# Mechanisms for a No-Regret Agent: Beyond the Common Prior



**Modibo Camara**  Jason Hartline  Aleck Johnsen

Northwestern University

Cornell ORIE Seminar

Appeared in the 2020 Symposium on Foundations of Computer Science (FOCS)

# Introduction

## Motivation

Policy success or failure often influenced by...

- ▶ Environmental details (e.g. consumer demand, labor supply)
- ▶ Individuals' beliefs about environment (e.g. inflation expectations)

Inherently dynamic

- ▶ Economic conditions evolve in unpredictable ways
- ▶ Individuals learn and adapt over time

Ideal policies would adapt to environment over time

- ▶ E.g. minimum wage adapted to current, local market conditions

## Motivation

Policy success or failure often influenced by...

- ▶ Environmental details (e.g. consumer demand, labor supply)
- ▶ Individuals' beliefs about environment (e.g. inflation expectations)

Inherently dynamic

- ▶ Economic conditions evolve in unpredictable ways
- ▶ Individuals learn and adapt over time

Ideal policies would adapt to environment over time

- ▶ E.g. minimum wage adapted to current, local market conditions

## Motivation

Policy success or failure often influenced by...

▶ Environmental details (e.g. consumer demand, labor supply)
▶ Individuals' beliefs about environment (e.g. inflation expectations)

Inherently dynamic

▶ Economic conditions evolve in unpredictable ways
▶ Individuals learn and adapt over time

Ideal policies would adapt to environment over time

▶ E.g. minimum wage adapted to current, local market conditions

## High Level Questions

**Can we develop dynamic policies that adapt to their environment over time?**

Without making assumptions on how the environment evolves?

  ▸ I.e. prior-free or adversarial

With permissive assumptions on agent behavior?

  ▸ Ex-ante optimal: Bayesian agents *want to* satisfy assumptions

  ▸ Ex-post feasible: non-Bayesian algorithm guaranteed to satisfy assumptions

## High Level Questions

**Can we develop dynamic policies that adapt to their environment over time?**

Without making assumptions on how the environment evolves?

▶ I.e. prior-free or adversarial

With permissive assumptions on agent behavior?

▶ Ex-ante optimal: Bayesian agents *want to* satisfy assumptions

▶ Ex-post feasible: non-Bayesian algorithm guaranteed to satisfy assumptions

## High Level Questions

**Can we develop dynamic policies that adapt to their environment over time?**

Without making assumptions on how the environment evolves?

- ▶ I.e. prior-free or adversarial

With permissive assumptions on agent behavior?

- ▶ Ex-ante optimal: Bayesian agents *want to* satisfy assumptions
- ▶ Ex-post feasible: non-Bayesian algorithm guaranteed to satisfy assumptions

## Framework

### Repeated interaction between policymaker and single agent

Hidden state of nature observed after each interaction

Minimum viable case?

## Framework

Repeated interaction between policymaker and single agent

Hidden state of nature observed after each interaction

Minimum viable case?

## Framework

Repeated interaction between policymaker and single agent

Hidden state of nature observed after each interaction

Minimum viable case?

## Contribution

### Standard behavioral assumptions are insufficient, allow for odd behaviors

Refine existing assumptions to counterfactual calibration

▶ Agent fully and consistently exploits any private information

Propose calibrated policy that adapts over time using historical data

Conditions under policymaker's regret is bounded relative to best static policy

## Contribution

Standard behavioral assumptions are insufficient, allow for odd behaviors

Refine existing assumptions to counterfactual calibration

▶ Agent fully and consistently exploits any private information

Propose calibrated policy that adapts over time using historical data

Conditions under policymaker's regret is bounded relative to best static policy

## Contribution

Standard behavioral assumptions are insufficient, allow for odd behaviors

Refine existing assumptions to counterfactual calibration
- ▶ Agent fully and consistently exploits any private information

Propose calibrated policy that adapts over time using historical data

Conditions under policymaker's regret is bounded relative to best static policy

## Contribution

Standard behavioral assumptions are insufficient, allow for odd behaviors

Refine existing assumptions to counterfactual calibration
  ▶ Agent fully and consistently exploits any private information

Propose calibrated policy that adapts over time using historical data

Conditions under policymaker's regret is bounded relative to best static policy

## Related Literature

**Robust dynamic mechanisms**, e.g. Chassang 2013, Penta 2015, Madarász and Prat 2016, Mirrokni, Paes Leme, Tang, and Zuo 2020, Carroll 2020, ...

**Data-driven auction design**, e.g. Blum and Hartline 2005, Elkind 2007, Balcan, Blum, Hartline, and Mansour 2008, Cole and Roughgarden 2014, ...

**Regret-based behavioral assumptions**, e.g. Foster and Vohra 1997, Nekipelov, Syrgkanis, and Tardos 2015, Braverman, Mao, Schneider, and Weinberg 2018, ...

**Statistical learning in incomplete information games**, e.g. Liang 2020, Immorlica, Mao, Slivkins, and Wu 2020, Cummings, Devanur, Huang, and Wang 2020, ...

## Related Literature

**Robust dynamic mechanisms**, e.g. Chassang 2013, Penta 2015, Madarász and Prat 2016, Mirrokni, Paes Leme, Tang, and Zuo 2020, Carroll 2020, ...

**Data-driven auction design**, e.g. Blum and Hartline 2005, Elkind 2007, Balcan, Blum, Hartline, and Mansour 2008, Cole and Roughgarden 2014, ...

