

DATA WITHOUT BORDERS

GAME-THEORETIC CHALLENGES IN DEMOCRATIZING DATA

STATISTICS SEMINAR. APRIL 8, 2024

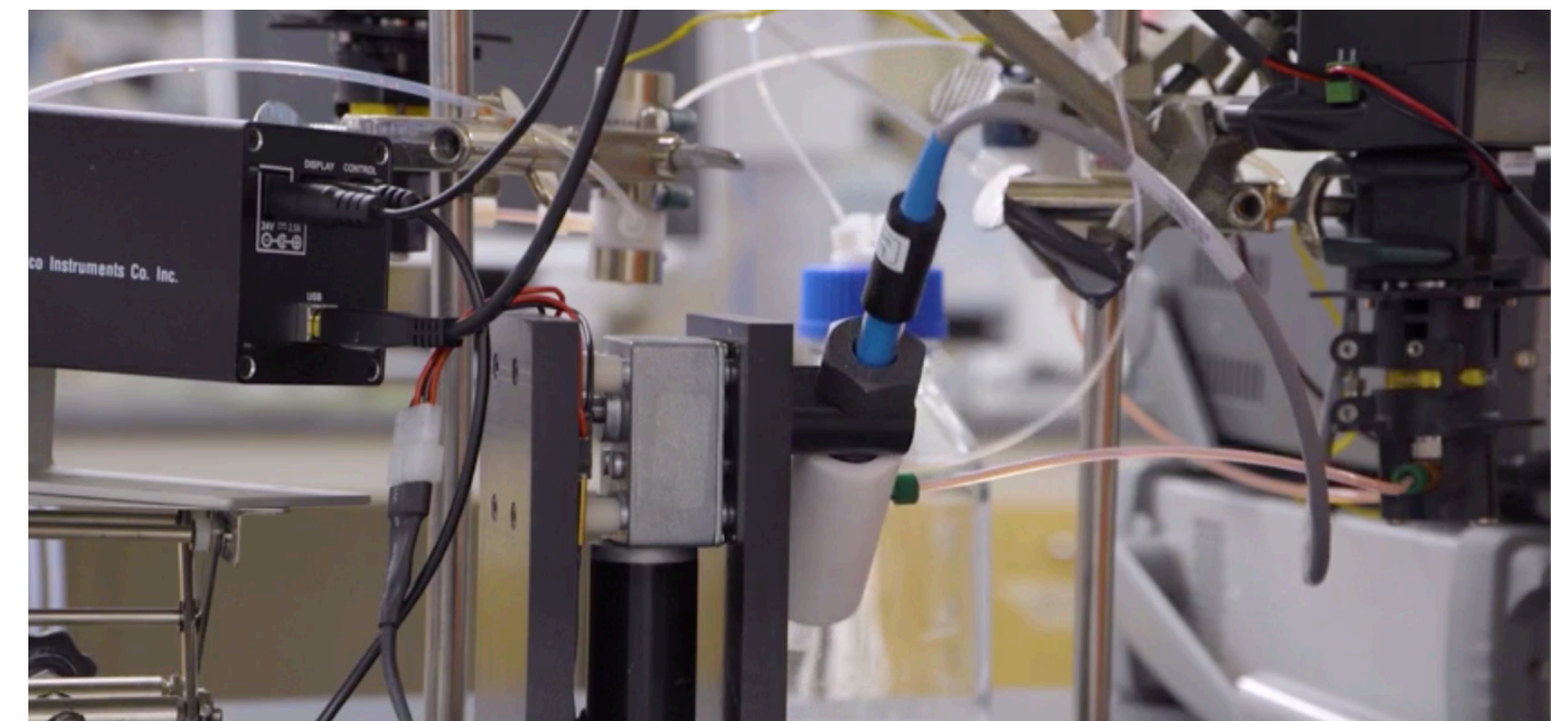
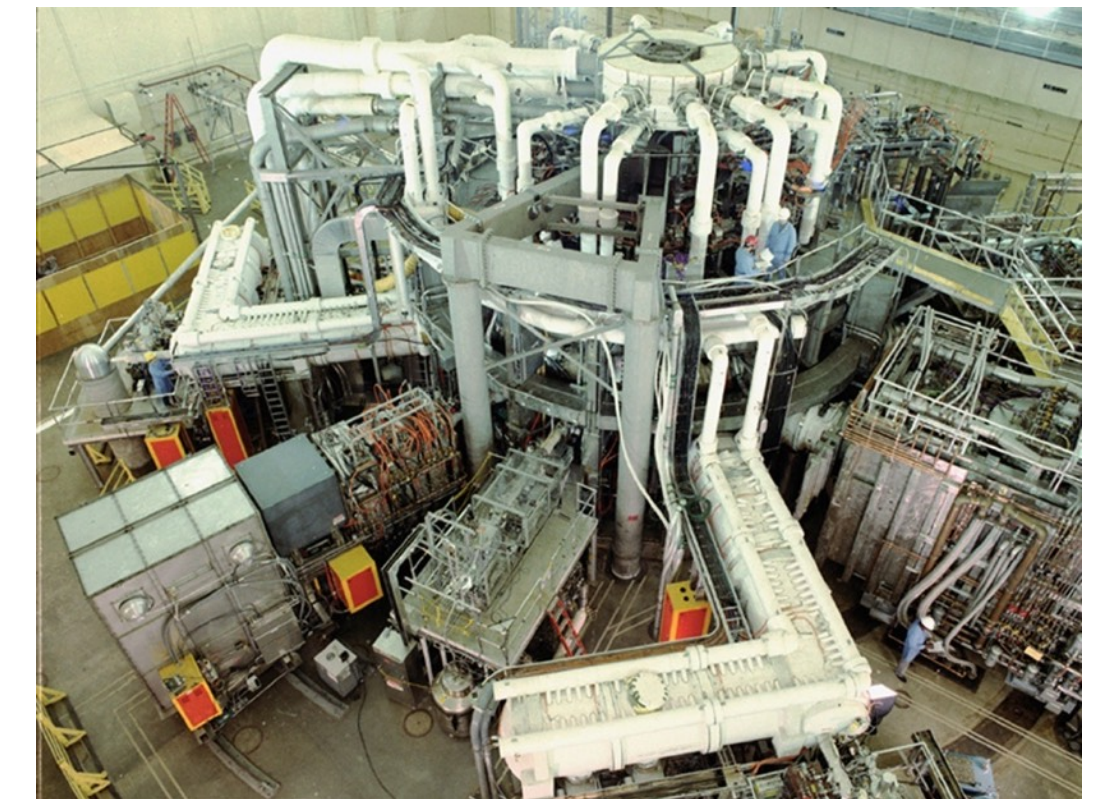
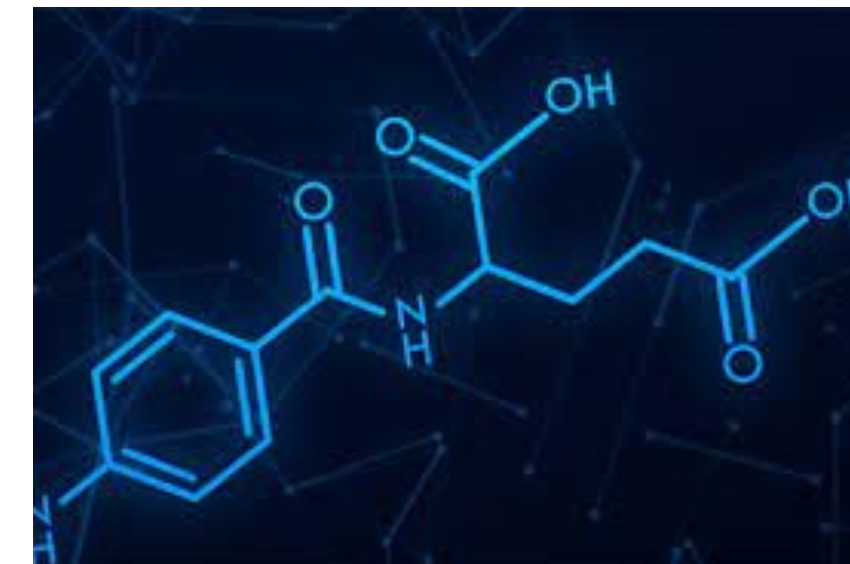
KIRTHEVASAN KANDASAMY

DEPARTMENT OF COMPUTER SCIENCES, UW-MADISON

BASED ON JOINT WORK WITH: YIDING CHEN, JERRY ZHU, AND OTHERS

DATA DRIVEN METHODS ARE UBIQUITOUS

- ▶ Consumer facing businesses
- ▶ Industrial processes
- ▶ Scientific research
- ▶ Transport/logistics



- ▶ Data is the *new oil*.
- ▶ Data is the *new gold*.

*The Economist, NY Times, Forbes, Wired, Deloitte, EY
and several more ...*

- ▶ Data is the *new oil*.
- ▶ Data is the *new gold*.

*The Economist, NY Times, Forbes, Wired, Deloitte, EY
and several more ...*

- ▶ But data is different to other types of resources
 - ▶ Data is **costly** to produce, **free** to replicate.

Everyone collects data, everyone shares their data with others.

- Cost incurred by one organization to produce data can benefit others.
- Better for the organizations, better for society at large.



Small organizations with little data:

A B C D E F

Small organizations with little data:

A B C D E F

Large organization with lots of data:

Small organizations with little data:

A B C D E F

Large organization with lots of data:



Small organizations with little data:

A B C D E F

Large organization with lots of data:



By sharing data with each other, small organizations can compete with larger organizations.

