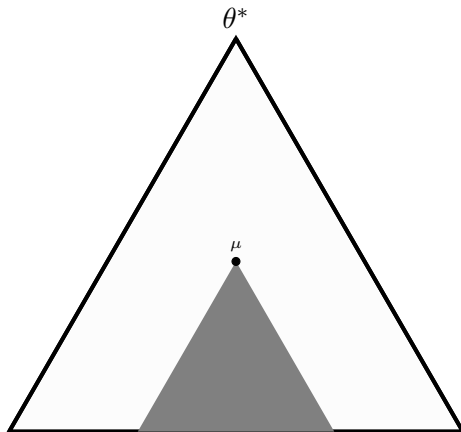
# C-Deception Again



Figure: Beliefs in shaded region are more constrained deceptive than $\mu$ given $\theta^*$.

# C-Deception Result

### Definition (Linear Family with outsider option)

If $R$ has preferences in the family of linear preferences with an outside option $u_0$, then $R$ selects an action $y \in \Delta(\Theta \cup [0,1])$ to maximize

$$\sum_\theta \mu(\theta)\frac{y(\theta)}{\beta(\theta)} + y_0 u_0.$$

### Proposition

*The belief $\mu'$ is more constrained deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging than $\mu$ relative to the family of linear preferences with an outside option given $\theta^*$.*

Nesting

# Kullback-Leibler (KL) Deception

Kullback–Leibler divergence between distributions $\mu$ and $\mu'$:

$$D_{KL}(\mu' \mid\mid \mu) = \sum \mu'(\theta) \log \frac{\mu'(\theta)}{\mu(\theta)}.$$

[Assume that $\mu(\theta) = 0$ implies $\mu'(\theta) = 0$ and follow the convention that $x \log x = 0$ when $x = 0$.]

### Definition (KL-Deception)

The belief $\mu'$ is more **Kullback-Leibler (KL)-deceptive** than $\mu$ given $\theta^*$ if $\mu(\theta^* > 0$ and

$$D_{KL}(I(\cdot \mid \theta^*) \mid\mid \mu'(\cdot) < D_{KL}(I(\cdot \mid \theta^*) \mid\mid \mu'(\cdot).$$

# Simpler

A message is KL-deceptive if and only if there is another message that induces beliefs placing higher probability on the true state.
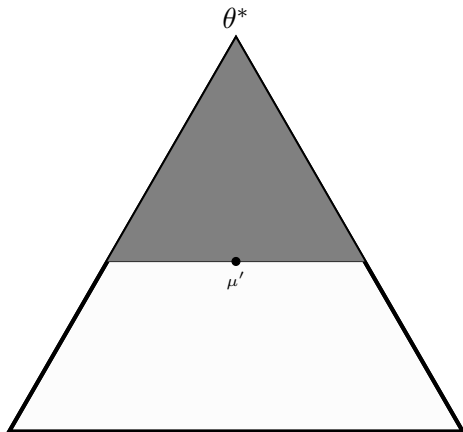
# KL-Deception Illustrated



Figure: The belief $\mu'$ is more KL-deceptive than $\mu$ given $\theta^*$ if and only if $\mu$ is in shaded region.
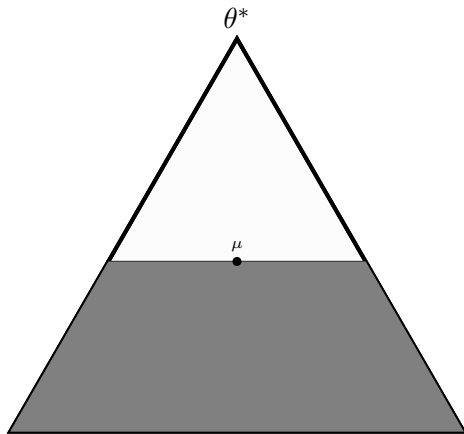
# Again



Figure: Beliefs in shaded region are more *KL*-deceptive than $\mu$ given $\theta^*$.

# *KL*-Deception Result

### Definition (Exponential Family)

The family of preferences is logarithmic if $g(y) = \log y$.

### Proposition

*The belief $\mu'$ is more Kullback-Leibler deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging than $\mu$ relative to the family of logarithmic preferences given $\theta^*$.*

Nesting

# Preliminary

1. $Y$ and $\Theta$ linearly ordered.
2. Elements of $\Theta$ denoted by $\theta_1, \ldots, \theta_N$ where $\theta_i < \theta_j$ if and only if $i < j$.

### Definition (Increasing Differences)

Assume that $\Theta$ and $Y$ are completely ordered. The function $u : \Theta \times Y \to \mathbb{R}$ satisfies **increasing differences** if $u(\theta_j, y) - u(\theta_i, y)$ is increasing in $y$ whenever $j > i$.

Given a distribution $\mu$, denote by $C(j; \mu)$ the cumulative probability determined by $\mu$. That is, $C(j; \mu) = \sum_{i \leq j} \mu(\theta_i)$.