Regret-based behavioral assumptions, e.g. Foster and Vohra 1997, Nekipelov, Syrgkanis, and Tardos 2015, Braverman, Mao, Schneider, and Weinberg 2018, ...

Statistical learning in incomplete information games, e.g. Liang 2020, Immorlica, Mao, Slivkins, and Wu 2020, Cummings, Devanur, Huang, and Wang 2020, ...

## Related Literature

**Robust dynamic mechanisms**, e.g. Chassang 2013, Penta 2015, Madarász and Prat 2016, Mirrokni, Paes Leme, Tang, and Zuo 2020, Carroll 2020, ...

**Data-driven auction design**, e.g. Blum and Hartline 2005, Elkind 2007, Balcan, Blum, Hartline, and Mansour 2008, Cole and Roughgarden 2014, ...

**Regret-based behavioral assumptions**, e.g. Foster and Vohra 1997, Nekipelov, Syrgkanis, and Tardos 2015, Braverman, Mao, Schneider, and Weinberg 2018, ...

**Statistical learning in incomplete information games**, e.g. Liang 2020, Immorlica, Mao, Slivkins, and Wu 2020, Cummings, Devanur, Huang, and Wang 2020, ...

## Related Literature

**Robust dynamic mechanisms**, e.g. Chassang 2013, Penta 2015, Madarász and Prat 2016, Mirrokni, Paes Leme, Tang, and Zuo 2020, Carroll 2020, ...

**Data-driven auction design**, e.g. Blum and Hartline 2005, Elkind 2007, Balcan, Blum, Hartline, and Mansour 2008, Cole and Roughgarden 2014, ...

**Regret-based behavioral assumptions**, e.g. Foster and Vohra 1997, Nekipelov, Syrgkanis, and Tardos 2015, Braverman, Mao, Schneider, and Weinberg 2018, ...

**Statistical learning in incomplete information games**, e.g. Liang 2020, Immorlica, Mao, Slivkins, and Wu 2020, Cummings, Devanur, Huang, and Wang 2020, ...

Introduction
00000

Model
00000

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

Model

Introduction
00000

Model
●0000

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

## Stage Game

Policymaker sets policy $p \in \mathcal{P}$, sends message $m \in \mathcal{M}$

Agent chooses response $r \in \mathcal{R}$

Hidden state of nature $s \in \mathcal{S}$

Payoffs $u^P(p, r, s)$ and $u^A(p, r, s)$

**Assumption**: $\mathcal{P}, \mathcal{M}, \mathcal{R}, \mathcal{S}$ finite

Introduction
ooooo

Model
●oooo

Agent's Behavior
ooooooooooo

Calibrated Policy
ooooooooo

Conclusion
ooo

# Stage Game

Policymaker sets policy $p \in \mathcal{P}$, sends message $m \in \mathcal{M}$

Agent chooses response $r \in \mathcal{R}$

Hidden state of nature $s \in \mathcal{S}$

Payoffs $u^P(p, r, s)$ and $u^A(p, r, s)$

**Assumption**: $\mathcal{P}, \mathcal{M}, \mathcal{R}, \mathcal{S}$ finite

## Stage Game

Policymaker sets policy $p \in \mathcal{P}$, sends message $m \in \mathcal{M}$

Agent chooses response $r \in \mathcal{R}$

Hidden state of nature $s \in \mathcal{S}$

Payoffs $u^P(p, r, s)$ and $u^A(p, r, s)$

**Assumption**: $\mathcal{P}, \mathcal{M}, \mathcal{R}, \mathcal{S}$ finite

## Stage Game

Policymaker sets policy $p \in \mathcal{P}$, sends message $m \in \mathcal{M}$

Agent chooses response $r \in \mathcal{R}$

Hidden state of nature $s \in \mathcal{S}$

Payoffs $u^P(p, r, s)$ and $u^A(p, r, s)$

**Assumption**: $\mathcal{P}, \mathcal{M}, \mathcal{R}, \mathcal{S}$ finite

## Stage Game

Policymaker sets policy $p \in \mathcal{P}$, sends message $m \in \mathcal{M}$

Agent chooses response $r \in \mathcal{R}$

Hidden state of nature $s \in \mathcal{S}$

Payoffs $u^P(p, r, s)$ and $u^A(p, r, s)$

**Assumption**: $\mathcal{P}, \mathcal{M}, \mathcal{R}, \mathcal{S}$ finite

## Repeated Game

Stage game repeated $T$ times

**Period t − 1**  →  **Period t**                                         →  **Period t + 1**

⋮
1. Policymaker sets policy $p_t$ and message $m_t$                    ⋮
2. Agent chooses response $r_t$
3. State $s_t$ is observed