Ethical/Legal

Privacy

Ownership of data

Ethical/Legal

Privacy

Ownership of data

Security

Data breaches

Adversarial attacks

Ethical/Legal

Privacy
Ownership of data

Security

Data breaches
Adversarial attacks

Logistical

Inter-operability
Communication costs

Ethical/Legal

Privacy
Ownership of data

Security

Data breaches
Adversarial attacks

Logistical

Inter-operability
Communication costs

Incentives

Free-riding
Competition
Data monetization
Data valuation

Ethical/Legal

Privacy
Ownership of data

Security

Data breaches
Adversarial attacks

Logistical

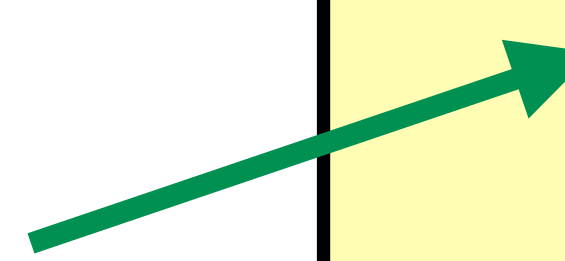
Inter-operability
Communication costs

Incentives

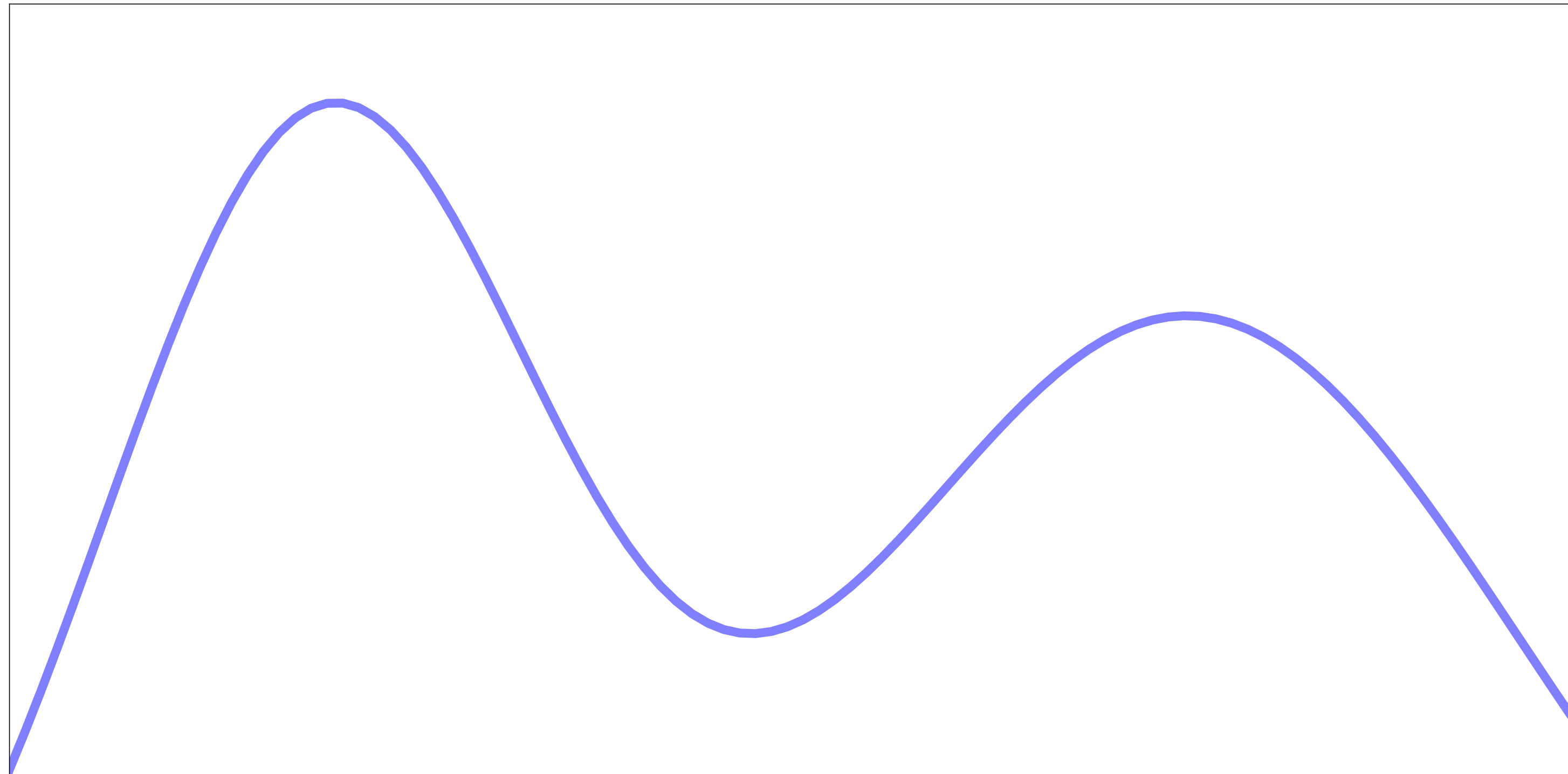
Free-riding
Competition

Data monetization
Data valuation

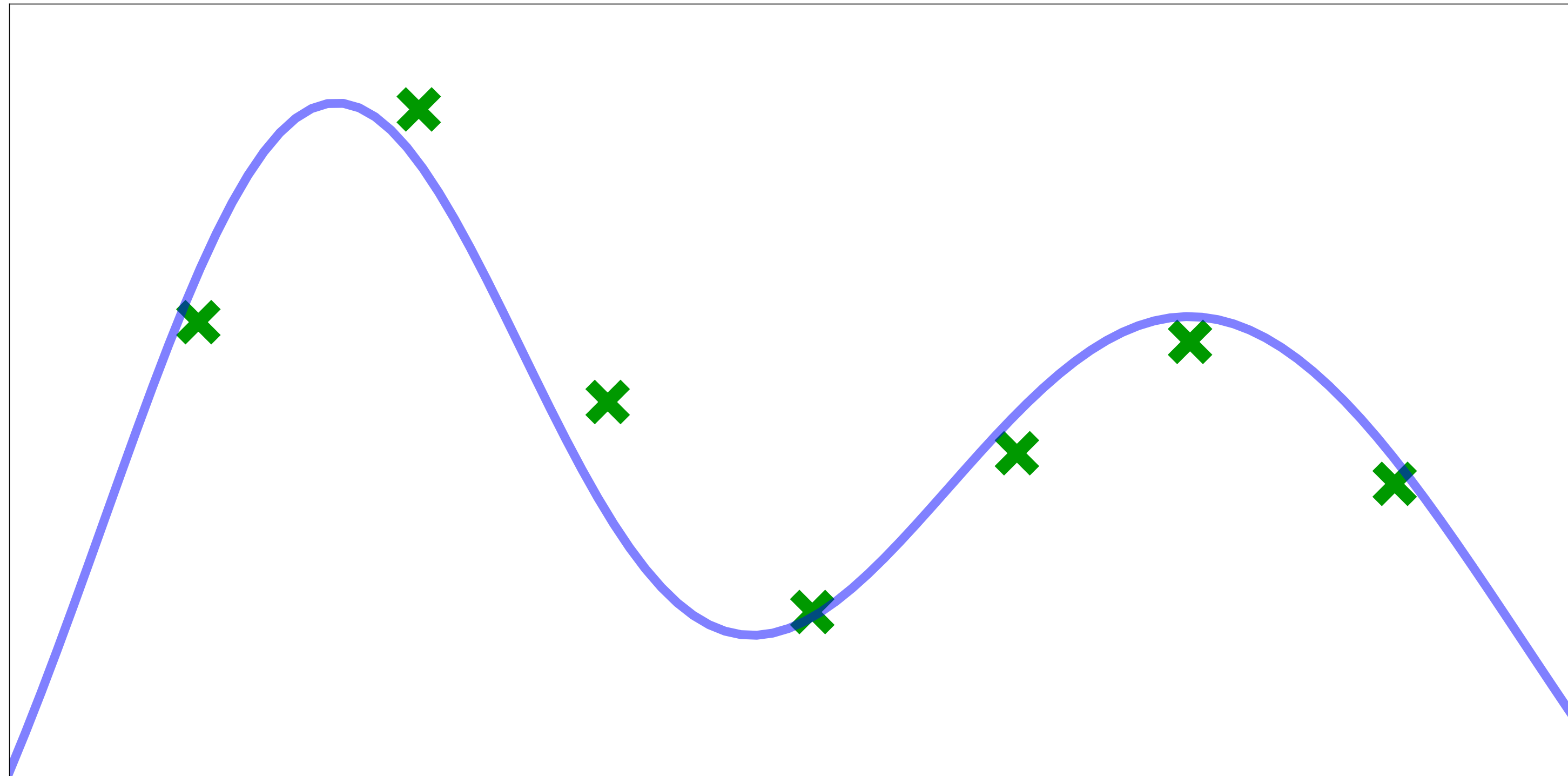
This talk



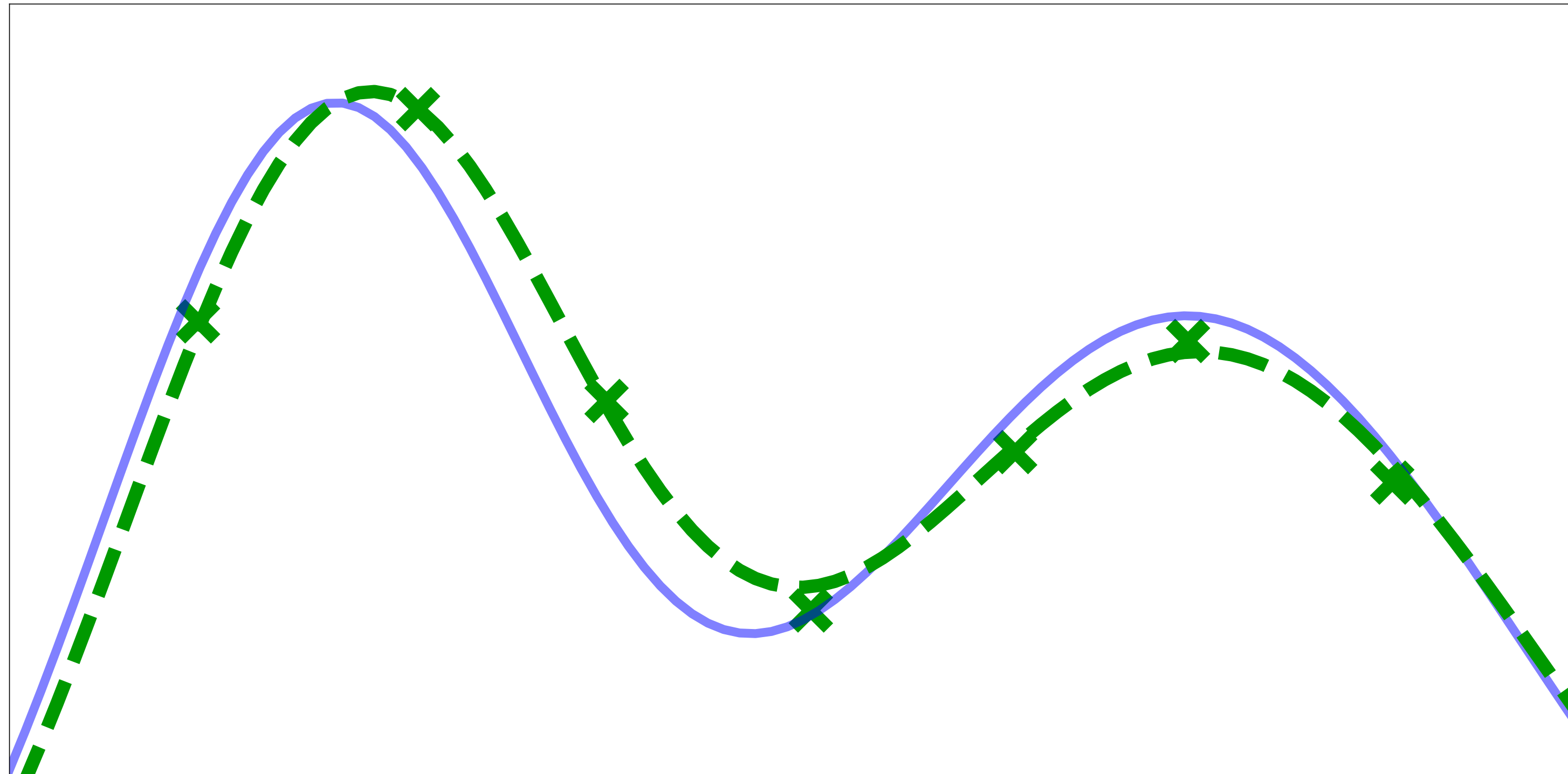
agent's penalty = estimation error + cost of data collection



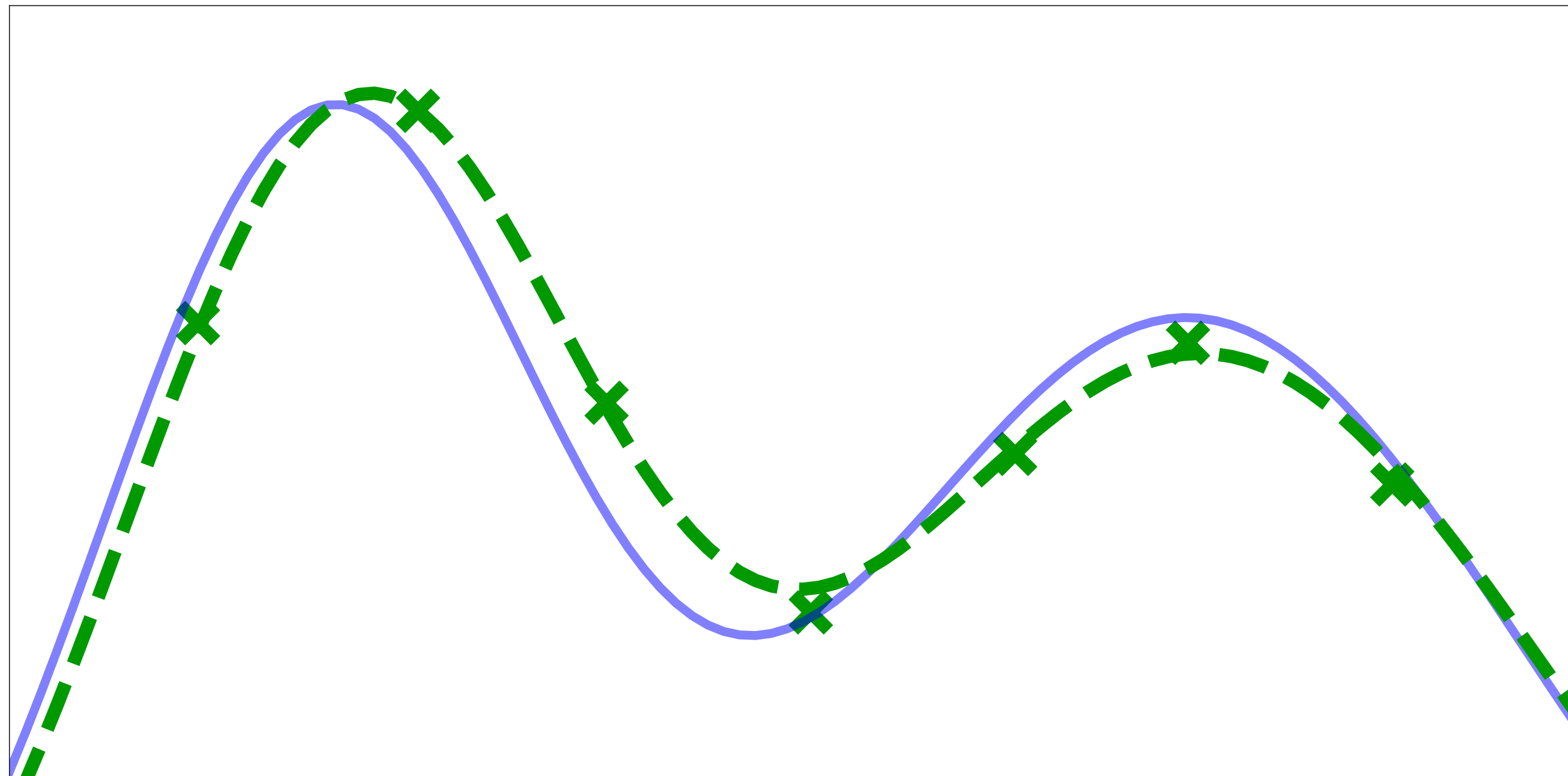
agent's penalty = estimation error + cost of data collection



agent's penalty = estimation error + cost of data collection



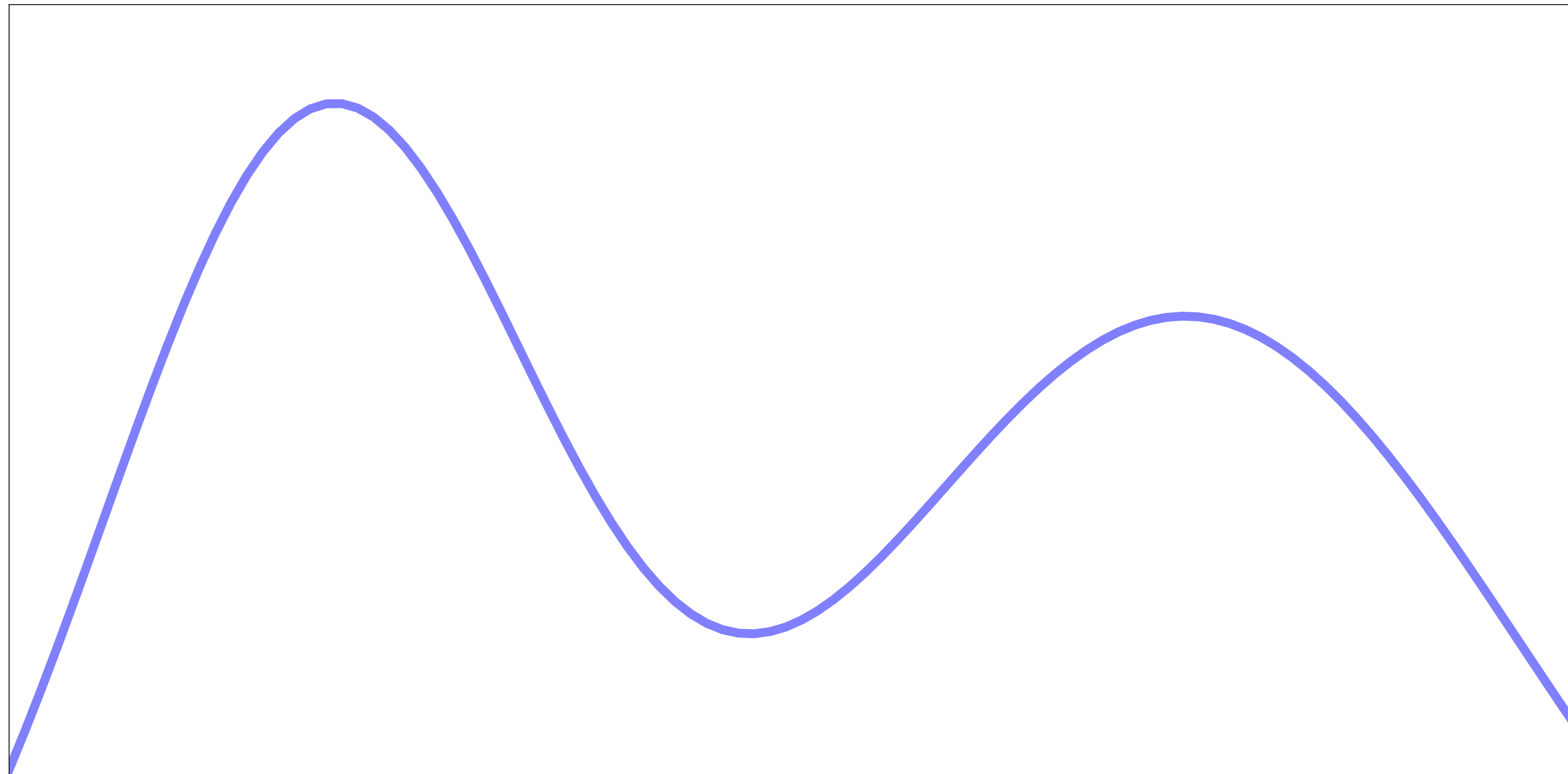
agent's penalty = estimation error + cost of data collection



When working on her own, an agent will collect enough data until the cost offsets the (diminishing) increase in value from data.

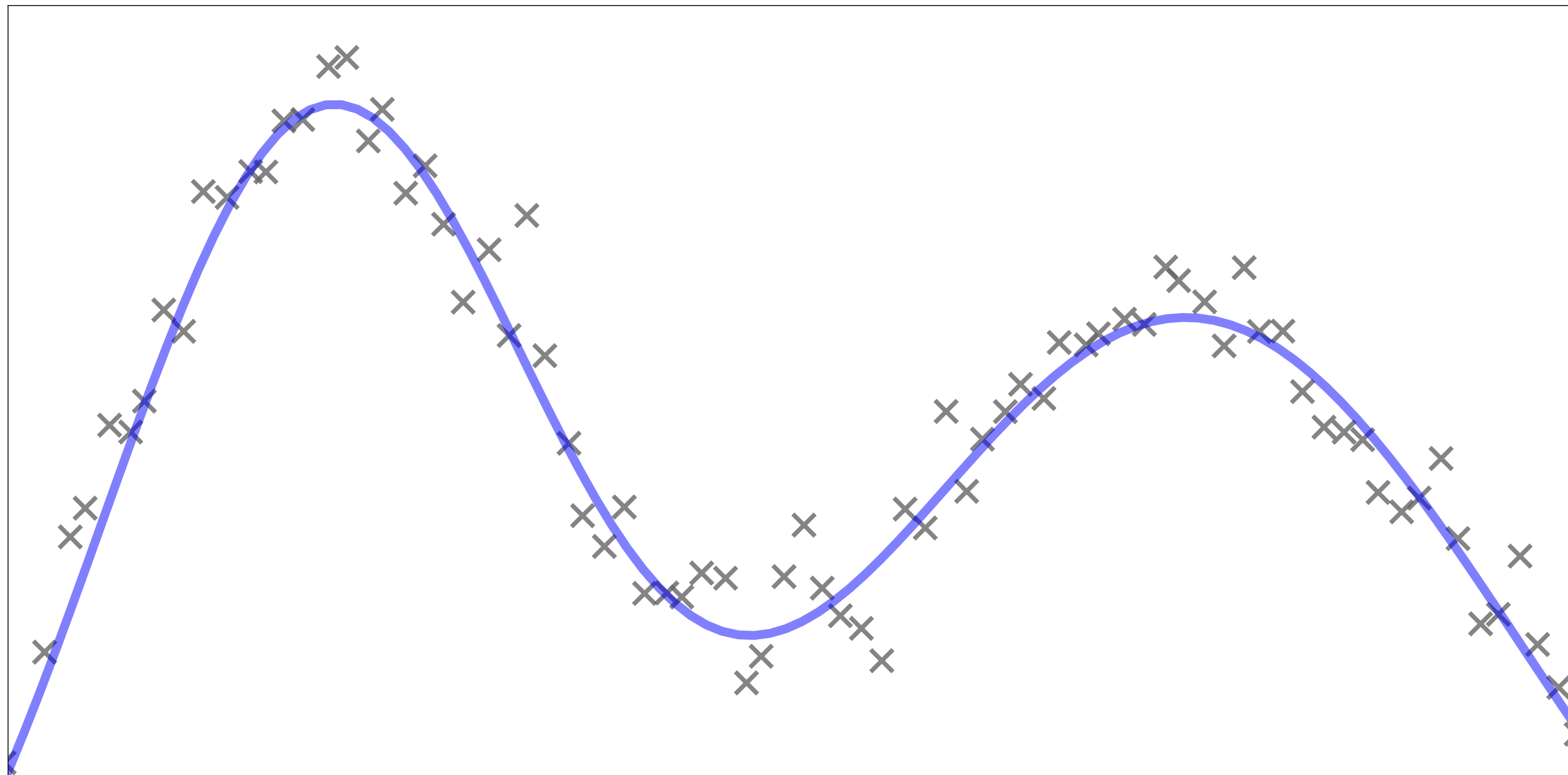
Multiple agents share data via a *naive* pool-and-share protocol:

- ▶ Everyone collects data, everyone gets a copy of the others' data.



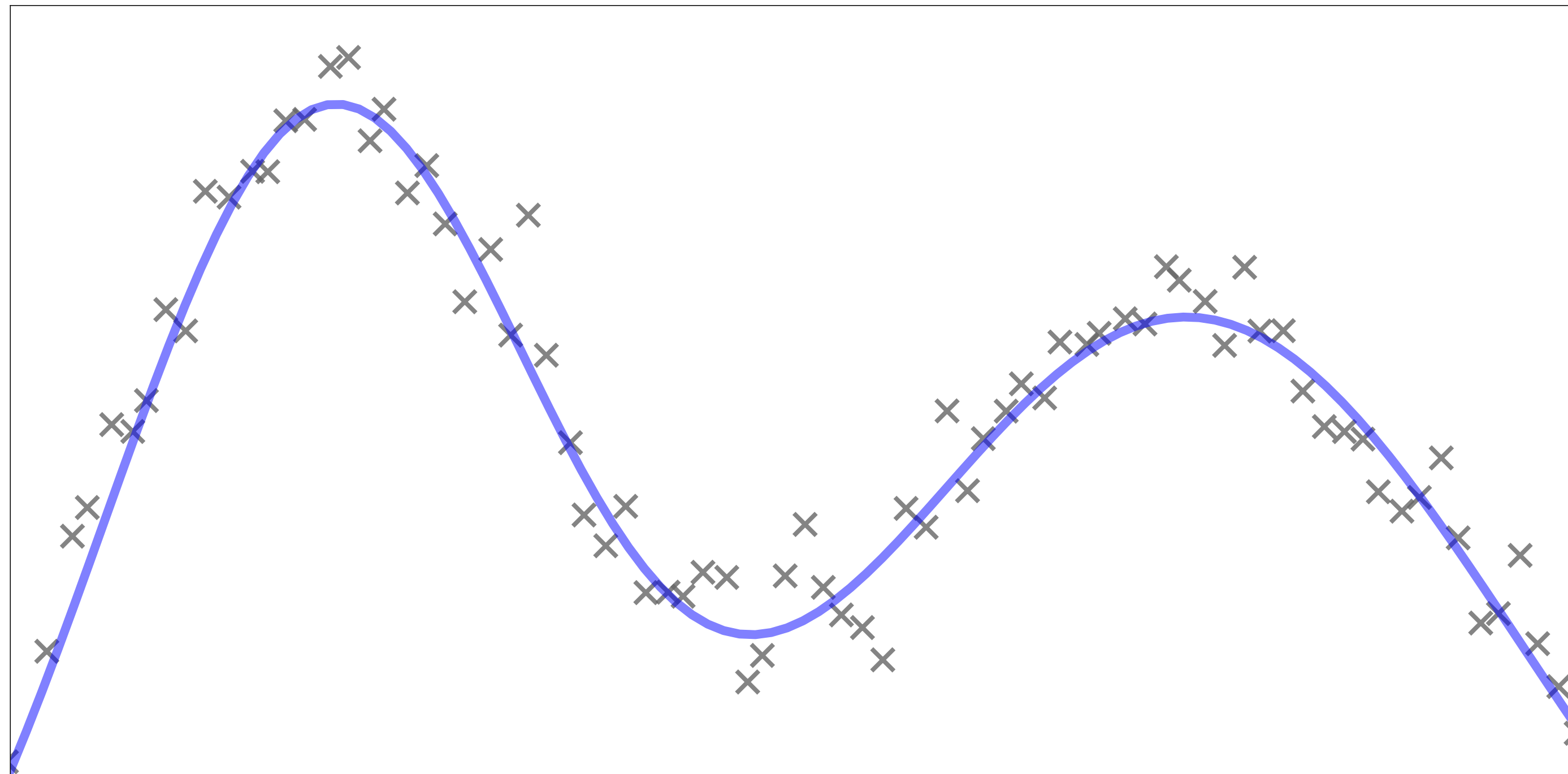
Multiple agents share data via a *naive* pool-and-share protocol:

- ▶ Everyone collects data, everyone gets a copy of the others' data.