### Definition (First-Order Stochastic Dominance)

$\mu$ **dominates** the probability distribution $\mu'$ if $\mu \neq \mu'$ and $C(j; \mu') - C(j; \mu) \geq 0$ for all $j$.

# Monotone Deception Defined

### Definition (Monotone Deception)

The belief $\mu'$ is more **monotonically ($M$)-deceptive** than $\mu$ given $\theta^*$ if one of the following conditions hold:

1. $\mu'$ is more strongly deceptive than $\mu$ given $\theta^*$;
2. $\mu$ dominates $\mu'$ and $\mu(\theta) = \mu'(\theta) = 0$ for $\theta > \theta^*$;
3. $\mu'$ dominates $\mu$ and $\mu(\theta) = \mu'(\theta) = 0$ for $\theta < \theta^*$.

# Comments

1. The first condition holds and $\mu(\theta)) = \mu'(\theta) = 0$ for $\theta > \theta^*$, then the second condition holds.
2. Likewise for the third condition.
3. Similarly, if the first condition holds and $\mu(\theta) = \mu(\theta) = 0$ for $\theta < \theta^*$, then the third condition holds.
4. So: $\mu'$ is more $M$-deceptive than $\mu$ if:
   4.1 $\theta^*$ is neither the highest nor the lowest state in the support of $\mu'$ and $\mu$ and $\mu'$ is more strongly deceptive;
   4.2 $\theta^*$ is the lowest state given positive probability by $\mu$ and $\mu'$ is stochastically greater.
   4.3 $\theta^*$ is the highest state given positive probability by $\mu$ and $\mu'$ is stochastically lower.
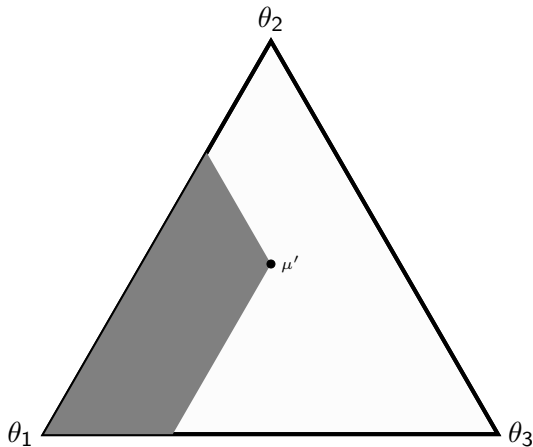
# M deception illustrated



Figure: Belief $\mu'$ is more M-deceptive than $\mu$ given $\theta_1$ if and only if $\mu$ is in shaded region.
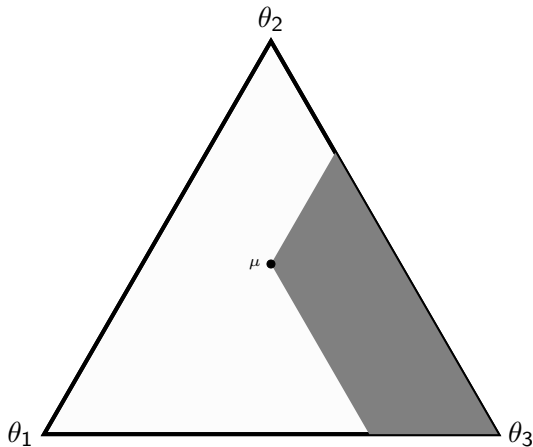
# M deception illustrated again



Figure: The shaded area is the set of beliefs that are more M-deceptive than $\mu$ given $\theta_1$.

# *M* Deception Result

### Proposition

*The belief $\mu'$ is more monotone deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging relative to the class of ID preferences.*

[*ID* preferences are concave, increasing, and satisfy increasing differences.]

# Remarks

1. $M$ deception restrictive (because strong deception is).
2. But: $\theta^* = \theta_1$ case is of particular interest.
   - $R$ would not purchase the item knowing given $\theta_1$.
   - $S$ would like to convince $R$ to buy.
   - Result associates damage with convincing $R$ to buy more than he wants.
   - $m$ is deceptive if it "exaggerates" the true state and there is a more moderate exaggeration available.
3. Different classes of preferences:
   3.1 Tail states don't changes optimal action.
   3.2 Quadratic loss.

Nesting

# Zenith Deception

### Definition (Zenith ($Z$) Deception)

The belief $\mu'$ is more **zenith deceptive** than $\mu$ given $\theta^*$ and $\mu$ if $\mu'(\theta^*) < \max_\theta \mu'(\theta)$ and $\mu(\theta) = \max_\theta \mu(\theta)$.