Policymaker's mechanism : history $\rightarrow p_t, m_t$

Agent's strategy : history, $p_t, m_t \rightarrow r_t$

Introduction
00000

Model
0●000

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

## Repeated Game

Stage game repeated $T$ times

**Period t − 1**  →  **Period t**                                  →   **Period t + 1**

     ⋮          1. Policymaker sets policy $p_t$ and message $m_t$        ⋮

                 2. Agent chooses response $r_t$

                 3. State $s_t$ is observed

Policymaker's mechanism : $\text{history} \to p_t, m_t$

Agent's strategy : $\text{history}, p_t, m_t \to r_t$

## Repeated Game

Stage game repeated $T$ times

**Period t − 1** → **Period t** → **Period t + 1**

⋮

1. Policymaker sets policy $p_t$ and message $m_t$
2. Agent chooses response $r_t$
3. State $s_t$ is observed

⋮

Policymaker's mechanism : $\text{history} \to p_t, m_t$

Agent's strategy : $\text{history}, p_t, m_t \to r_t$

## Policymaker's Regret

$r_{1:T}$ = agent's responses to actual policies $p_{1:T}$

$r_{1:T}^p$ = agent's responses to fixed policy $p$

**Definition**: policymaker's regret compares actual policies with best fixed policy, i.e.

$$\underbrace{\max_p \frac{1}{T} \sum_{t=1}^{T} u^P(p, r_t^p, s_t)}_{\text{utility under best fixed policy } p} - \underbrace{\frac{1}{T} \sum_{t=1}^{T} u^P(p_t, r_t, s_t)}_{\text{utility under actual policies } p_{1:T}}$$

**Note**: standard no-regret guarantees do not apply since $r_t^p \neq r_t$.

## Policymaker's Regret

$r_{1:T}$ = agent's responses to actual policies $p_{1:T}$

$r^p_{1:T}$ = agent's responses to fixed policy $p$

**Definition**: policymaker's regret compares actual policies with best fixed policy, i.e.

$$\underbrace{\max_p \frac{1}{T} \sum_{t=1}^{T} u^P(p, r^p_t, s_t)}_{\text{utility under best fixed policy } p} - \underbrace{\frac{1}{T} \sum_{t=1}^{T} u^P(p_t, r_t, s_t)}_{\text{utility under actual policies } p_{1:T}}$$

**Note**: standard no-regret guarantees do not apply since $r^p_t \neq r_t$.

## Policymaker's Regret

$r_{1:T}$ = agent's responses to actual policies $p_{1:T}$

$r_{1:T}^p$ = agent's responses to fixed policy $p$

**Definition**: policymaker's regret compares actual policies with best fixed policy, i.e.

$$\underbrace{\max_p \frac{1}{T} \sum_{t=1}^{T} u^P(p, r_t^p, s_t)}_{\text{utility under best fixed policy } p} - \underbrace{\frac{1}{T} \sum_{t=1}^{T} u^P(p_t, r_t, s_t)}_{\text{utility under actual policies } p_{1:T}}$$

**Note**: standard no-regret guarantees do not apply since $r_t^p \neq r_t$.

## Policymaker's Regret

$r_{1:T}$ = agent's responses to actual policies $p_{1:T}$

$r_{1:T}^{p}$ = agent's responses to fixed policy $p$

**Definition**: policymaker's regret compares actual policies with best fixed policy, i.e.

$$\underbrace{\max_{p} \frac{1}{T} \sum_{t=1}^{T} u^P(p, r_t^p, s_t)}_{\text{utility under best fixed policy } p} - \underbrace{\frac{1}{T} \sum_{t=1}^{T} u^P(p_t, r_t, s_t)}_{\text{utility under actual policies } p_{1:T}}$$

**Note**: standard no-regret guarantees do not apply since $r_t^p \neq r_t$.

Introduction
00000

Model
000●0

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

## Price Regulation

*Running Example*

Firm produces indivisible good at cost for price to maximize profit.

- Response $r$ : cost $\rightarrow$ price
- State $s =$ firm's cost, buyer's value

Policymaker regulates price to maximize welfare.

- Policy $p =$ (price floor, price cap)

Outcome is sale $= \mathbf{1}(\text{value} \geq \text{price}) \cdot \mathbf{1}(\text{price floor} \leq \text{price} \leq \text{price cap})$

- Profit $= \text{sale} \cdot (\text{price} - \text{cost})$
- Welfare $= \text{sale} \cdot (\text{value} - \text{cost})$

## Price Regulation                                    *Running Example*

Firm produces indivisible good at cost for price to maximize profit.

- ► Response $r$ : cost → price
- ► State $s =$ firm's cost, buyer's value

Policymaker regulates price to maximize welfare.

- ► Policy $p =$ (price floor, price cap)

Outcome is sale $= \mathbf{1}(\text{value} \geq \text{price}) \cdot \mathbf{1}(\text{price floor} \leq \text{price} \leq \text{price cap})$

- ► Profit $=$ sale $\cdot$ (price − cost)
- ► Welfare $=$ sale $\cdot$ (value − cost)

## Price Regulation *Running Example*

Firm produces indivisible good at $\text{cost}$ for $\text{price}$ to maximize profit.

- Response $r : \text{cost} \rightarrow \text{price}$
- State $s = $ firm's $\text{cost}$, buyer's $\text{value}$

Policymaker regulates $\text{price}$ to maximize welfare.