Multiple agents share data via a *naive* pool-and-share protocol:

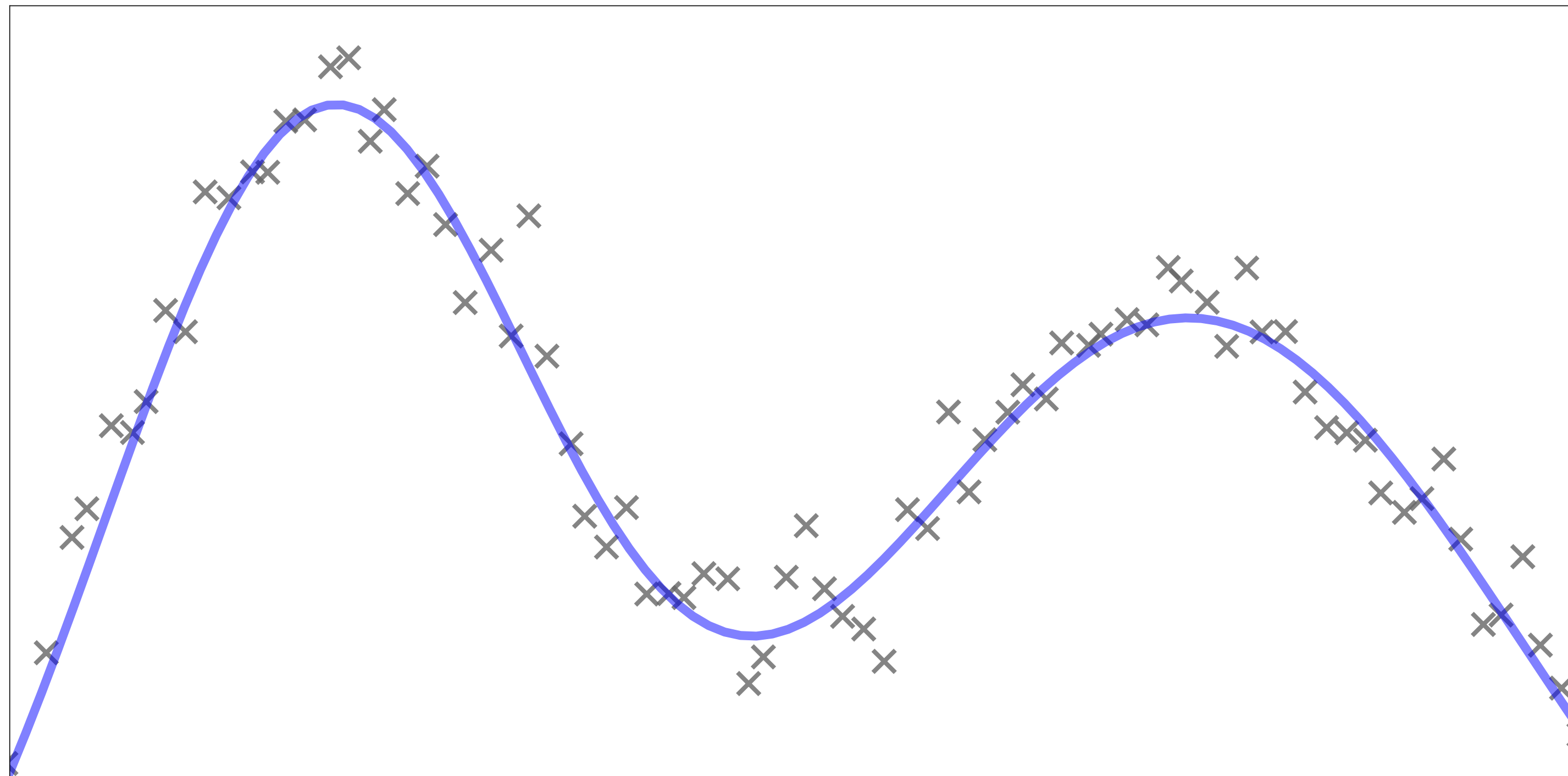
- ▶ Everyone collects data, everyone gets a copy of the others' data.



If others are already contributing large amounts of data, an agent has no incentive to collect/contribute data of her own.

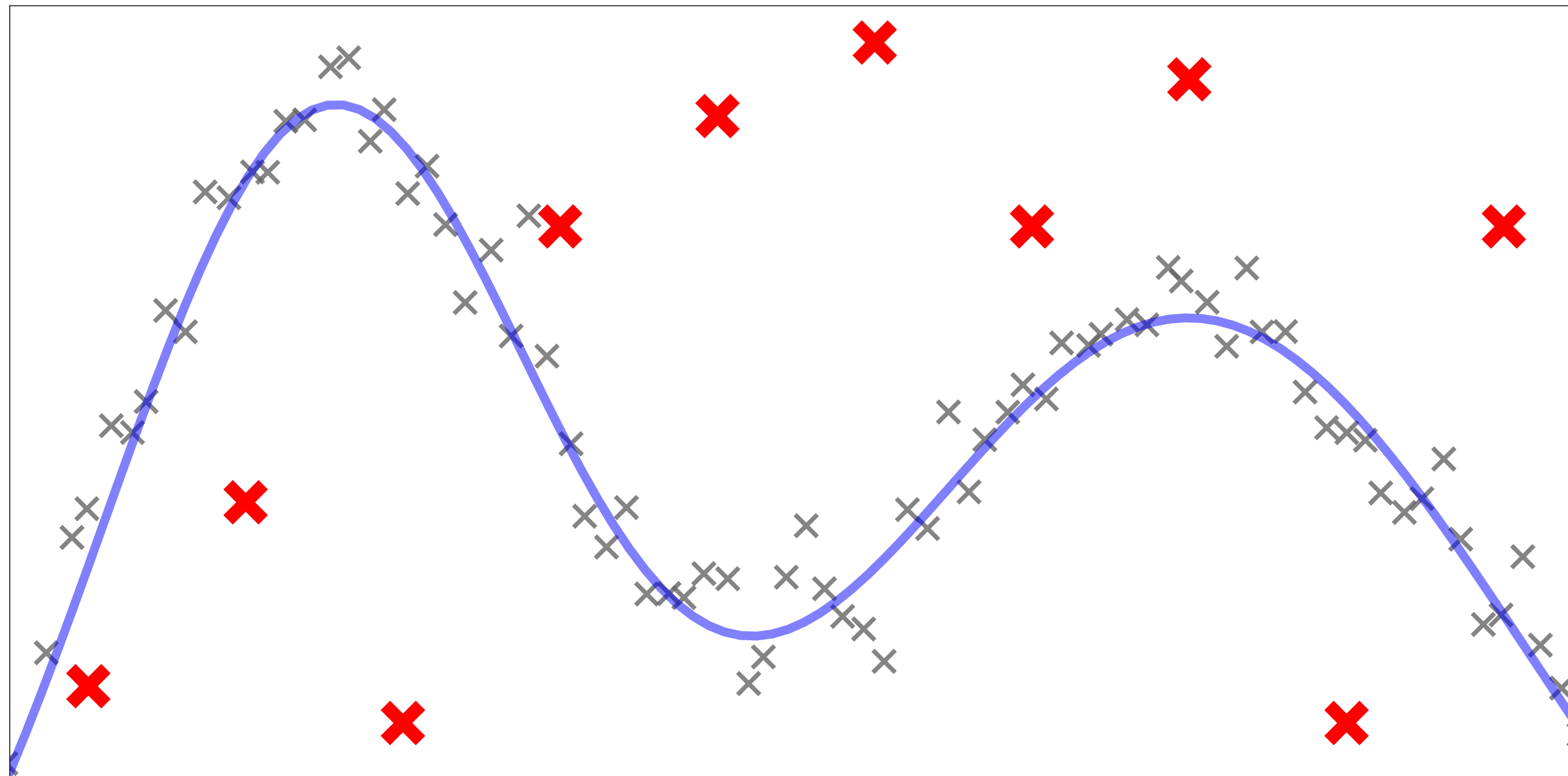
A seemingly plausible work-around (but does not work):

Pool-and-share but only if the agent contributes sufficient data



A seemingly plausible work-around (but does not work):

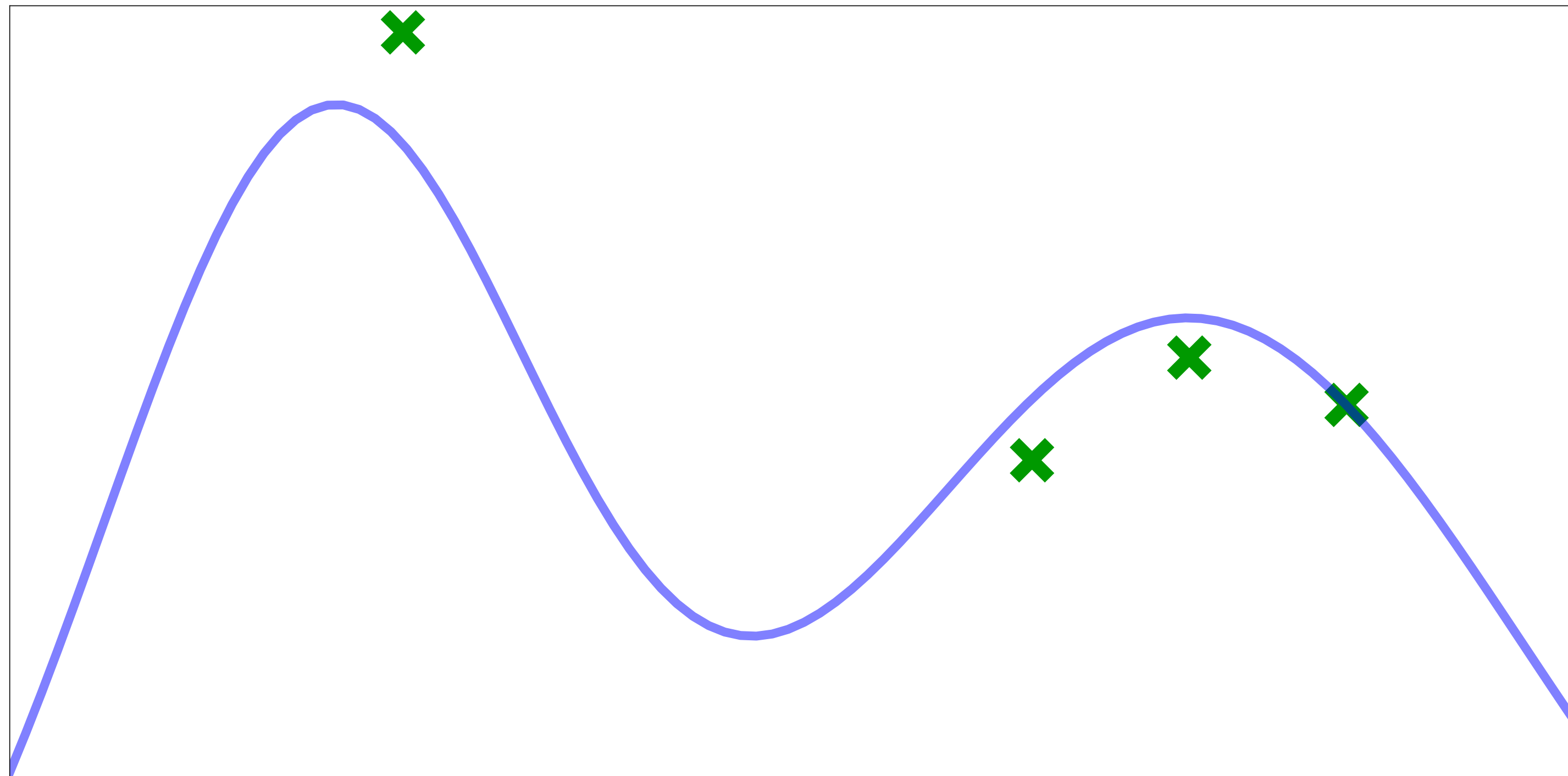
Pool-and-share but only if the agent contributes sufficient data



- ▶ Agent can submit fabricated data and then discard it when learning.

A seemingly plausible work-around (but does not work):

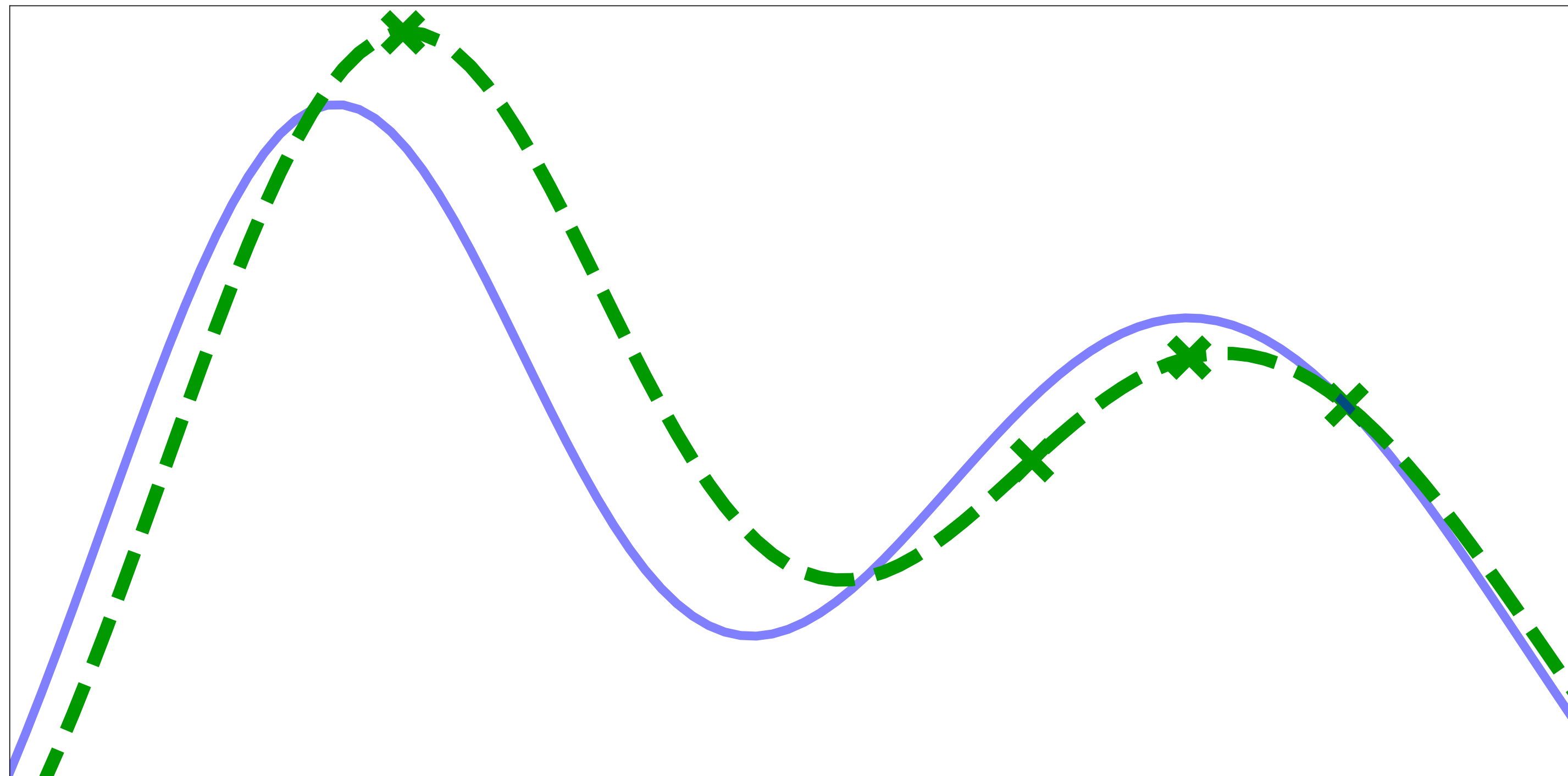
Pool-and-share but only if the agent contributes sufficient data



- ▶ Agent can submit fabricated data and then discard it when learning.