There is an $\mu$ makes $\theta^*$ most likely (and $\mu'$ don't make $\theta^*$ most likely).
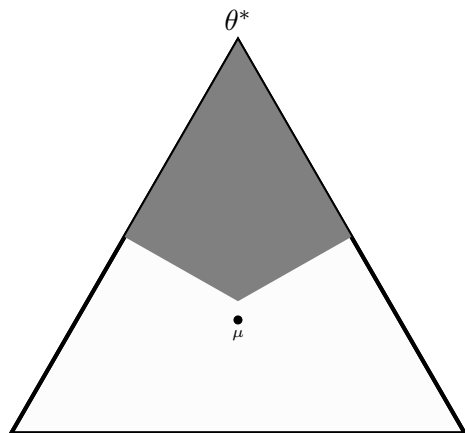
# Z-Deception Illustrated



Figure: The belief $\mu'$ is more zenith deceptive than $\mu$ given $\theta^*$ if and only if $\mu$ is in the shaded region. ($\mu(\theta^*) = \max_\theta \mu(\theta)$, and $\mu'(\theta^*) < \max_\theta \mu'(\theta)$.)
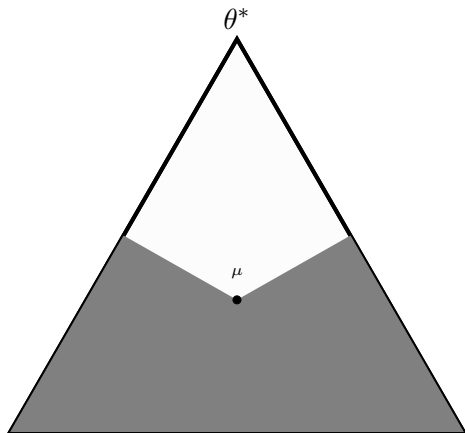
Figure: Beliefs in shaded region are more zenith deceptive than $\mu$ given $\theta^*$. ($\mu(\theta^*) = \max_\theta \mu(\theta)$.)

# Z-Deception Result

### Definition (Uniformly Linear Family)

The family $(g, \beta)$ of preferences is uniformly linear if $g(\cdot)$ is linear $c > 0$ such that $\beta(\theta) = c$ for all $\theta$.

### Proposition

*The belief $\mu'$ is more zenith deceptive than $\mu$ given $\theta^*$ if and only if it is more damaging than $\mu$ relative to the family of uniform linear preferences given $\theta^*$.*

Nesting

# Volkswagen

- Volkswagen advertised that they produced diesel cars that were not dangerous to the environment.
- Vehicles were equipped with illegal emission defeat devices during government tests. (Not disclosed.)
- Message: announcements about safety.
- Alternative message: truth, saying nothing.
- Decision: whether to buy.
- Deception: thinking the car was environmentally friendly.
- Damage: buying wrong car in response to the message.

Comments

# Machinima

- ▶ Machinima paid people who created videos posted on youtube to include Xbox footage in their reviews.
- ▶ Machinima asked the youtubers not to reveal payments.
- ▶ FTC argued that it was misleading to represent paid endorsers as independent reviewers.
- ▶ No evidence that the reviews were false, but FTC claimed withholding information about payments may influence the interpretation of the videos.
- ▶ Message: ads
- ▶ Alternative: disclosing payments (or not making them)
- ▶ Decision: What system to buy
- ▶ Deception: Unjustified confidence in quality of system
- ▶ Damage: Purchase of wrong system

Comments

# POM Wonderful

▶ The ads claimed that the POM (pomegranate juice) could prevent or reduce the risk of heart disease, prostate cancer, and erectile dysfunction.

▶ POM provided supporting evidence.

▶ FTC asserted evidence was inadequate (lacking proper controls; statistically insignificant results).

▶ Message: ads

▶ Alternative message: comments on other characteristics; complete description of evidence.

▶ Decision: whether to buy (or how much to buy).

▶ Deception: inaccurate impression about health benefits.

▶ Damage: buying too much.

Comments

# Kellogg's Mini-wheats

- ▶ Ads claimed that children who ate Frosted Mini-Wheats were 20% more attentive than those who skipped breakfast.
- ▶ Kellogg referred to a study to back up the claims.
- ▶ FTC argued that claims, while not literally false, were misleading. (Half of subjects showed no increase; only 10% significant gains.)
- ▶ Message: ads
- ▶ Alternative message: comments on other characteristics; complete description of evidence
- ▶ Decision: whether to buy (or how much to buy)
- ▶ Deception: inaccurate impression about benefits.
- ▶ Damage: buying too much.

Comments

# Red Bull

- ▶ Red Bull energy drink does **not** give you wings (literally).
- ▶ Law suit (and settlement) based on lack of evidence that it gives you (figurative) wings.

Comments

# Chat GPT

1. Digital Economics and Platform Markets
2. Behavioral Economics and Nudging
3. Environmental and Resource Economics
4. Health Economics and Healthcare Markets
5. Innovation and Intellectual Property
6. Data Economics and Privacy
7. Development Economics
8. Economics of Information and Learning

Tribute