- Policy $p = (\text{price floor}, \text{price cap})$

Outcome is $\text{sale} = \mathbf{1}(\text{value} \geq \text{price}) \cdot \mathbf{1}(\text{price floor} \leq \text{price} \leq \text{price cap})$

- Profit $= \text{sale} \cdot (\text{price} - \text{cost})$
- Welfare $= \text{sale} \cdot (\text{value} - \text{cost})$

Introduction
00000

Model
0000●

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

# Repeated Price Regulation

*Running Example*

### Sequence of buyers $t$ with $\mathrm{value}_t$

State $s_t = (\mathrm{value}_t, \mathrm{cost}_t)$ observed after period $t$

▶ Replace $\mathrm{value}_t$ with $\mathrm{sale}_t$ if needed

$\mathrm{price}_t, \mathrm{price\ floor}_t, \mathrm{price\ cap}_t$ depend on observed history

Introduction
OOOOO

Model
OOOO●

Agent's Behavior
OOOOOOOOOOO

Calibrated Policy
OOOOOOOO

Conclusion
OOO

# Repeated Price Regulation

*Running Example*

Sequence of buyers $t$ with $\mathrm{value}_t$

State $s_t = (\mathrm{value}_t, \mathrm{cost}_t)$ observed after period $t$
- Replace $\mathrm{value}_t$ with $\mathrm{sale}_t$ if needed

$\mathrm{price}_t, \mathrm{price\ floor}_t, \mathrm{price\ cap}_t$ depend on observed history

Introduction
00000

Model
0000●

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

## Repeated Price Regulation

*Running Example*

Sequence of buyers $t$ with $\text{value}_t$

State $s_t = (\text{value}_t, \text{cost}_t)$ observed after period $t$
- Replace $\text{value}_t$ with $\text{sale}_t$ if needed

$\text{price}_t, \text{price floor}_t, \text{price cap}_t$ depend on observed history

Introduction
00000

Model
00000

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

Agent's Behavior

## Preview of Definitions

## Dealing with Information

**Tortoise travels 1km in 1h** : uninformed agent satisfies no-regret.



**Hare travels 1km in 1h** : informed agent satisfies no-regret.

## Dealing with Information

**Tortoise travels 1km in 1h** : uninformed agent satisfies no-regret.



START ▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪ FINISH

**Hare travels 1km in 1h** : informed agent satisfies no-regret.



START FINISH

Introduction
○○○○○

Model
○○○○○

Agent's Behavior
○●○○○○○○○○○

Calibrated Policy
○○○○○○○○

Conclusion
○○○

## Dealing with Information

**Tortoise travels 1km in 1h** : uninformed agent satisfies no-regret.



$\boxed{\text{START}}$ ▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪ $\boxed{\text{FINISH}}$

**Hare travels 1km in 2min** : informed agent satisfies no-regret conditioned on
her information.



$\boxed{\text{START}}$ ▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪ $\boxed{\text{FINISH}}$

## Agent's Conditional Regret

**Definition**: agent's regret conditioned on information $I_t$ is

$$
\max_{h:\mathcal{I}\to\mathcal{R}} \quad \underbrace{\frac{1}{T}\sum_{t=1}^{T} u^A(p_t, h(I_t), s_t)}_{\text{utility under modified responses } h(I_t)} \quad - \quad \underbrace{\frac{1}{T}\sum_{t=1}^{T} u^A(p_t, r_t, s_t)}_{\text{utility under actual responses } r_{1:T}}
$$

where modification rule $h$ maps information $I_t$ to response $h(I_t)$.

## No-Regret

**Definition**: regret = regret conditioned on public information

$$I_t = (p_t, m_t)$$

**Definition**: no-regret = regret $\to 0$ as $T \to \infty$

## No-Regret

**Definition**: regret = regret conditioned on public information

$$I_t = (p_t, m_t)$$

**Definition**: no-regret = regret $\to 0$ as $T \to \infty$

## Behavior Reveals Information

Game of **rock-paper-scissors** between agent and nature.

$s_t =$ $\boxed{R}$ $\boxed{P}$ $\boxed{S}$ $\boxed{R}$ $\boxed{P}$ $\boxed{S}$ $\boxed{R}$ $\boxed{P}$ $\boxed{S}$ $\boxed{R}$ $\boxed{P}$ $\boxed{S}$

An uninformed strategy; no correlation between response and state

$r_t =$ $\boxed{S}$ $\boxed{R}$ $\boxed{R}$ $\boxed{P}$ $\boxed{P}$ $\boxed{P}$ $\boxed{S}$ $\boxed{P}$ $\boxed{R}$ $\boxed{R}$ $\boxed{P}$ $\boxed{R}$

$u_t^A =$ $\boxed{-1}$ $\boxed{-1}$ $\boxed{1}$ $\boxed{1}$ $\boxed{0}$ $\boxed{-1}$ $\boxed{-1}$ $\boxed{0}$ $\boxed{1}$ $\boxed{0}$ $\boxed{0}$ $\boxed{1}$

## Behavior Reveals Information

Game of **rock-paper-scissors** between agent and nature.

$$s_t = \quad \boxed{R} \quad \boxed{P} \quad \boxed{S} \quad \boxed{R} \quad \boxed{P} \quad \boxed{S} \quad \boxed{R} \quad \boxed{P} \quad \boxed{S} \quad \boxed{R} \quad \boxed{P} \quad \boxed{S}$$

An uninformed strategy; no correlation between response and state

$$r_t = \quad \boxed{S} \quad \boxed{R} \quad \boxed{R} \quad \boxed{P} \quad \boxed{P} \quad \boxed{P} \quad \boxed{S} \quad \boxed{P} \quad \boxed{R} \quad \boxed{R} \quad \boxed{P} \quad \boxed{R}$$

$$u_t^A = \quad \boxed{-1} \quad \boxed{-1} \quad \boxed{1} \quad \boxed{1} \quad \boxed{0} \quad \boxed{-1} \quad \boxed{-1} \quad \boxed{0} \quad \boxed{1} \quad \boxed{0} \quad \boxed{0} \quad \boxed{1}$$

## Behavior Reveals Information

Game of **rock-paper-scissors** between agent and nature.