A seemingly plausible work-around (but does not work):

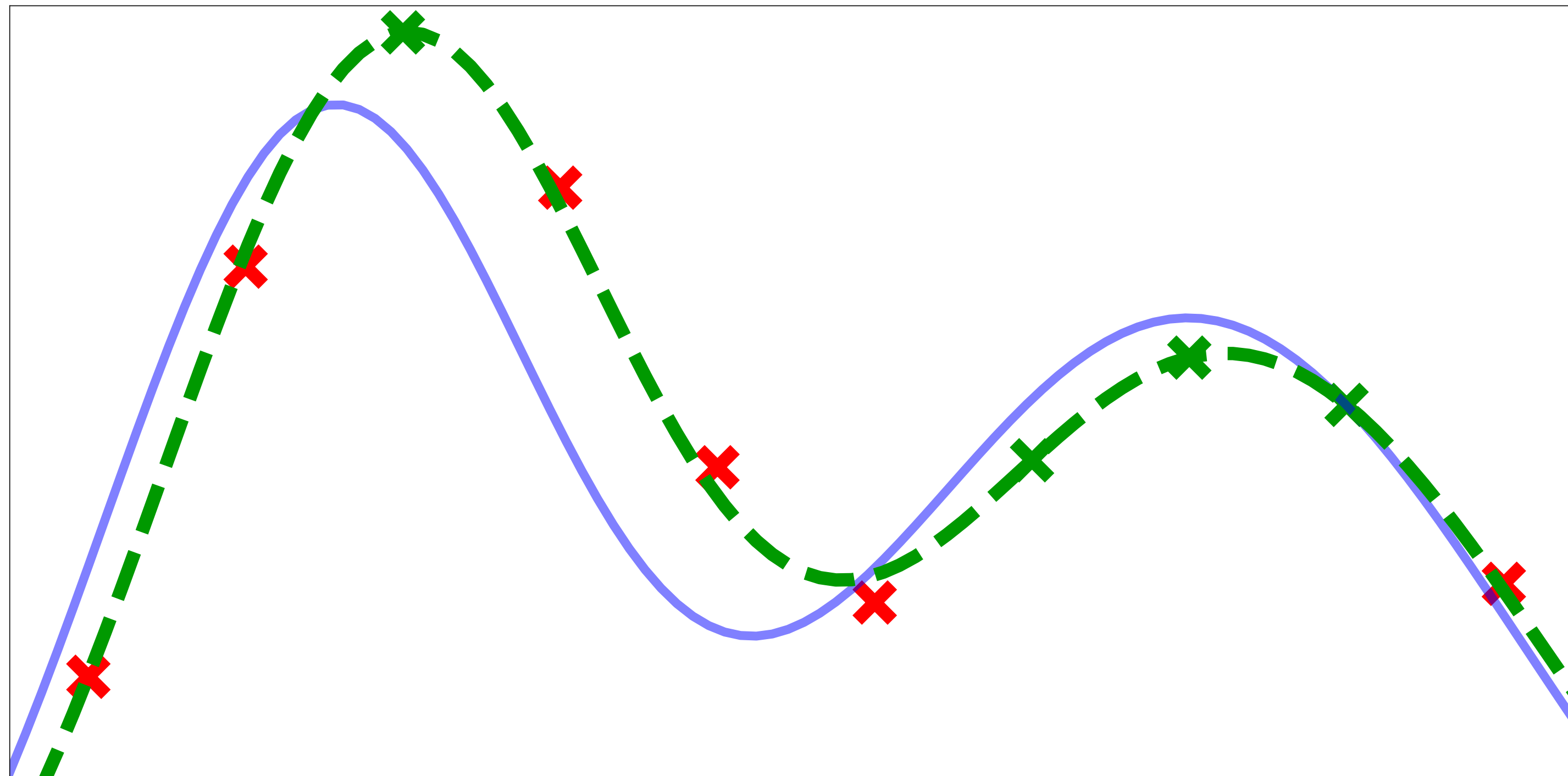
Pool-and-share but only if the agent contributes sufficient data



- ▶ Agent can submit fabricated data and then discard it when learning.

A seemingly plausible work-around (but does not work):

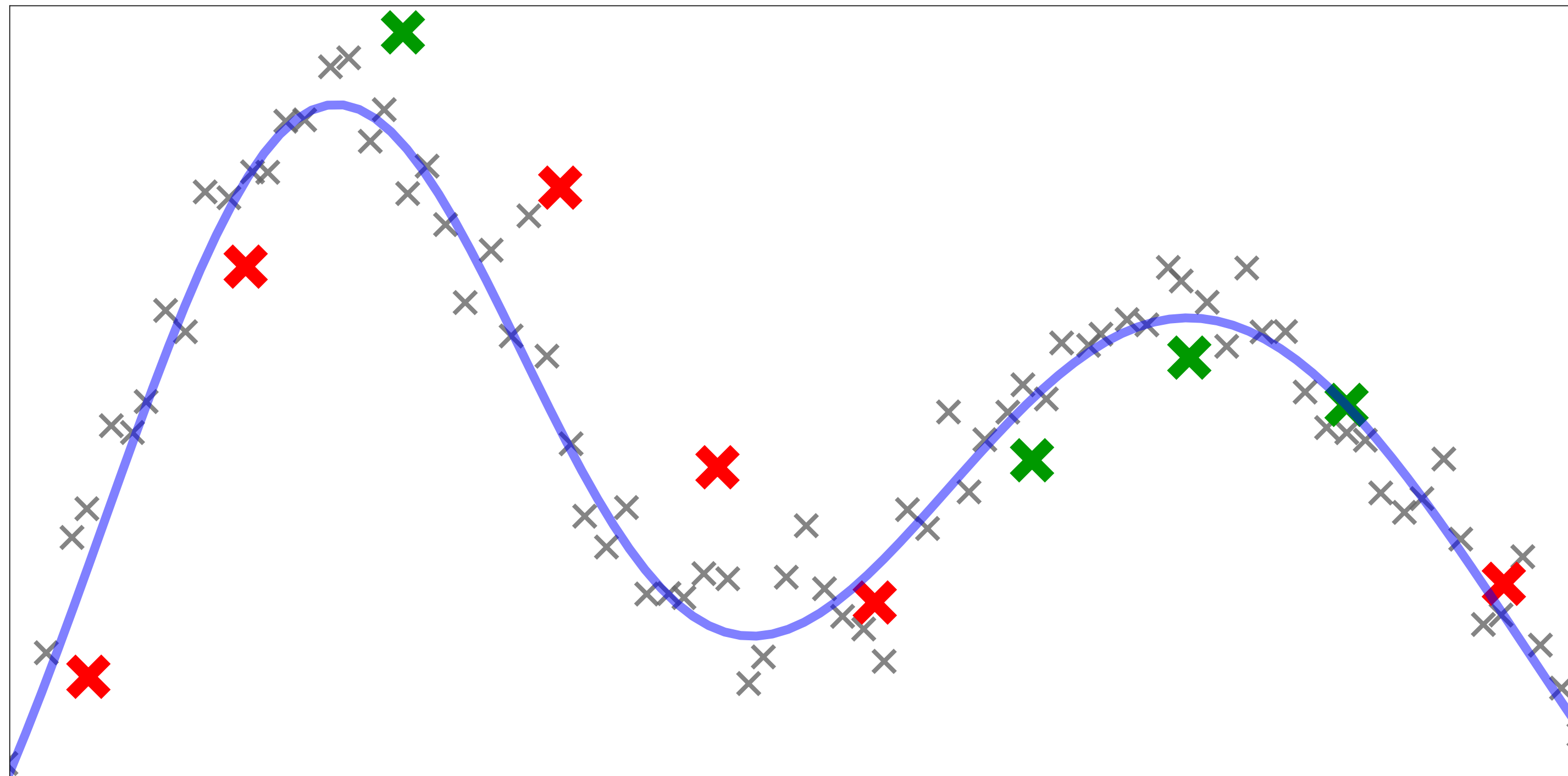
Pool-and-share but only if the agent contributes sufficient data



- ▶ Agent can submit fabricated data and then discard it when learning.
- ▶ Agent may fabricate based on a small sample she has collected, so it may not always be easy to detect.

A seemingly plausible work-around (but does not work):

Pool-and-share but only if the agent contributes sufficient data



- ▶ Agent can submit fabricated data and then discard it when learning.
- ▶ Agent may fabricate based on a small sample she has collected, so it may not always be easy to detect.

Ethical/Legal

Privacy
Ownership of data

Security

Data breaches
Adversarial attacks

Logistical

Inter-operability
Communication costs

Incentives

Free-riding
Competition

Data monetization
Data valuation

BUT THERE IS A DEMAND FOR DATA SHARING IN THE REAL WORLD

Data sharing platforms/consortia



An open standard for secure data sharing

Marketplaces for data and ML models

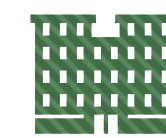


Mechanisms for data sharing and federated learning



Data marketplaces

Contributors



Marketplace

Consumers



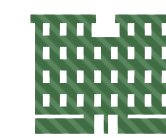
Mechanisms for data sharing and federated learning



Goal: Incentivize agents to collect as much data and share it honestly.

Data marketplaces

Contributors



Marketplace

Consumers



Mechanisms for data sharing and federated learning

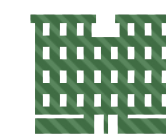


Goal: Incentivize agents to collect as much data and share it honestly.

- Do not simply pool and share data!

Data marketplaces

Contributors

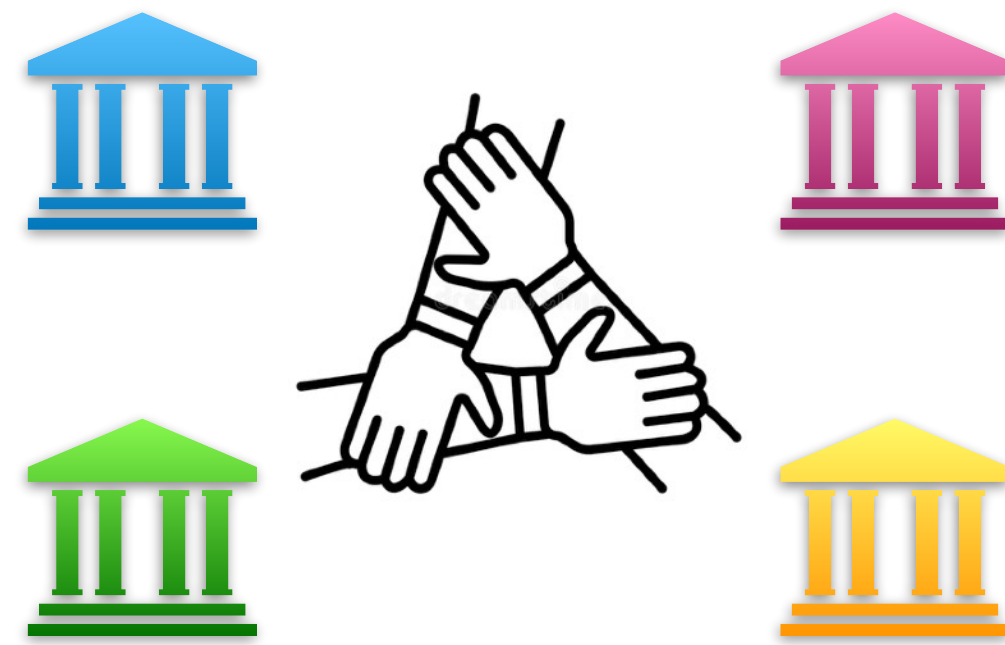


Marketplace

Consumers



Mechanisms for data sharing and federated learning



Goal: Incentivize agents to collect as much data and share it honestly.

- Do not simply pool and share data!
- Cross-check for quality of the data contributed.

Data marketplaces

Contributors



Marketplace

Consumers



Mechanisms for data sharing and federated learning

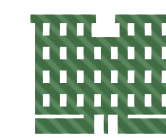


Goal: Incentivize agents to collect as much data and share it honestly.

- Do not simply pool and share data!
- Cross-check for quality of the data contributed.
- **More/better data contributed \implies more/better data received.**

Data marketplaces

Contributors



Marketplace

Consumers



Mechanisms for data sharing and federated learning

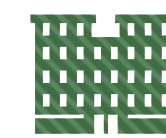


Goal: Incentivize agents to collect as much data and share it honestly.

- Do not simply pool and share data!
- Cross-check for quality of the data contributed.
- More/better data contributed \implies more/better data received.

Data marketplaces

Contributors



Marketplace

Consumers



Goal: Incentivize contributors to honestly contribute lots of data. Fairly reward them for effort via payments from consumers.

Mechanisms for data sharing and federated learning

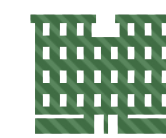


Goal: Incentivize agents to collect as much data and share it honestly.

- Do not simply pool and share data!
- Cross-check for quality of the data contributed.
- **More/better data contributed \implies more/better data received.**

Data marketplaces

Contributors



Marketplace

Consumers



Goal: Incentivize contributors to honestly contribute lots of data. Fairly reward them for effort via payments from consumers.

- **Higher quality data \implies higher revenue for data contributors.**

Mechanisms for data sharing and federated learning

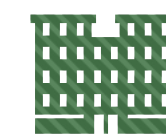


Goal: Incentivize agents to collect as much data and share it honestly.

- Do not simply pool and share data!
- Cross-check for quality of the data contributed.
- **More/better data contributed \implies more/better data received.**

Data marketplaces

Contributors



Marketplace

Consumers



Goal: Incentivize contributors to honestly contribute lots of data. Fairly reward them for effort via payments from consumers.

- **Higher quality data \implies higher revenue for data contributors.**
- A mediator checks for the quality of the data from contributors

Mechanisms for data sharing and federated learning

Sim, Zhang, Chan, Low 2020

Xu, Lyu, Ma et al 2021

Blum, Haghtalab, Phillips, Shao 2021

Karimireddy, Guo, Jordan 2022

Fraboni, Vidal, Lorenzi 2021

Lin, Du, Liu 2019

Ding, Fang, Huang 2020

Liu, Tian, Chen et al 2022

Data marketplaces

Cai, Daskalakis, Papadimitriou 2015

Agarwal, Dahleh, Sarkar, 2019

Agarwal, Dahleh, Horel, Rui, 2020

Jia, Dao, Wang et al, 2019

Wang, Rausch, Zhang et al 2020

Key difference:

- ▶ All these works assume agents will always truthfully submit the data they have, i.e without fabrication/alteration.

1. Mechanism design for collaborative normal mean estimation

(Chen, Zhu, Kandasamy, *NeurIPS 2023*)

- ▶ **Intuitions, overview of results**
- ▶ **Problem formalism**
- ▶ **Mechanism and theoretical analysis**

2. Extensions & Future work

- ▶ **Collaborative supervised learning, design of experiments**
- ▶ **Data marketplaces**

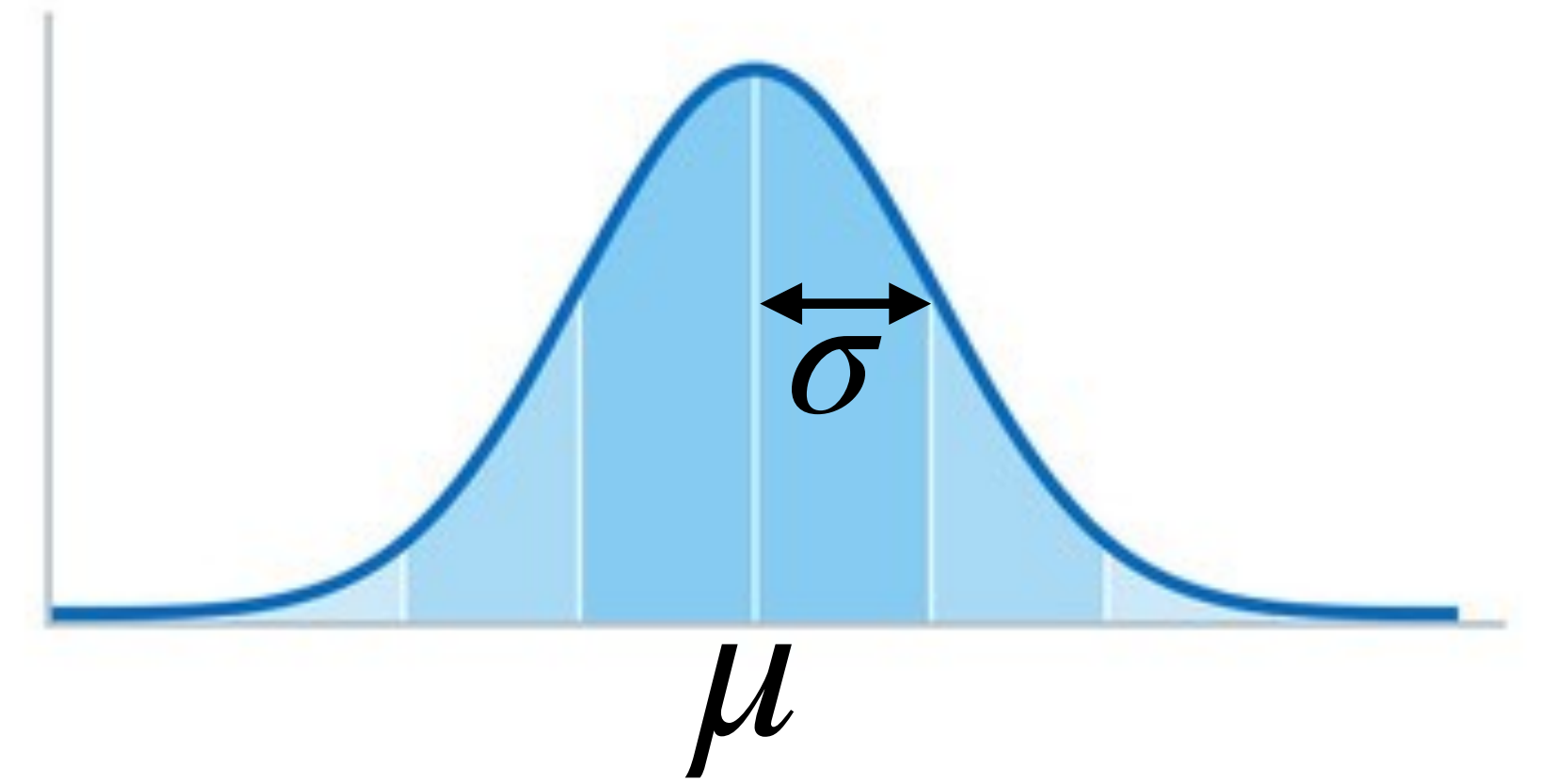
1. Mechanism design for collaborative normal mean estimation

(Chen, Zhu, Kandasamy, *NeurIPS 2023*)

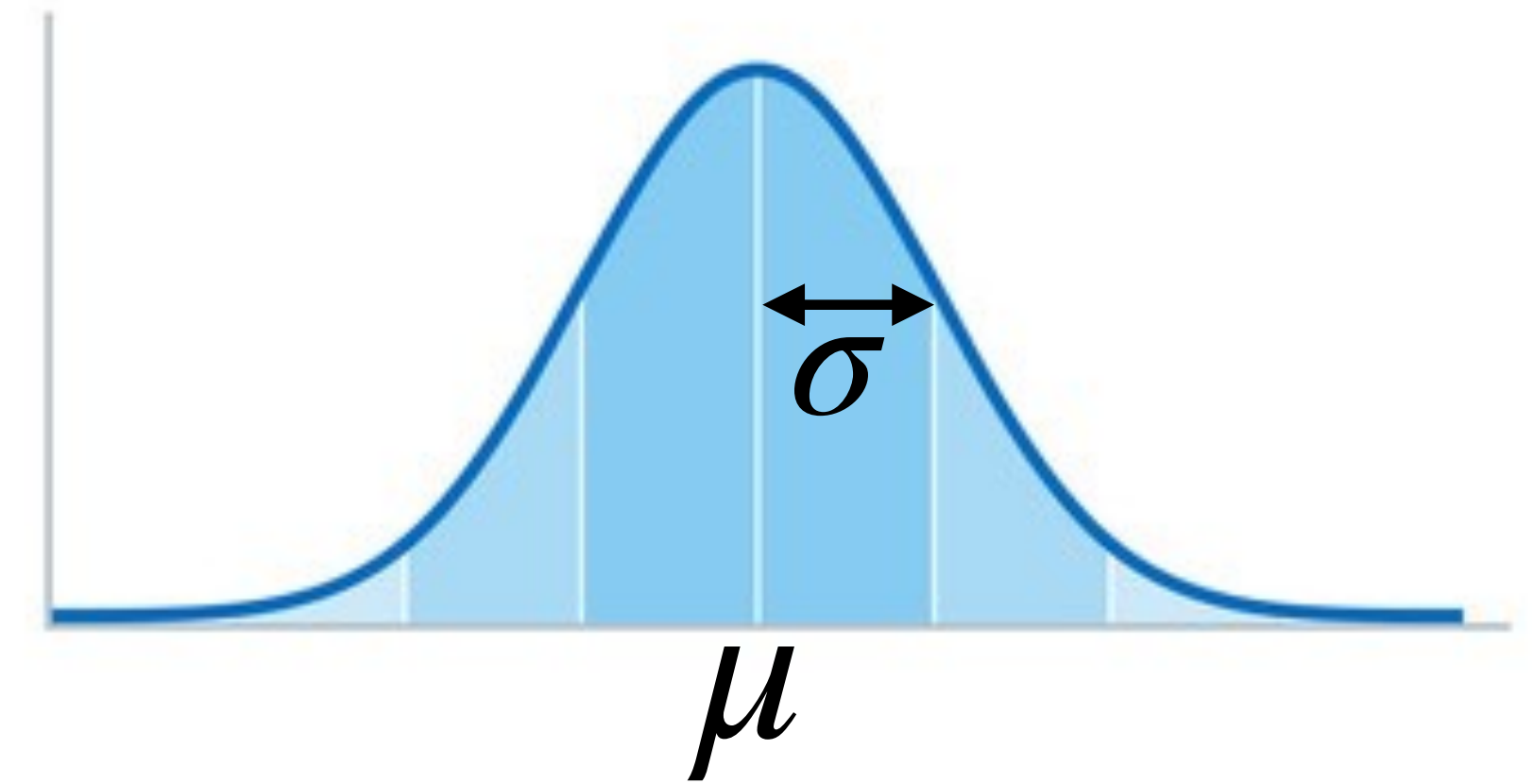
- ▶ **Intuitions and Challenges**
- ▶ Problem formalism
- ▶ Mechanism and theoretical analysis

2. Extensions & Future work

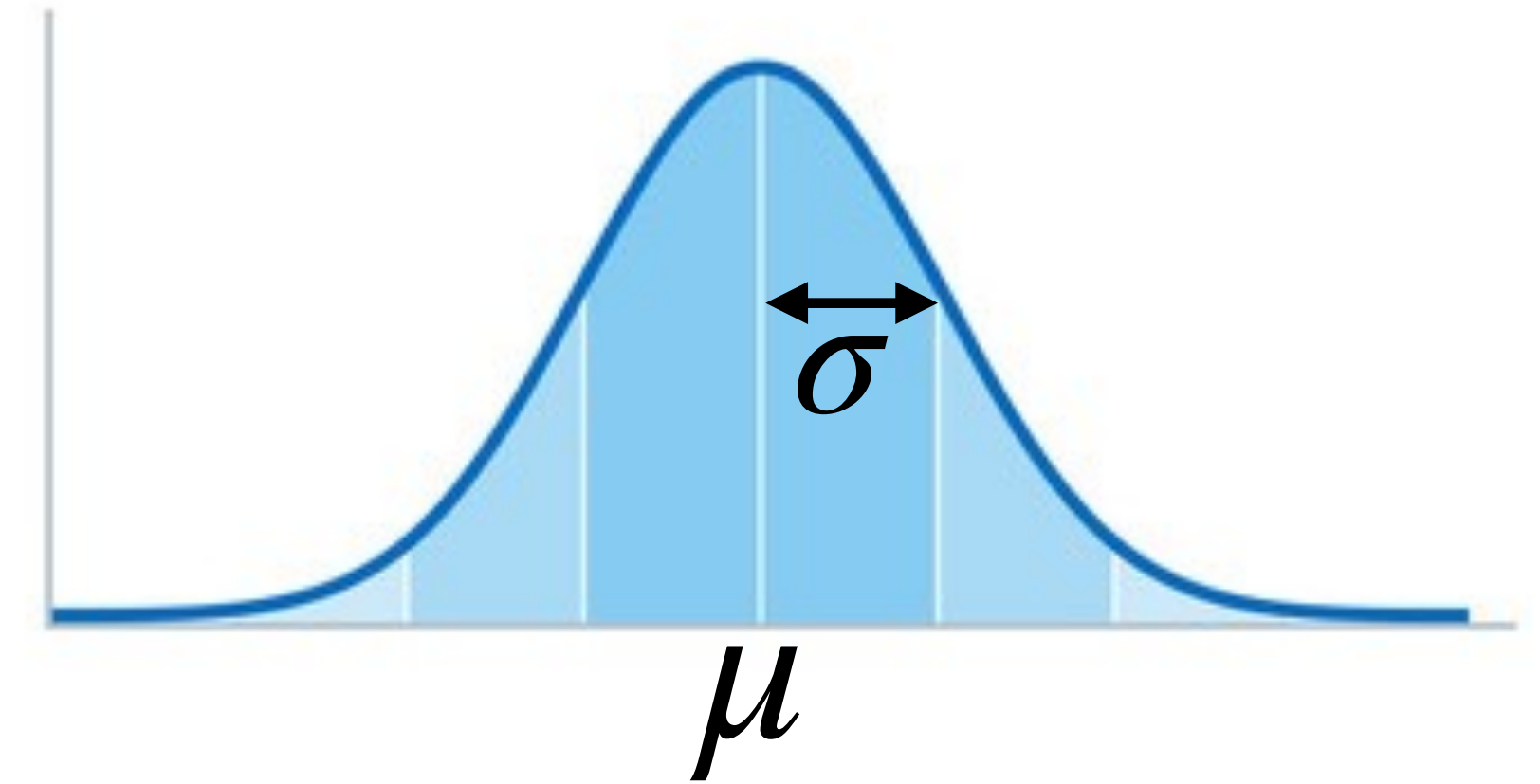
- ▶ Collaborative supervised learning, design of experiments
- ▶ Data marketplaces



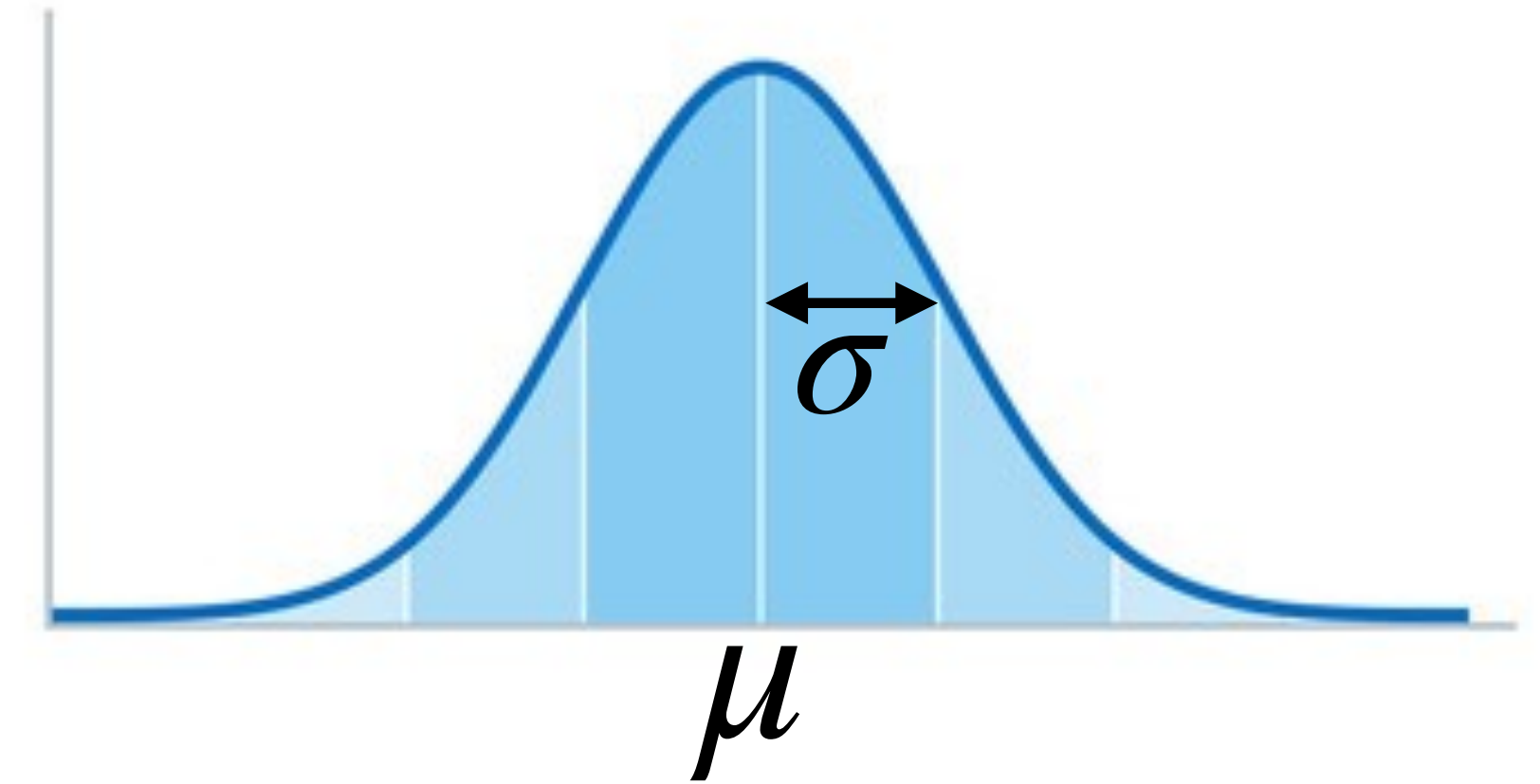
- Estimate the mean μ of a normal distribution with *known* variance σ^2 .



- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .



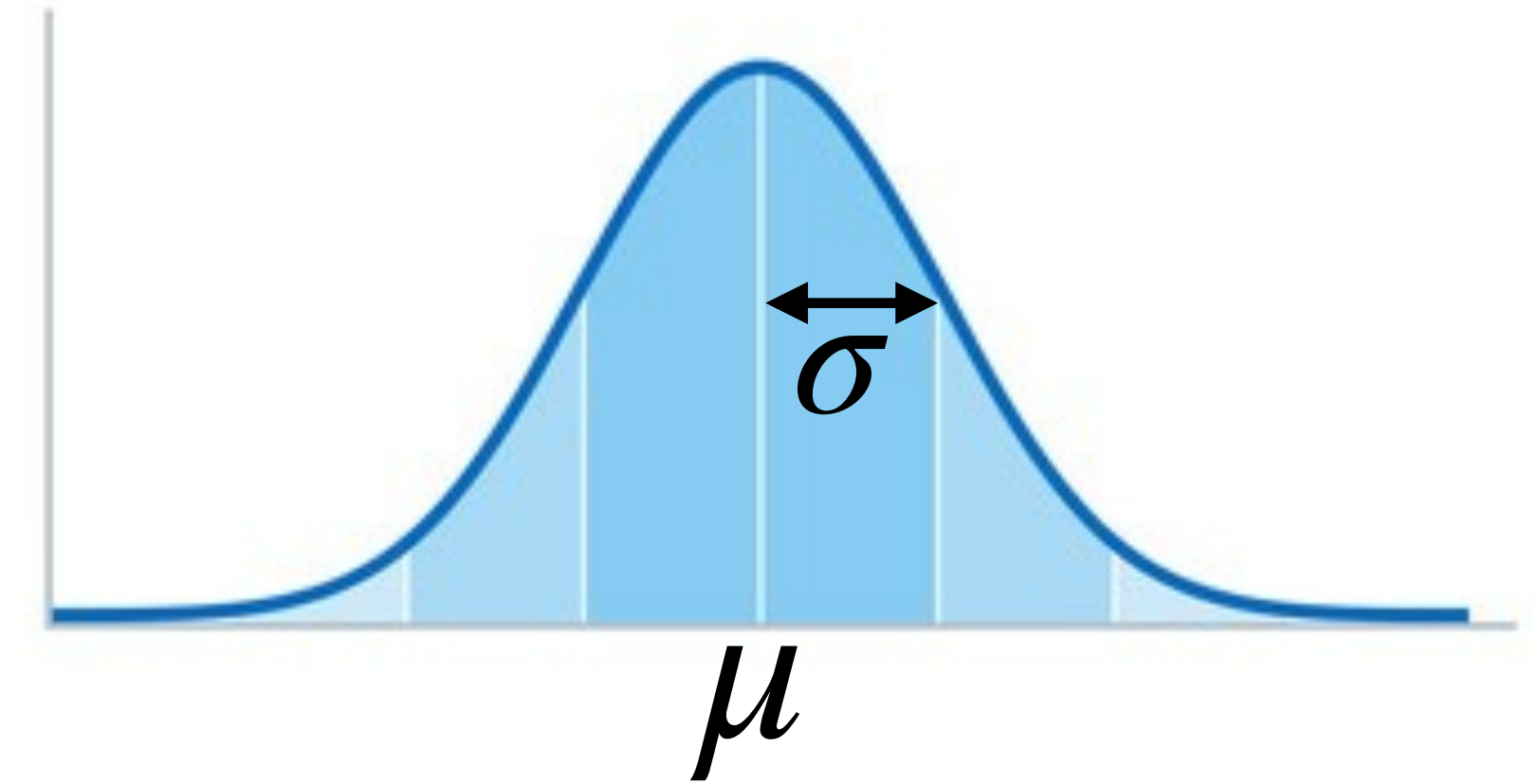
- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .
- Each agent wishes to minimize
penalty = estimation error + data collection cost



- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .
- Each agent wishes to minimize