$s_t =$ | R | | P | | S | | R | | P | | S | | R | | P | | S | | R | | P | | S |

An informed strategy; perfect correlation between response and state

$r_t =$ | R | | P | | S | | R | | P | | S | | R | | P | | S | | R | | P | | S |
$u_t^A =$ | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 |

## Calibration

**Definition**: internal regret = regret conditioned on information revealed on-path, i.e.

$$I_t = (p_t, m_t, r_t)$$

**Definition**: calibration = internal regret $\rightarrow 0$ as $T \rightarrow \infty$

## Calibration

**Definition**: internal regret $=$ regret conditioned on information revealed on-path, i.e.

$$I_t = (p_t, m_t, r_t)$$

**Definition**: calibration $=$ internal regret $\rightarrow 0$ as $T \rightarrow \infty$

## No-Regret to Calibration

**Example:** no-regret allows strange behavior that calibration rules out

$s_t =$   R   R   R   R   R   R   P   P   P   P   P   P

Tortoise strategy: uninformed, optimal, no-regret

$r_t =$   P   P   P   P   P   P   P   P   P   P   P   P

$u_t^A =$   1   1   1   1   1   1   0   0   0   0   0   0

# No-Regret to Calibration

**Example:** no-regret allows strange behavior that calibration rules out

$s_t =$   R   R   R   R   R   R   P   P   P   P   P   P

Tortoise strategy: uninformed, optimal, no-regret

$r_t =$   P   P   P   P   P   P   P   P   P   P   P   P
$u_t^A =$   1   1   1   1   1   1   0   0   0   0   0   0

## No-Regret to Calibration

**Example:** no-regret allows strange behavior that calibration rules out

$s_t =$    R   R   R   R   R   R   P   P   P   P   P   P

Lazy hare strategy: informed, suboptimal, no-regret

$r_t =$    R   R   R   R   R   R   S   S   S   S   S   S

$u_t^A =$   0   0   0   0   0   0   1   1   1   1   1   1

## Impossibility Result

### Proposition

*There exists a strategy for the agent where no mechanism can guarantee non-trivial bound on policymaker's regret across all $s_{1:T}$ where agent's strategy is calibrated.*

**Takeaway**: calibration is not enough for low-regret policy design

▶ Even if we know the agent's strategy in advance

## Impossibility Result

### Proposition

*There exists a strategy for the agent where no mechanism can guarantee non-trivial bound on policymaker's regret across all $s_{1:T}$ where agent's strategy is calibrated.*

**Takeaway**: calibration is not enough for low-regret policy design

▶ Even if we know the agent's strategy in advance

## Counterfactual Calibration

**Definition**: counterfactual internal regret = regret conditioned on information revealed on- and off-path, i.e.

$$I_t = \left(p_t, m_t, r_t, r_t^{p_1}, \ldots, r_t^{p_n}\right)$$

**Definition**: counterfactual calibration = CIR $\rightarrow 0$ as $T \rightarrow \infty$

## Counterfactual Calibration

**Definition**: counterfactual internal regret = regret conditioned on information
revealed on- and off-path, i.e.

$$I_t = \left(p_t, m_t, r_t, r_t^{p_1}, \ldots, r_t^{p_n}\right)$$

**Definition**: counterfactual calibration = CIR $\to 0$ as $T \to \infty$

## Calibration to Counterfactual Calibration

**Example:** calibration allows strange behavior that counterfactual calibration rules out

$s_t =$  | R | R | R | R | R | R | P | P | P | P | P | P |

If policymaker follows mechanism, play tortoise strategy

$r_t =$  | P | P | P | P | P | P | P | P | P | P | P | P |

$u_t^A =$  | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

If policymaker follows fixed policy, play active hare strategy

$r_t =$  | P | P | P | P | P | P | S | S | S | S | S | S |

$u_t^A =$  | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

## Calibration to Counterfactual Calibration

**Example:** calibration allows strange behavior that counterfactual calibration rules out

$s_t =$   R   R   R   R   R   R   P   P   P   P   P   P

If policymaker follows mechanism, play tortoise strategy

$r_t =$   P   P   P   P   P   P   P   P   P   P   P   P

$u_t^A =$   1   1   1   1   1   1   0   0   0   0   0   0

If policymaker follows fixed policy, play active hare strategy

$r_t =$   P   P   P   P   P   P   S   S   S   S   S   S

$u_t^A =$   1   1   1   1   1   1   1   1   1   1   1   1

## Calibration to Counterfactual Calibration

**Example:** calibration allows strange behavior that counterfactual calibration rules out

$s_t =$   | R | R | R | R | R | R | P | P | P | P | P | P |

If policymaker follows mechanism, play tortoise strategy

$r_t =$   | P | P | P | P | P | P | P | P | P | P | P | P |
$u_t^A =$   | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

If policymaker follows fixed policy, play active hare strategy

$r_t =$   | P | P | P | P | P | P | S | S | S | S | S | S |
$u_t^A =$   | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

20 / 32

# Non-clairvoyance

**Definition**: non-clairvoyance = counterfactual calibration & non-negative regret

**Intuition:**

1. Counterfactual calibration $\implies$ agent fully exploits her private info

2. Non-negative regret $\implies$ agent doesn't outperform best use of public info

3. Therefore, her private info must not be useful

# Non-clairvoyance

**Definition**: non-clairvoyance = counterfactual calibration & non-negative regret

**Intuition:**

1. Counterfactual calibration $\implies$ agent fully exploits her private info
2. Non-negative regret $\implies$ agent doesn't outperform best use of public info
3. Therefore, her private info must not be useful