penalty = estimation error + data collection cost

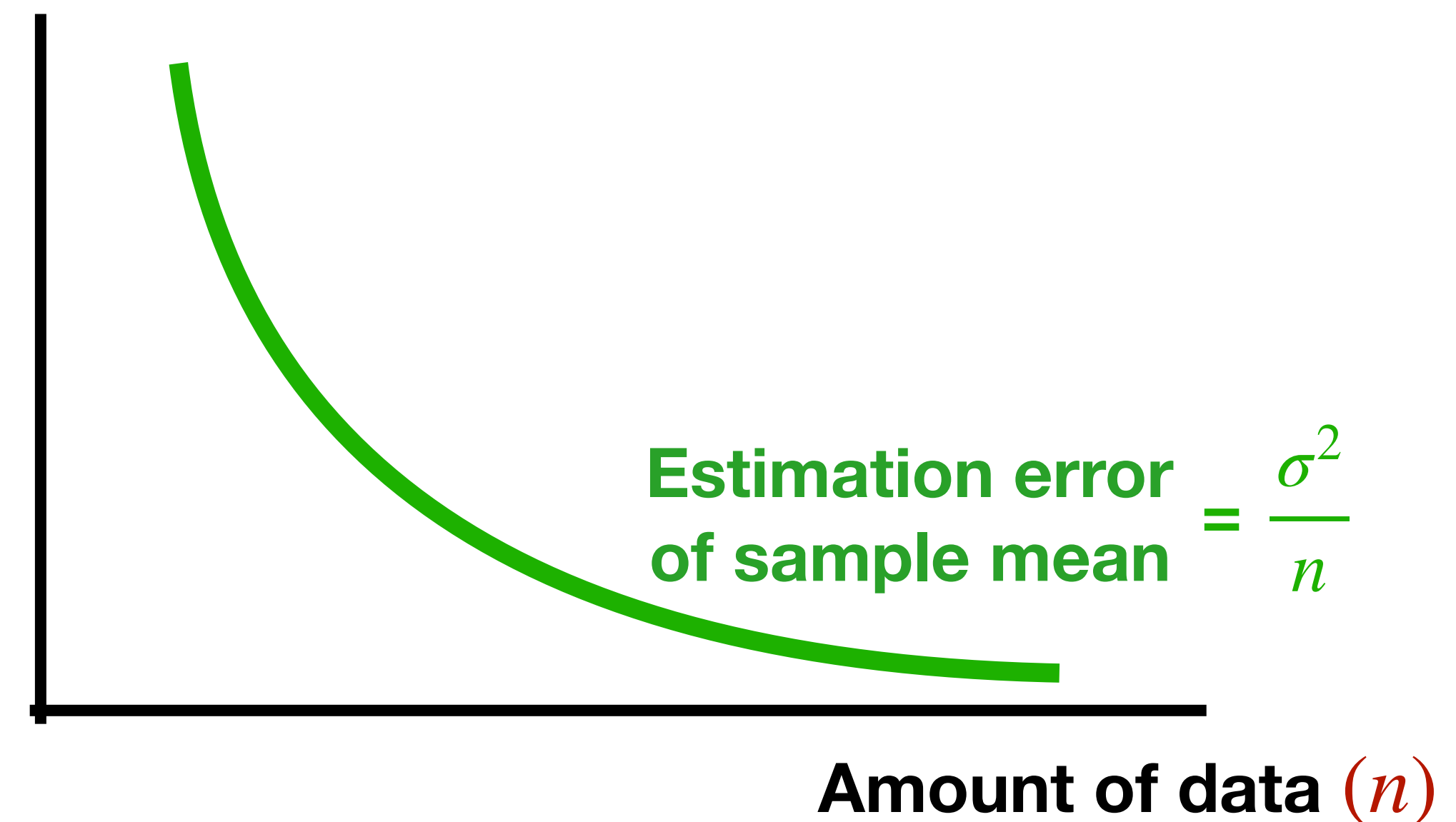
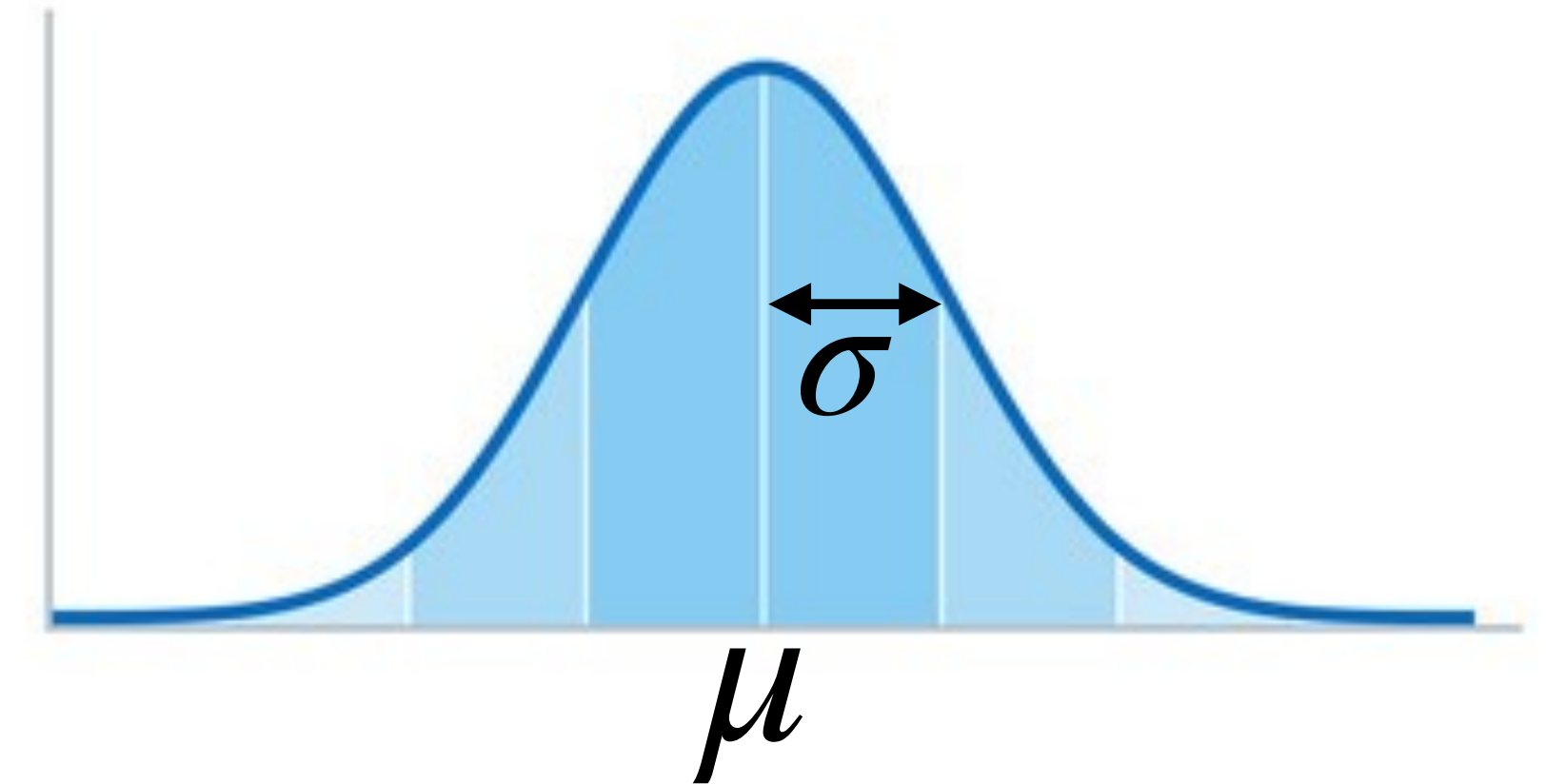
$$= \frac{\sigma^2}{n} + cn$$



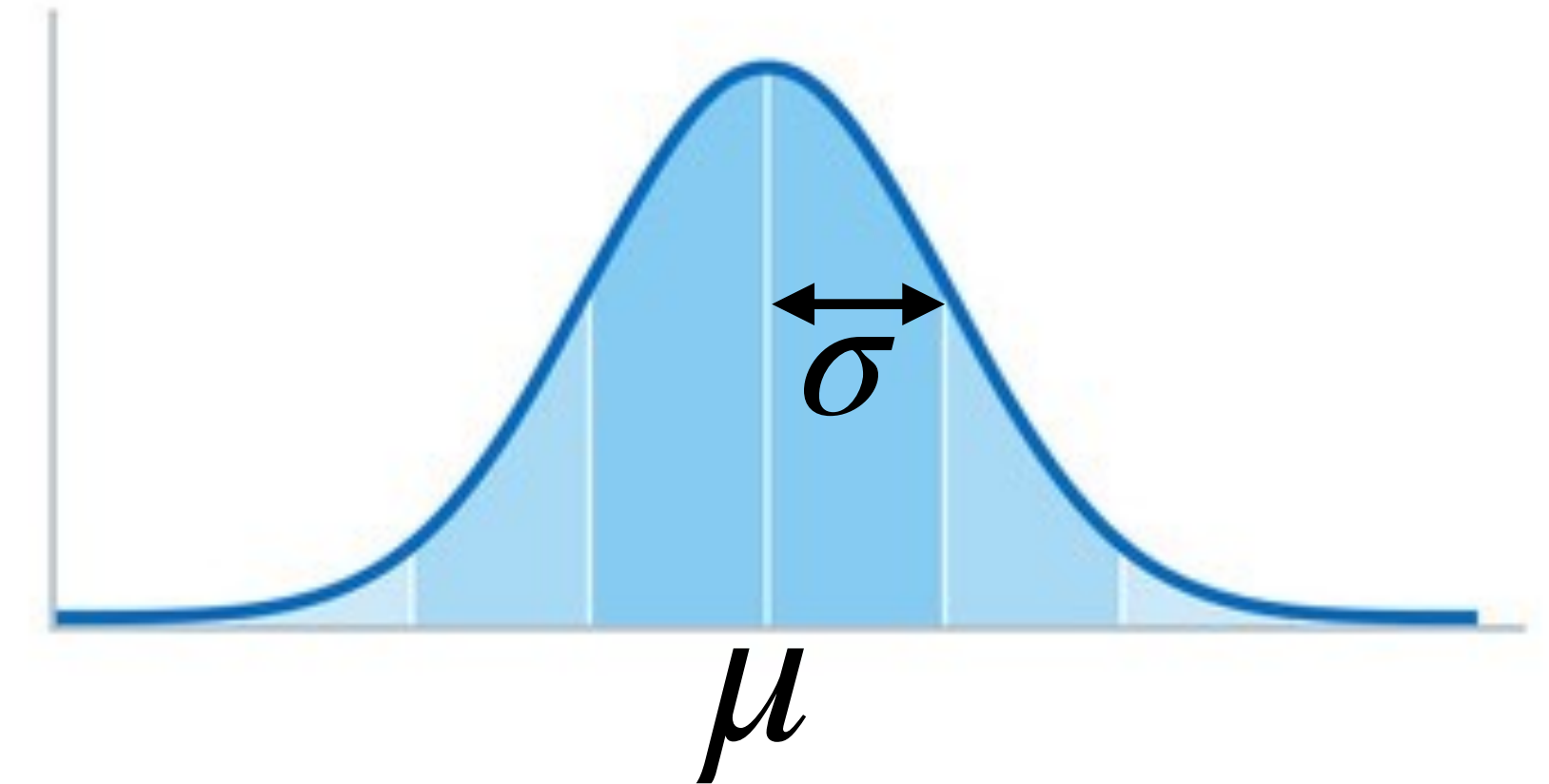
- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .
- Each agent wishes to minimize

penalty = estimation error + data collection cost

$$= \frac{\sigma^2}{n} + cn$$

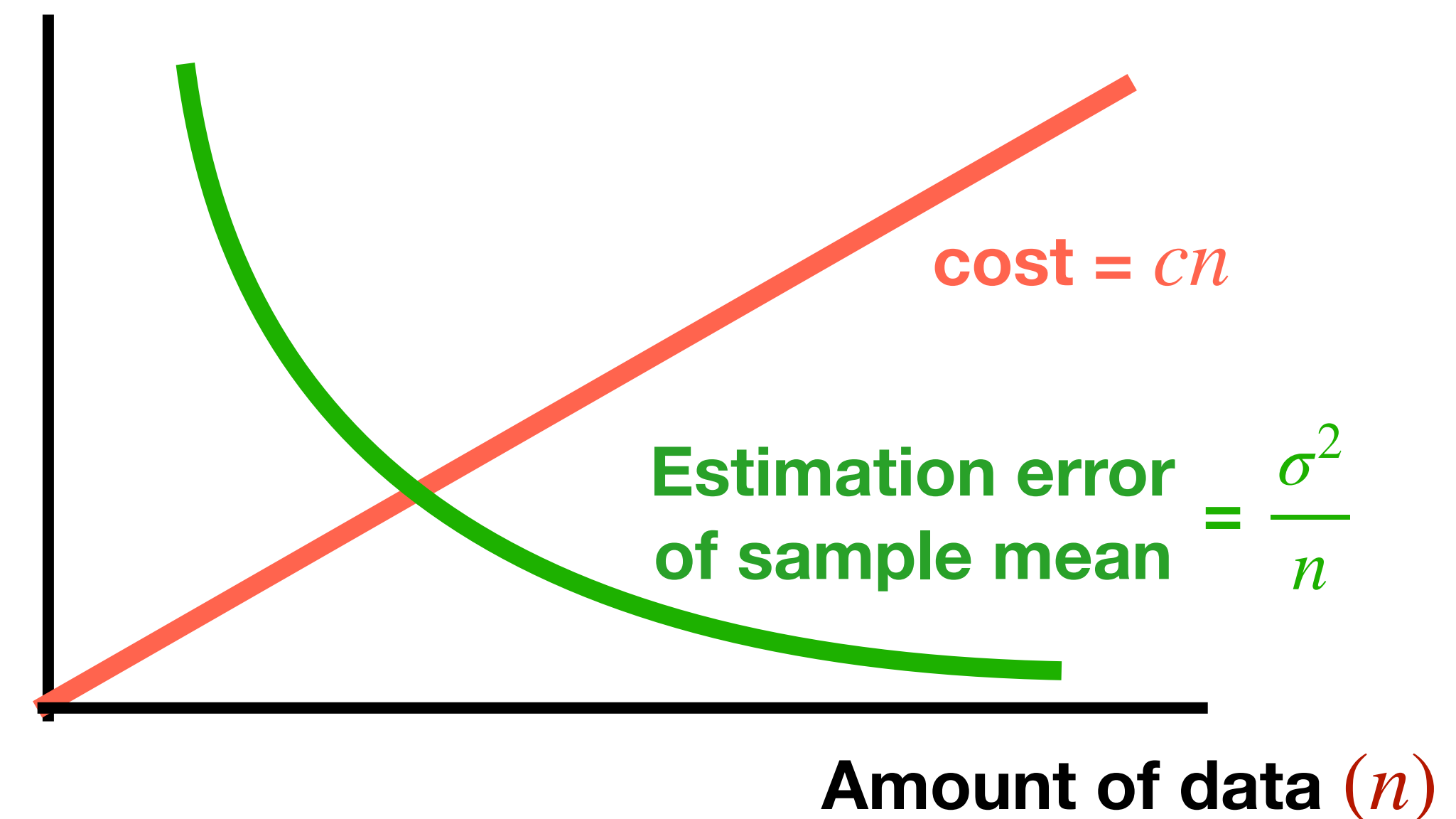


- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .
- Each agent wishes to minimize

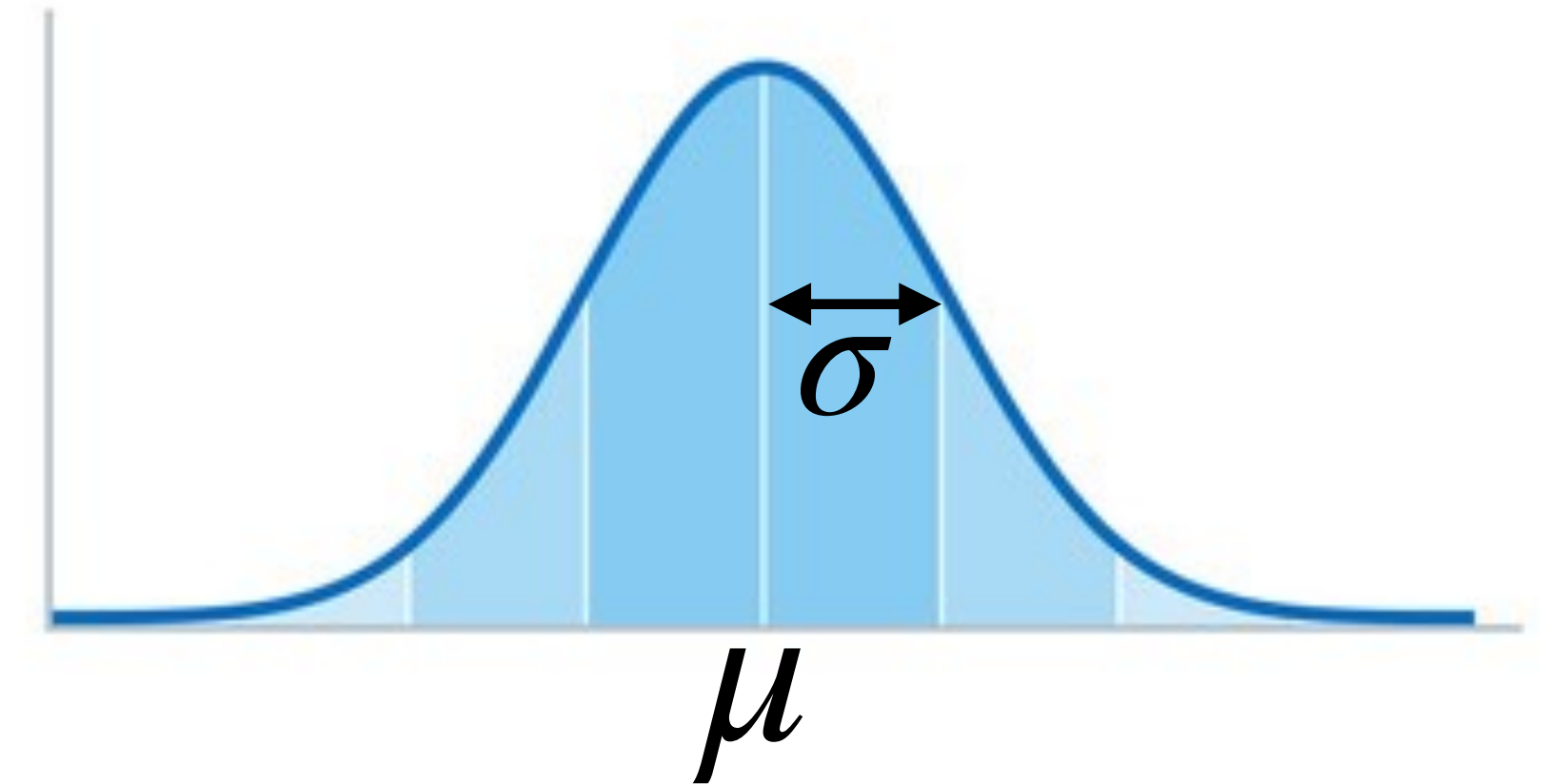


penalty = estimation error + data collection cost

$$= \frac{\sigma^2}{n} + cn$$

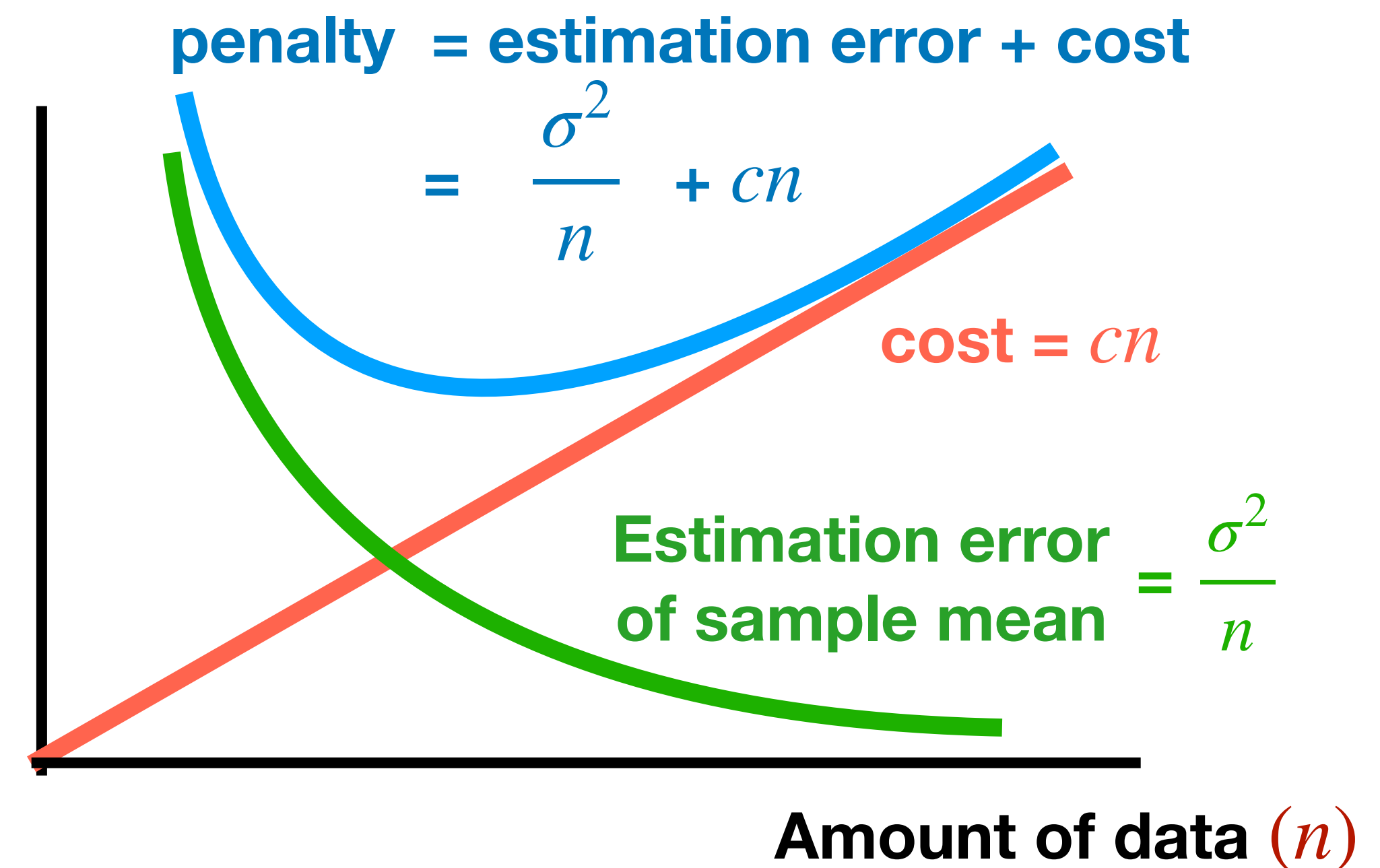


- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .
- Each agent wishes to minimize