Introduction
00000

Model
00000

Agent's Behavior
0000000000●

Calibrated Policy
00000000

Conclusion
000

# Non-clairvoyance

**Definition**: non-clairvoyance = counterfactual calibration & non-negative regret

**Intuition:**

1. Counterfactual calibration $\implies$ agent fully exploits her private info
2. Non-negative regret $\implies$ agent doesn't outperform best use of public info
3. Therefore, her private info must not be useful

# Non-clairvoyance

**Definition**: non-clairvoyance = counterfactual calibration & non-negative regret

**Intuition:**

1. Counterfactual calibration $\implies$ agent fully exploits her private info
2. Non-negative regret $\implies$ agent doesn't outperform best use of public info
3. Therefore, her private info must not be useful

Introduction
00000

Model
00000

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

Calibrated Policy

## Preview of Mechanism

In each period $t$...

1. Form probabilistic forecast of state $s_t$

2. Assume agent shares forecast

3. Choose $\epsilon$-robust policy based on forecast

**Informal result**: works well when agent is non-clairvoyant

## Preview of Mechanism

In each period $t$...

1. Form probabilistic forecast of state $s_t$

2. Assume agent shares forecast

3. Choose $\epsilon$-robust policy based on forecast

**Informal result**: works well when agent is non-clairvoyant

## Preview of Mechanism

In each period $t$...

1. Form probabilistic forecast of state $s_t$
2. Assume agent shares forecast
3. Choose $\epsilon$-robust policy based on forecast

**Informal result**: works well when agent is non-clairvoyant

## Preview of Mechanism

In each period $t$...

1. Form probabilistic forecast of state $s_t$
2. Assume agent shares forecast
3. Choose $\epsilon$-robust policy based on forecast

**Informal result**: works well when agent is non-clairvoyant

## $\epsilon$-Robustness $\qquad$ *Stage Game*

$\pi_t^s =$ state distribution

**Definition**: mixed response $\pi^r$ is $\epsilon$-best reply to policy $p$ if

$$\underbrace{\max_{r'} \mathrm{E}_{\pi^s}\left[u^A(p, r', s)\right]}_{\text{optimal utility}} - \underbrace{\mathrm{E}_{\pi^s, \pi^r}\left[u^A(p, r, s)\right]}_{\text{actual utility}} \leq \epsilon$$

**Definition**: policymaker's worst-case utility given $\epsilon$-best reply is

$$\mathrm{WC}_\epsilon(p, \pi^s,) = \min_{\pi^r} \ \mathrm{E}_{\pi^s, \pi^r}\left[u^P(p, r, s)\right] \quad \text{s.t. } \pi^r \text{ is } \epsilon\text{-best reply}$$

## $\epsilon$-Robustness $\qquad$ *Stage Game*

$\pi_t^s =$ state distribution

**Definition**: mixed response $\pi^r$ is $\epsilon$-best reply to policy $p$ if

$$\underbrace{\max_{r'} \mathrm{E}_{\pi^s}\left[u^A(p, r', s)\right]}_{\text{optimal utility}} - \underbrace{\mathrm{E}_{\pi^s, \pi^r}\left[u^A(p, r, s)\right]}_{\text{actual utility}} \leq \epsilon$$

**Definition**: policymaker's worst-case utility given $\epsilon$-best reply is

$$\mathrm{WC}_\epsilon(p, \pi^s, ) = \min_{\pi^r} \ \mathrm{E}_{\pi^s, \pi^r}\left[u^P(p, r, s)\right] \quad \text{s.t. } \pi^r \text{ is } \epsilon\text{-best reply}$$

## $\epsilon$-Robustness                                                      *Stage Game*

$\pi_t^s =$ state distribution

**Definition**: mixed response $\pi^r$ is $\epsilon$-best reply to policy $p$ if

$$\underbrace{\max_{r'} \mathrm{E}_{\pi^s}\left[u^A(p, r', s)\right]}_{\text{optimal utility}} - \underbrace{\mathrm{E}_{\pi^s, \pi^r}\left[u^A(p, r, s)\right]}_{\text{actual utility}} \leq \epsilon$$

**Definition**: policymaker's worst-case utility given $\epsilon$-best reply is

$$\mathrm{WC}_\epsilon(p, \pi^s, ) = \min_{\pi^r} \ \mathrm{E}_{\pi^s, \pi^r}\left[u^P(p, r, s)\right] \quad \text{s.t. } \pi^r \text{ is } \epsilon\text{-best reply}$$

## Cost of $\epsilon$-Robustness

*Stage Game*

**Definition**: policymaker's best-case utility given $\epsilon$-best reply is

$$\mathrm{BC}_\epsilon(p, \pi^s) = \max_{\pi^r} \, \mathrm{E}_{\pi^s, \pi^r}\left[u^P(p, r, s)\right] \quad \text{s.t. } \pi^r \text{ is } \epsilon\text{-best reply}$$

**Definition**: cost of $\epsilon$-robustness is

$$\mathrm{CoR}_\epsilon(p, \pi^s) := \mathrm{BC}_\epsilon(p, \pi^s) - \mathrm{WC}_\epsilon(p, \pi^s)$$

Introduction
00000

Model
00000

Agent's Behavior
00000000000

Calibrated Policy
00●00000

Conclusion
000

## Cost of $\epsilon$-Robustness *Stage Game*

**Definition**: policymaker's best-case utility given $\epsilon$-best reply is

$$\mathrm{BC}_\epsilon(p, \pi^s) = \max_{\pi^r} \mathrm{E}_{\pi^s, \pi^r}\left[u^P(p, r, s)\right] \quad \text{s.t. } \pi^r \text{ is } \epsilon\text{-best reply}$$

**Definition**: cost of $\epsilon$-robustness is

$$\mathrm{CoR}_\epsilon(p, \pi^s) := \mathrm{BC}_\epsilon(p, \pi^s) - \mathrm{WC}_\epsilon(p, \pi^s)$$

## Calibrated Policy

### $\epsilon =$ robustness parameter

$\tilde{\pi}_t^s =$ calibrated forecast with grid width $\delta$

**Policy $p_t$** $:= \epsilon$-robust policy assuming forecast is correct, i.e.

$$p_t \in \arg \max_p \mathrm{WC}_\epsilon(p, \tilde{\pi}_t^s)$$

**Message $m_t$** $:=$ forecast $\tilde{\pi}_t^s$

## Calibrated Policy

$\epsilon =$ robustness parameter

$\tilde{\pi}_t^s =$ calibrated forecast with grid width $\delta$

**Policy $p_t$** $:= \epsilon$-robust policy assuming forecast is correct, i.e.

$$p_t \in \arg \max_p \mathrm{WC}_\epsilon(p, \tilde{\pi}_t^s)$$

**Message $m_t$** $:=$ forecast $\tilde{\pi}_t^s$

## Calibrated Policy

$\epsilon$ = robustness parameter

$\tilde{\pi}_t^s$ = calibrated forecast with grid width $\delta$

**Policy $p_t$** := $\epsilon$-robust policy assuming forecast is correct, i.e.

$$p_t \in \arg \max_p \mathrm{WC}_\epsilon(p, \tilde{\pi}_t^s)$$

**Message $m_t$** := forecast $\tilde{\pi}_t^s$

## Calibrated Policy

$\epsilon =$ robustness parameter

$\tilde{\pi}_t^s =$ calibrated forecast with grid width $\delta$

**Policy $p_t$** $:= \epsilon$-robust policy assuming forecast is correct, i.e.

$$p_t \in \arg \max_p \mathrm{WC}_\epsilon(p, \tilde{\pi}_t^s)$$

**Message $m_t$** $:=$ forecast $\tilde{\pi}_t^s$

## Assumptions

1. Agent is non-clairvoyant
   counterfactual calibration + non-negative regret

2. Information useless to agent under any policy $\Rightarrow$ not harmful to policymaker
   technical assumption on the stage game

## Assumptions

1. Agent is non-clairvoyant
   counterfactual calibration $+$ non-negative regret

2. Information useless to agent under any policy $\Rightarrow$ not harmful to policymaker
   technical assumption on the stage game

Introduction
○○○○○

Model
○○○○○

Agent's Behavior
○○○○○○○○○○○

Calibrated Policy
○○○○○●○○

Conclusion
○○○

# Regret Bound

## Theorem

*Policymaker's regret from calibrated policy is less than*

$$
\underbrace{\frac{1}{T}\sum_{t=1}^{T}\mathrm{CoR}_\epsilon(p_t, \tilde{\pi}_t^s)}_{\substack{\text{cost of robustness}}} + \underbrace{\frac{1}{\epsilon}}_{\substack{\text{sensitivity}}}\left(\underbrace{O\left(\mathrm{CIR}_T\right)}_{\substack{\text{agent miscalibration}}} + \underbrace{\tilde{O}\left(\frac{\sqrt{|\mathcal{S}|\cdot N_\delta(\Delta(\mathcal{S}))}}{T^{1/4}} + \sqrt{\delta}\right)}_{\substack{\text{forecast miscalibration}}}\right)
$$

$\tilde{\pi}_t^s$ = forecast with     $\mathrm{CIR}_T$ = counterfactual     $\mathcal{S}$ = state space
grid width $\delta$          internal regret       $\Delta(\mathcal{S})$ = state distributions
                                             $N_\delta(\cdot)$ = $\delta$-covering number

**Tradeoff**: $\epsilon \uparrow \implies$ sensitivity to miscalibration $\downarrow$ & cost of robustness $\uparrow$

# Regret Bound

## Theorem

*Policymaker's regret from calibrated policy is less than*

$$\underbrace{\frac{1}{T}\sum_{t=1}^{T}\mathrm{CoR}_\epsilon(p_t, \tilde{\pi}_t^s)}_{\substack{\text{cost of robustness}}} + \frac{1}{\epsilon}\left( \underbrace{O\left(\mathrm{CIR}_T\right)}_{\substack{\text{agent miscalibration}}} + \underbrace{\tilde{O}\left(\frac{\sqrt{|\mathcal{S}|\cdot N_\delta(\Delta(\mathcal{S}))}}{T^{1/4}} + \sqrt{\delta}\right)}_{\substack{\text{forecast miscalibration}}} \right)$$

$\tilde{\pi}_t^s = $ forecast with  grid width $\delta$

$\mathrm{CIR}_T = $ counterfactual  internal regret

$\mathcal{S} = $ state space  
$\Delta(\mathcal{S}) = $ state distributions  
$N_\delta(\cdot) = \delta$-covering number