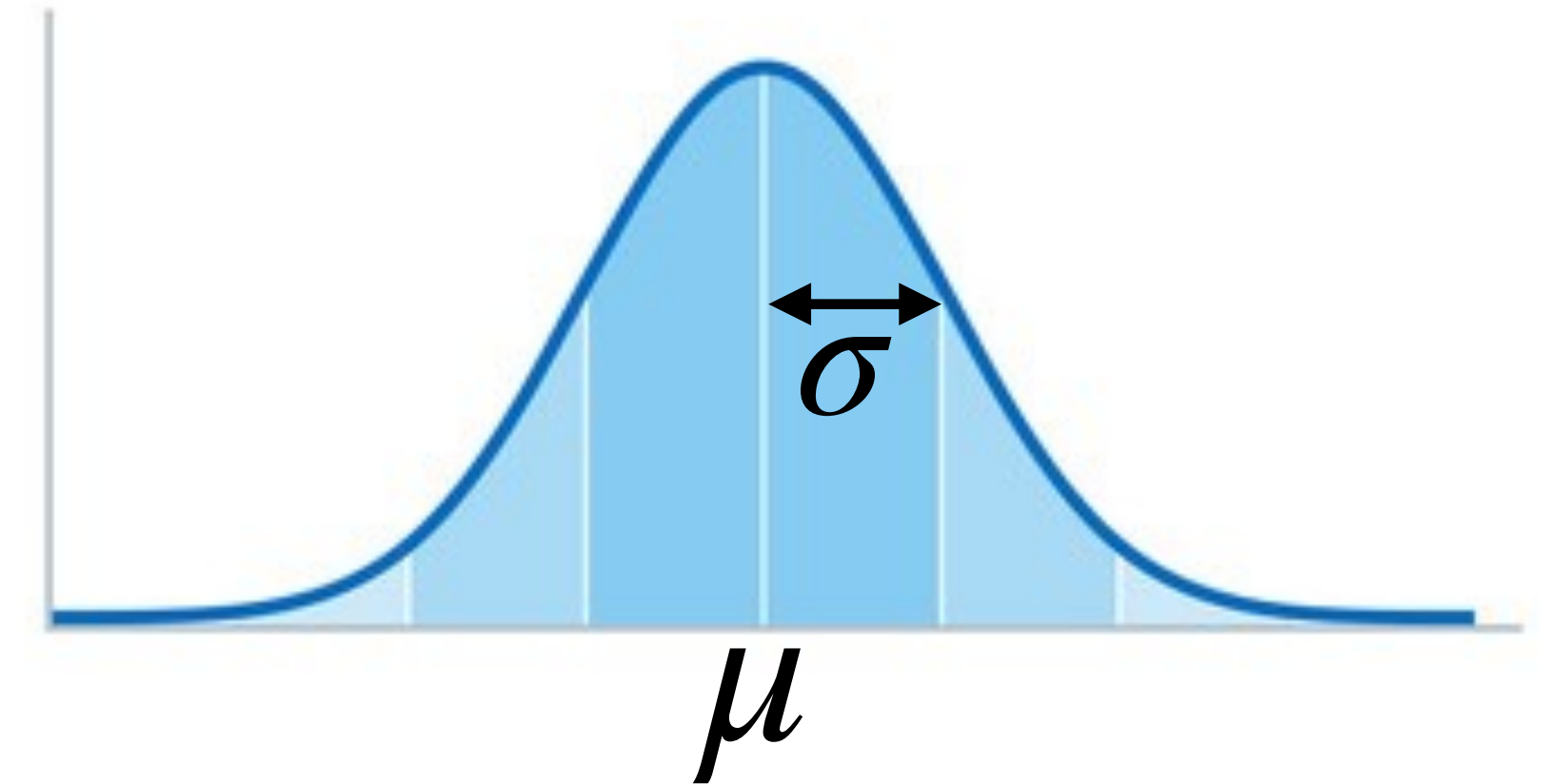


penalty = estimation error + data collection cost

$$= \frac{\sigma^2}{n} + cn$$



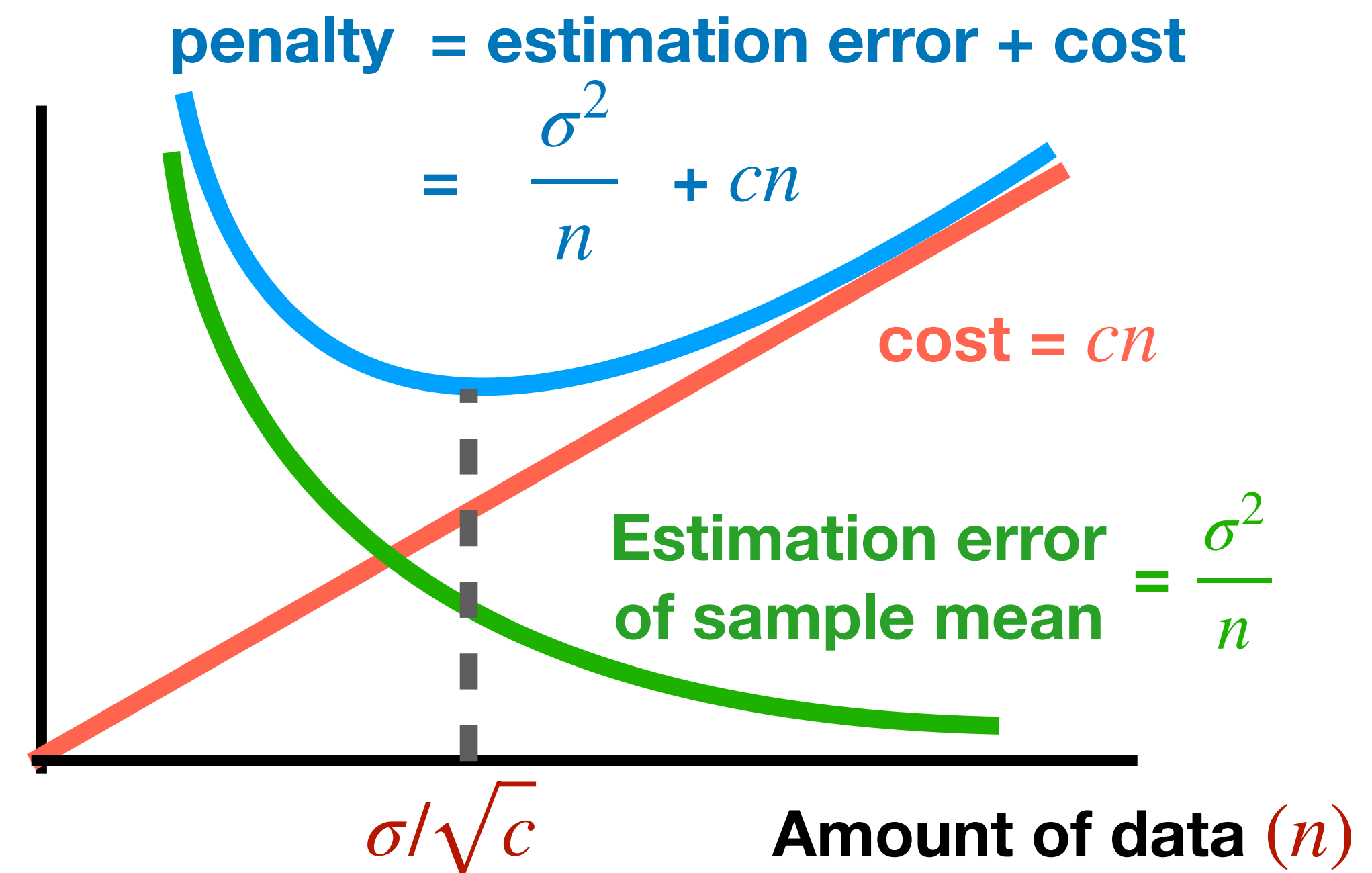
- Estimate the mean μ of a normal distribution with *known* variance σ^2 .
- An agent can collect samples at *known* unit cost c .
- Each agent wishes to minimize



penalty = estimation error + data collection cost

$$= \frac{\sigma^2}{n} + cn$$

- If working on her own, agent will collect σ/\sqrt{c} points to minimize penalty.



- Now consider m agents collecting and sharing their data.

- Now consider m agents collecting and sharing their data.
- Social penalty of all m agents if they collectively collect n_{tot} points.

$$\text{social penalty} = \text{estimation error of all agents} + \text{data collection cost} = m \times \frac{\sigma^2}{n_{\text{tot}}} + cn_{\text{tot}}$$

- Now consider m agents collecting and sharing their data.
- Social penalty of all m agents if they collectively collect n_{tot} points.

$$\text{social penalty} = \text{estimation error of all agents} + \text{data collection cost} = m \times \frac{\sigma^2}{n_{\text{tot}}} + cn_{\text{tot}}$$

- To minimize penalty, they should collect $n_{\text{tot}}^{\star} = \frac{\sigma\sqrt{m}}{\sqrt{c}}$ points.

- Now consider m agents collecting and sharing their data.
- Social penalty of all m agents if they collectively collect n_{tot} points.

$$\text{social penalty} = \text{estimation error of all agents} + \text{data collection cost} = m \times \frac{\sigma^2}{n_{\text{tot}}} + cn_{\text{tot}}$$

- To minimize penalty, they should collect $n_{\text{tot}}^{\star} = \frac{\sigma\sqrt{m}}{\sqrt{c}}$ points.

- Each agent needs to collect only $n^{\star} = \frac{\sigma}{\sqrt{mc}}$ points

Only $\times 1/\sqrt{m}$ when compared to working on her own (σ/\sqrt{c} points).

- Now consider m agents collecting and sharing their data.
- Social penalty of all m agents if they collectively collect n_{tot} points.

$$\text{social penalty} = \text{estimation error of all agents} + \text{data collection cost} = m \times \frac{\sigma^2}{n_{\text{tot}}} + cn_{\text{tot}}$$

- To minimize penalty, they should collect $n_{\text{tot}}^* = \frac{\sigma\sqrt{m}}{\sqrt{c}}$ points.

- Each agent needs to collect only $n^* = \frac{\sigma}{\sqrt{mc}}$ points

Only $\times 1/\sqrt{m}$ when compared to working on her own (σ/\sqrt{c} points).

- But she has $\times \sqrt{m}$ data.

	Amount of data she needs to collect (n_i)	Amount of data available to her (n_{tot})	Penalty $\frac{\sigma^2}{n_{\text{tot}}} + cn_i$
Working on her own			
Working together			

	Amount of data she needs to collect (n_i)	Amount of data available to her (n_{tot})	Penalty $\frac{\sigma^2}{n_{\text{tot}}} + cn_i$
Working on her own	$\frac{\sigma}{\sqrt{c}}$		
Working together	$\frac{\sigma}{\sqrt{cm}}$		

	Amount of data she needs to collect (n_i)	Amount of data available to her (n_{tot})	Penalty $\frac{\sigma^2}{n_{\text{tot}}} + cn_i$
Working on her own	$\frac{\sigma}{\sqrt{c}}$	$\frac{\sigma}{\sqrt{c}}$	
Working together	$\frac{\sigma}{\sqrt{cm}}$	$\frac{\sigma\sqrt{m}}{\sqrt{c}}$	

	Amount of data she needs to collect (n_i)	Amount of data available to her (n_{tot})	Penalty $\frac{\sigma^2}{n_{\text{tot}}} + cn_i$
Working on her own	$\frac{\sigma}{\sqrt{c}}$	$\frac{\sigma}{\sqrt{c}}$	$2\sigma\sqrt{c}$
Working together	$\frac{\sigma}{\sqrt{cm}}$	$\frac{\sigma\sqrt{m}}{\sqrt{c}}$	$\frac{2\sigma\sqrt{c}}{\sqrt{m}}$

	Amount of data she needs to collect (n_i)	Amount of data available to her (n_{tot})	Penalty $\frac{\sigma^2}{n_{\text{tot}}} + cn_i$
Working on her own	$\frac{\sigma}{\sqrt{c}}$	$\frac{\sigma}{\sqrt{c}}$	$2\sigma\sqrt{c}$
Working together	$\frac{\sigma}{\sqrt{cm}}$	$\frac{\sigma\sqrt{m}}{\sqrt{c}}$	$\frac{2\sigma\sqrt{c}}{\sqrt{m}}$

Agents can reduce data collection costs, and improve estimation error by sharing data with others.

- ▶ Naive mechanism 1: “pool and share”

- ▶ Naive mechanism 1: “pool and share”
 - ▶ Selfish agents will *free-ride*: not collecting any data, but using the data that the others have contributed.

$$\text{penalty} = \frac{\sigma^2}{n_{\text{tot}}} + c \times n_i$$

- ▶ Naive mechanism 1: “pool and share”
 - ▶ Selfish agents will *free-ride*: not collecting any data, but using the data that the others have contributed.

$$\text{penalty} = \frac{\sigma^2}{n_{\text{tot}}} + c \times n_i = \frac{\sigma^2}{(m-1) \times \frac{\sigma}{\sqrt{mc}}} + c \times 0$$

- ▶ Naive mechanism 1: “pool and share”
 - ▶ Selfish agents will *free-ride*: not collecting any data, but using the data that the others have contributed.

$$\text{penalty} = \frac{\sigma^2}{n_{\text{tot}}} + c \times n_i = \frac{\sigma^2}{(m-1) \times \frac{\sigma}{\sqrt{mc}}} + c \times 0 \approx \frac{\sigma\sqrt{c}}{\sqrt{m}}$$

- ▶ Naive mechanism 1: “pool and share”
 - ▶ Selfish agents will *free-ride*: not collecting any data, but using the data that the others have contributed.

$$\text{penalty} = \frac{\sigma^2}{n_{\text{tot}}} + c \times n_i = \frac{\sigma^2}{(m-1) \times \frac{\sigma}{\sqrt{mc}}} + c \times 0 \approx \frac{\sigma\sqrt{c}}{\sqrt{m}} = \frac{1}{2} \times \underbrace{\frac{2\sigma\sqrt{c}}{\sqrt{m}}}_{\text{penalty for a well-behaved agent}}$$

- ▶ Naive mechanism 1: “pool and share”
 - ▶ Selfish agents will *free-ride*: not collecting any data, but using the data that the others have contributed.

$$\text{penalty} = \frac{\sigma^2}{n_{\text{tot}}} + c \times n_i = \frac{\sigma^2}{(m-1) \times \frac{\sigma}{\sqrt{mc}}} + c \times 0 \approx \frac{\sigma\sqrt{c}}{\sqrt{m}} = \frac{1}{2} \times \underbrace{\frac{2\sigma\sqrt{c}}{\sqrt{m}}}_{\text{penalty for a well-behaved agent}}$$

- ▶ Naive mechanism 2: “pool and share, but only if you contribute enough data”

- ▶ Naive mechanism 1: “pool and share”
 - ▶ Selfish agents will *free-ride*: not collecting any data, but using the data that the others have contributed.

$$\text{penalty} = \frac{\sigma^2}{n_{\text{tot}}} + c \times n_i = \frac{\sigma^2}{(m-1) \times \frac{\sigma}{\sqrt{mc}}} + c \times 0 \approx \frac{\sigma\sqrt{c}}{\sqrt{m}} = \frac{1}{2} \times \underbrace{\frac{2\sigma\sqrt{c}}{\sqrt{m}}}_{\text{penalty for a well-behaved agent}}$$

- ▶ Naive mechanism 2: “pool and share, but only if you contribute enough data”
 - ▶ Agents can fabricate data, and then discard it after receiving others’ data.

Each agent i will:

- Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\}$.

Agents may collect any number of points, and lie (e.g withhold, lie, fabricate) about what they collect.

Each agent i will:

- Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\}$.

Agents may collect any number of points, and lie (e.g withhold, lie, fabricate) about what they collect.

The mechanism:

- To each agent, allocates a noisy version A_i of the others' data. The noise is proportional to how much the agent's submission Y_i differs from the others' submissions $\{Y_j\}_{j \neq i}$.

Each agent i will:

- Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\}$.

Agents may collect any number of points, and lie (e.g withhold, lie, fabricate) about what they collect.

The mechanism:

- To each agent, allocates a noisy version A_i of the others' data. The noise is proportional to how much the agent's submission Y_i differs from the others' submissions $\{Y_j\}_{j \neq i}$.

Each agent i will:

- Estimate μ using all the information they have (X_i, Y_i, A_i) .

Each agent i will:

- Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\}$.

Agents may collect any number of points, and lie (e.g withhold, lie, fabricate) about what they collect.

The mechanism:

- To each agent, allocates a noisy version A_i of the others' data. The noise is proportional to how much the agent's submission Y_i differs from the others' submissions $\{Y_j\}_{j \neq i}$.

Each agent i will:

- Estimate μ using all the information they have (X_i, Y_i, A_i) .

▶ *We design a (minimax) optimal estimator to enforce truthful reporting.*

This mechanism is

This mechanism is

- ▶ **Nash Incentive-compatible:** Provided that other agents are well-behaved, the best strategy for an agent is to,

This mechanism is

- ▶ **Nash Incentive-compatible:** Provided that other agents are well-behaved, the best strategy for an agent is to,
 - ▶ Collect a sufficient amount ($n^* = \sigma/\sqrt{mc}$) of data.