**Tradeoff**: $\epsilon \uparrow \implies$ sensitivity to miscalibration $\downarrow$ & cost of robustness $\uparrow$

## Robustness Lemma

### Lemma

For any distribution $\pi^s$, policy $p$, and constants $\epsilon' > \epsilon > 0$,

$$\mathrm{WC}_{\epsilon'}(p, \pi^s) \geq \mathrm{WC}_{\epsilon}(p, \pi^s) - O\left(\frac{\epsilon' - \epsilon}{\epsilon}\right)$$

# Calibrated Price Regulation                    *Running Example*

Typical tradeoff:

- price $\text{cap}_t$ too large $\implies$ $\text{price}_t$ too large, fewer sales
- price $\text{cap}_t$ too small $\implies$ risk of $\text{cost}_t > \text{price cap}_t$, firm shutdown

How to balance tradeoff depends on market conditions

- Predict market conditions $=$ forecast of $(\text{value}_t, \text{cost}_t)$
- Even if $\text{value}_{1:t-1}$ not observed, feasible if $(\text{sale}_{1:t-1}, \text{cost}_{1:t-1})$ observed

Calibrated policy assumes forecast true and optimizes in stage game.

- Firm's beliefs $\approx$ forecast $\implies$ $\epsilon$-optimal pricing w.r.t. forecast

# Calibrated Price Regulation                    *Running Example*

Typical tradeoff:

- price $\text{cap}_t$ too large $\implies$ $\text{price}_t$ too large, fewer sales
- price $\text{cap}_t$ too small $\implies$ risk of $\text{cost}_t > \text{price cap}_t$, firm shutdown

How to balance tradeoff depends on market conditions

- Predict market conditions = forecast of $(\text{value}_t, \text{cost}_t)$
- Even if $\text{value}_{1:t-1}$ not observed, feasible if $(\text{sale}_{1:t-1}, \text{cost}_{1:t-1})$ observed

Calibrated policy assumes forecast true and optimizes in stage game.

- Firm's beliefs $\approx$ forecast $\implies$ $\epsilon$-optimal pricing w.r.t. forecast

## Calibrated Price Regulation                                    *Running Example*

Typical tradeoff:

- $\text{price cap}_t$ too large $\implies$ $\text{price}_t$ too large, fewer sales
- $\text{price cap}_t$ too small $\implies$ risk of $\text{cost}_t > \text{price cap}_t$, firm shutdown

How to balance tradeoff depends on market conditions

- Predict market conditions = forecast of $(\text{value}_t, \text{cost}_t)$
- Even if $\text{value}_{1:t-1}$ not observed, feasible if $(\text{sale}_{1:t-1}, \text{cost}_{1:t-1})$ observed

Calibrated policy assumes forecast true and optimizes in stage game.

- Firm's beliefs $\approx$ forecast $\implies$ $\epsilon$-optimal pricing w.r.t. forecast

## Calibrated Price Regulation

*Running Example*

Typical tradeoff:

▶ price $\mathrm{cap}_t$ too large $\implies$ price$_t$ too large, fewer sales
▶ price $\mathrm{cap}_t$ too small $\implies$ risk of $\mathrm{cost}_t > $ price $\mathrm{cap}_t$, firm shutdown

How to balance tradeoff depends on market conditions

▶ Predict market conditions = forecast of $(\mathrm{value}_t, \mathrm{cost}_t)$
▶ Even if $\mathrm{value}_{1:t-1}$ not observed, feasible if $(\mathrm{sale}_{1:t-1}, \mathrm{cost}_{1:t-1})$ observed

Calibrated policy assumes forecast true and optimizes in stage game.

▶ Firm's beliefs $\approx$ forecast $\implies$ $\epsilon$-optimal pricing w.r.t. forecast

Introduction
00000

Model
00000

Agent's Behavior
00000000000

Calibrated Policy
00000000

Conclusion
000

# Conclusion

## Ways to Deal with Private Information

1. Assume it doesn't exist

2. Assume it exists but is well-understood

3. Optimize against worst-case private information (extension)

4. Adapt to private information over time (work-in-progress)

## Ways to Deal with Private Information

1. Assume it doesn't exist
2. Assume it exists but is well-understood
3. Optimize against worst-case private information (extension)
4. Adapt to private information over time (work-in-progress)

## Ways to Deal with Private Information

1. Assume it doesn't exist
2. Assume it exists but is well-understood
3. Optimize against worst-case private information (extension)
4. Adapt to private information over time (work-in-progress)

## Ways to Deal with Private Information

1. Assume it doesn't exist
2. Assume it exists but is well-understood
3. Optimize against worst-case private information (extension)
4. Adapt to private information over time (work-in-progress)

## Future Directions

Beyond our minimum viable case...

- ▶ What if feedback is imperfect?
- ▶ Can we incorporate dynamic incentives?
- ▶ Multiple agents?

Potential areas of application?

- ▶ Price regulation that adapts to changing costs and demand
- ▶ Minimum wages that adapt to changing labor markets
- ▶ Worker incentives that adapt to changing workplace

## Future Directions

Beyond our minimum viable case...

- ▶ What if feedback is imperfect?
- ▶ Can we incorporate dynamic incentives?
- ▶ Multiple agents?

Potential areas of application?

- ▶ Price regulation that adapts to changing costs and demand
- ▶ Minimum wages that adapt to changing labor markets
- ▶ Worker incentives that adapt to changing workplace

**Thank you!**