This mechanism is

- ▶ **Nash Incentive-compatible:** Provided that other agents are well-behaved, the best strategy for an agent is to,
 - ▶ Collect a sufficient amount ($n^* = \sigma/\sqrt{mc}$) of data.
 - ▶ Submit it truthfully.

This mechanism is

- ▶ **Nash Incentive-compatible:** Provided that other agents are well-behaved, the best strategy for an agent is to,
 - ▶ Collect a sufficient amount ($n^* = \sigma/\sqrt{mc}$) of data.
 - ▶ Submit it truthfully.
 - ▶ Use the recommended minimax-optimal estimator.

This mechanism is

- ▶ **Nash Incentive-compatible:** Provided that other agents are well-behaved, the best strategy for an agent is to,
 - ▶ Collect a sufficient amount ($n^* = \sigma/\sqrt{mc}$) of data.
 - ▶ Submit it truthfully.
 - ▶ Use the recommended minimax-optimal estimator.
- ▶ **Individually rational:** Provided that others are well-behaved, an agent does not do worse than the best she could do on her own.

This mechanism is

- ▶ **Nash Incentive-compatible:** Provided that other agents are well-behaved, the best strategy for an agent is to,
 - ▶ Collect a sufficient amount ($n^* = \sigma/\sqrt{mc}$) of data.
 - ▶ Submit it truthfully.
 - ▶ Use the recommended minimax-optimal estimator.
- ▶ **Individually rational:** Provided that others are well-behaved, an agent does not do worse than the best she could do on her own.
- ▶ **Efficient:** Social penalty at the Nash strategies is at most a factor 2 of the global minimum.

1. Mechanism design for collaborative normal mean estimation

(Chen, Zhu, Kandasamy, Neurips 2023)

- ▶ Intuitions and Challenges
- ▶ **Problem formalism**
- ▶ Mechanism and theoretical analysis

2. Extensions & Future work

- ▶ Collaborative supervised learning, design of experiments
- ▶ Data marketplaces

A mechanism M receives a dataset from each agent, and returns an *allocation* A_i to each agent.

A mechanism M receives a dataset from each agent, and returns an *allocation* A_i to each agent.

- ▶ The mechanism designer can choose a space of allocations \mathcal{A} to obtain desirable outcomes.

A mechanism M receives a dataset from each agent, and returns an *allocation* A_i to each agent.

- ▶ The mechanism designer can choose a space of allocations \mathcal{A} to obtain desirable outcomes.

E.g. A larger dataset, $\mathcal{A} = \bigcup_{k \geq 0} \mathbb{R}^k$

We can write the space of mechanisms \mathcal{M} as,

A mechanism M receives a dataset from each agent, and returns an *allocation* A_i to each agent.

- ▶ The mechanism designer can choose a space of allocations \mathcal{A} to obtain desirable outcomes.

E.g. A larger dataset, $\mathcal{A} = \bigcup_{k \geq 0} \mathbb{R}^k$

We can write the space of mechanisms \mathcal{M} as,

$$\mathcal{M} = \left\{ M = (\mathcal{A}, b); \quad \mathcal{A} \subset \text{universal set}, \quad b : \left(\bigcup_{n \geq 0} \mathbb{R}^n \right)^m \rightarrow \mathcal{A}^m \right\}$$

A mechanism M receives a dataset from each agent, and returns an *allocation* A_i to each agent.

- ▶ The mechanism designer can choose a space of allocations \mathcal{A} to obtain desirable outcomes.

E.g. A larger dataset, $\mathcal{A} = \bigcup_{k \geq 0} \mathbb{R}^k$

We can write the space of mechanisms \mathcal{M} as,

$$\mathcal{M} = \left\{ M = (\mathcal{A}, b); \quad \mathcal{A} \subset \text{universal set}, \quad b : \left(\bigcup_{n \geq 0} \mathbb{R}^n \right)^m \rightarrow \mathcal{A}^m \right\}$$

Datasets received from the m agents.



After the mechanism is published an agent will

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.
- ▶ Submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\} = f_i(X_i)$.

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.
- ▶ Submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\} = f_i(X_i)$.
 - ▶ f_i maps the dataset collected to possibly altered dataset of a potentially different size (e.g fabrication, withholding, altering etc).

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.
- ▶ Submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\} = f_i(X_i)$.
 - ▶ f_i maps the dataset collected to possibly altered dataset of a potentially different size (e.g fabrication, withholding, altering etc).
- ▶ On receiving her allocation A_i , she will estimate μ via an estimator $h_i(X_i, Y_i, A_i)$.

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.
- ▶ Submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\} = f_i(X_i)$.
 - ▶ f_i maps the dataset collected to possibly altered dataset of a potentially different size (e.g fabrication, withholding, altering etc).
- ▶ On receiving her allocation A_i , she will estimate μ via an estimator $h_i(X_i, Y_i, A_i)$.
 - ▶ An agent need not use the "straightforward" (e.g sample mean) estimator.

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.
- ▶ Submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\} = f_i(X_i)$.
 - ▶ f_i maps the dataset collected to possibly altered dataset of a potentially different size (e.g fabrication, withholding, altering etc).
- ▶ On receiving her allocation A_i , she will estimate μ via an estimator $h_i(X_i, Y_i, A_i)$.
 - ▶ An agent need not use the "straightforward" (e.g sample mean) estimator.

An agent's strategy $s_i = (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$

After the mechanism is published an agent will

- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$.
- ▶ Submit $Y_i = \{y_{i,1}, \dots, y_{i,n'_i}\} = f_i(X_i)$.
 - ▶ f_i maps the dataset collected to possibly altered dataset of a potentially different size (e.g fabrication, withholding, altering etc).
- ▶ On receiving her allocation A_i , she will estimate μ via an estimator $h_i(X_i, Y_i, A_i)$.
 - ▶ An agent need not use the "straightforward" (e.g sample mean) estimator.

An agent's strategy $s_i = (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$

$$\mathcal{F} = \text{submission functions} = \left\{ f : \bigcup_{n \geq 0} \mathbb{R}^n \rightarrow \bigcup_{n \geq 0} \mathbb{R}^n \right\}, \quad \mathcal{H} = \text{estimators} = \left\{ h : \bigcup_{n \geq 0} \mathbb{R}^n \times \bigcup_{n \geq 0} \mathbb{R}^n \times \mathcal{A} \rightarrow \mathbb{R} \right\}$$

An agent's penalty p_i in a mechanism M under a strategy profile $s = (s_1, \dots, s_m)$,

An agent's penalty p_i in a mechanism M under a strategy profile $s = (s_1, \dots, s_m)$,

$$\begin{aligned} p_i(M, s) &= \text{estimation error} && + && \text{data collection cost} \\ &= \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, Y_i, A_i) - \mu)^2 \right] && + && cn_i \end{aligned}$$

An agent's penalty p_i in a mechanism M under a strategy profile $s = (s_1, \dots, s_m)$,

$$\begin{aligned} p_i(M, s) &= \text{estimation error} && + && \text{data collection cost} \\ &= \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, Y_i, A_i) - \mu)^2 \right] && + && cn_i \end{aligned}$$

- ▶ We take a $\sup_{\mu \in \mathbb{R}} \dots$ since μ is unknown. Makes the problem well-defined.

An agent's penalty p_i in a mechanism M under a strategy profile $s = (s_1, \dots, s_m)$,

$$\begin{aligned} p_i(M, s) &= \text{estimation error} && + && \text{data collection cost} \\ &= \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, Y_i, A_i) - \mu)^2 \right] && + && cn_i \end{aligned}$$

- ▶ We take a $\sup_{\mu \in \mathbb{R}} \dots$ since μ is unknown. Makes the problem well-defined.
 - ▶ Otherwise, when $\mu = \mu'$, an agent can choose $n_i = 0$ and $h_i(\cdot, \cdot, \cdot) = \mu'$ to achieve zero penalty.

An agent's penalty p_i in a mechanism M under a strategy profile $s = (s_1, \dots, s_m)$,

$$\begin{aligned}
 p_i(M, s) &= \text{estimation error} && + && \text{data collection cost} \\
 &= \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, Y_i, A_i) - \mu)^2 \right] && + && cn_i
 \end{aligned}$$

- ▶ We take a $\sup_{\mu \in \mathbb{R}} \dots$ since μ is unknown. Makes the problem well-defined.
 - ▶ Otherwise, when $\mu = \mu'$, an agent can choose $n_i = 0$ and $h_i(\cdot, \cdot, \cdot) = \mu'$ to achieve zero penalty.
 - ▶ But this works only if agent knows $\mu = \mu'$ a priori.

A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

Desiderata:

A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

Desiderata:

1. Nash Incentive-compatible (NIC): s^\star is a Nash equilibrium, i.e

$p_i(M, (s_i^\star, s_{-i}^\star)) \leq p_i(M, (s'_i, s_{-i}^\star))$ for all agents i and all other strategies s'_i .

A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

Desiderata:

- 1. Nash Incentive-compatible (NIC):** s^\star is a Nash equilibrium, i.e. $p_i(M, (s_i^\star, s_{-i}^\star)) \leq p_i(M, (s'_i, s_{-i}^\star))$ for all agents i and all other strategies s'_i .
- 2. Individually rational:** An agent's penalty at s^\star is no worse than the lowest penalty she could achieve on her own, i.e. $p_i(M, s^\star) \leq 2\sigma/\sqrt{c}$.

A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

Desiderata:

- 1. Nash Incentive-compatible (NIC):** s^\star is a Nash equilibrium, i.e. $p_i(M, (s_i^\star, s_{-i}^\star)) \leq p_i(M, (s'_i, s_{-i}^\star))$ for all agents i and all other strategies s'_i .
- 2. Individually rational:** An agent's penalty at s^\star is no worse than the lowest penalty she could achieve on her own, i.e. $p_i(M, s^\star) \leq 2\sigma/\sqrt{c}$.
- 3. Approximately efficient:** The social penalty $P(M, s^\star) = \sum_i p_i(M, s^\star)$ is at most a constant factor of the global minimum, i.e.

$$P(M, s^\star) \leq C \cdot \min_{M', s'} p(M', s')$$

A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

Desiderata:

1. **Nash Incentive-compatible (NIC):** s^\star is a Nash equilibrium, i.e. $p_i(M, (s_i^\star, s_{-i}^\star)) \leq p_i(M, (s'_i, s_{-i}^\star))$ for all agents i and all other strategies s'_i .
2. **Individually rational:** An agent's penalty at s^\star is no worse than the lowest penalty she could achieve on her own, i.e. $p_i(M, s^\star) \leq 2\sigma/\sqrt{c}$.
3. **Approximately efficient:** The social penalty $P(M, s^\star) = \sum_i p_i(M, s^\star)$ is at most a constant factor of the global minimum, i.e.

$$P(M, s^\star) \leq C \cdot \min_{M', s'} p(M', s')$$

min without NIC, IR constraints



A mechanism will also publish a *recommended strategy profile* $s^\star = \{s_i^\star\}_{i \in [n]}$.

Desiderata:

- 1. Nash Incentive-compatible (NIC):** s^\star is a Nash equilibrium, i.e. $p_i(M, (s_i^\star, s_{-i}^\star)) \leq p_i(M, (s'_i, s_{-i}^\star))$ for all agents i and all other strategies s'_i .
- 2. Individually rational:** An agent's penalty at s^\star is no worse than the lowest penalty she could achieve on her own, i.e. $p_i(M, s^\star) \leq 2\sigma/\sqrt{c}$.
- 3. Approximately efficient:** The social penalty $P(M, s^\star) = \sum_i p_i(M, s^\star)$ is at most a constant factor of the global minimum, i.e.

$$P(M, s^\star) \leq C \cdot \underbrace{\min_{M', s'} p(M', s')}_{= 2\sigma\sqrt{mc}} \quad \left(\begin{array}{l} \text{min without NIC, IR} \\ \text{constraints} \end{array} \right) \quad (\text{pool-and-share})$$

1. Mechanism design for collaborative normal mean estimation

(Chen, Zhu, Kandasamy, Neurips 2023)

- ▶ Intuitions and Challenges
- ▶ Problem formalism
- ▶ **Mechanism and theoretical analysis**

2. Extensions & Future work

- ▶ Collaborative supervised learning, design of experiments
- ▶ Data marketplaces

Each agent i will

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' submissions Y_{-i} .

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma / \sqrt{cm}$ points from others' submissions Y_{-i} .
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$ # Variance proportional to difference

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' submissions Y_{-i} .
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$ # Variance proportional to difference
 - ▶ $Z_i \leftarrow \left\{ z + \epsilon_z, \text{ for all } z \in Y_{-i} \setminus Z_i, \text{ where } \epsilon_z \sim \mathcal{N}(0, \eta_i^2) \right\}$.

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' submissions Y_{-i} .
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$ # Variance proportional to difference
 - ▶ $Z_i \leftarrow \left\{ z + \epsilon_z, \text{ for all } z \in Y_{-i} \setminus Z_i, \text{ where } \epsilon_z \sim \mathcal{N}(0, \eta_i^2) \right\}$.
 - ▶ Set allocation to each agent, $A_i \leftarrow (Z_i, Z'_i, \eta_i^2)$.

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' submissions Y_{-i} .
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$ # Variance proportional to difference
 - ▶ $Z_i \leftarrow \left\{ z + \epsilon_z, \text{ for all } z \in Y_{-i} \setminus Z_i, \text{ where } \epsilon_z \sim \mathcal{N}(0, \eta_i^2) \right\}$.
 - ▶ Set allocation to each agent, $A_i \leftarrow (Z_i, Z'_i, \eta_i^2)$.

Each agent i will

Each agent i will

- ▶ Choose their strategy $s_i = (n_i, f_i, h_i)$
- ▶ Collect n_i points $X_i = \{x_{i,1}, \dots, x_{i,n_i}\}$ and submit $Y_i = f_i(X_i)$.

Mechanism

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma / \sqrt{cm}$ points from others' submissions Y_{-i} .
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$ # Variance proportional to difference
 - ▶ $Z_i \leftarrow \left\{ z + \epsilon_z, \text{ for all } z \in Y_{-i} \setminus Z_i, \text{ where } \epsilon_z \sim \mathcal{N}(0, \eta_i^2) \right\}$.
 - ▶ Set allocation to each agent, $A_i \leftarrow (Z_i, Z'_i, \eta_i^2)$.

Each agent i will

- ▶ Compute their estimate $h_i(X_i, Y_i, A_i)$

Mechanisms recommends that agents follow $s_i = s_i^\star = (n_i^\star, f_i^\star, h_i^\star)$,

Mechanisms recommends that agents follow $s_i = s_i^\star = (n_i^\star, f_i^\star, h_i^\star)$,

$$n_i^\star = \frac{\sigma}{\sqrt{cm}},$$

$$f_i^\star = \text{identity},$$

$$h_i^\star \left(X_i, Y_i, \underbrace{(Z_i, Z'_i, \eta_i^2)}_{A_i} \right) = \frac{\sum_{u \in X_i \cup Z_i} u + \frac{1}{1 + \eta_i^2 / \sigma^2} \sum_{u \in Z'_i} u}{|X_i \cup Z_i| + \frac{1}{1 + \eta_i^2 / \sigma^2} |Z'_i|}$$

Mechanisms recommends that agents follow $s_i = s_i^\star = (n_i^\star, f_i^\star, h_i^\star)$,

$$n_i^\star = \frac{\sigma}{\sqrt{cm}},$$

$$f_i^\star = \text{identity},$$

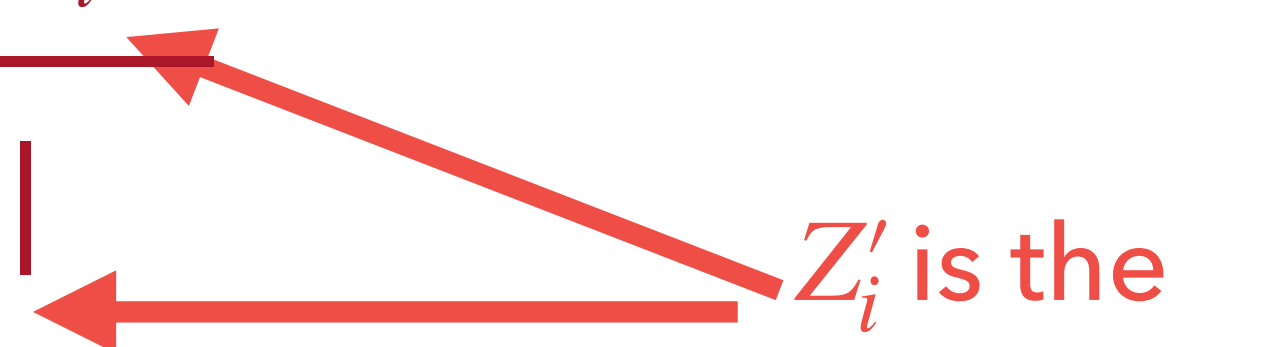
$$h_i^\star \left(X_i, Y_i, \underbrace{(Z_i, Z'_i, \eta_i^2)}_{A_i} \right) = \frac{\sum_{u \in X_i \cup Z_i} u + \frac{1}{1 + \eta_i^2 / \sigma^2} \sum_{u \in Z'_i} u}{|X_i \cup Z_i| + \frac{1}{1 + \eta_i^2 / \sigma^2} |Z'_i|}$$

That is collect a sufficient amount of data n_i^\star , submit it truthfully f_i^\star , and use a weighted average estimator h_i^\star .

Mechanisms recommends that agents follow $s_i = s_i^\star = (n_i^\star, f_i^\star, h_i^\star)$,

$$n_i^\star = \frac{\sigma}{\sqrt{cm}},$$

$$f_i^\star = \text{identity},$$

$$h_i^\star \left(X_i, Y_i, \underbrace{(Z_i, Z'_i, \eta_i^2)}_{A_i} \right) = \frac{\sum_{u \in X_i \cup Z_i} u + \frac{1}{1 + \eta_i^2 / \sigma^2} \sum_{u \in Z'_i} u}{|X_i \cup Z_i| + \frac{1}{1 + \eta_i^2 / \sigma^2} |Z'_i|}$$


Z'_i is the corrupted dataset.

That is collect a sufficient amount of data n_i^\star , submit it truthfully f_i^\star , and use a weighted average estimator h_i^\star .

Theorem: At the recommended strategy profile s^\star , the mechanism is Nash incentive-compatible, individually rational, and approximately efficient with $P(M, s^\star) \leq 2 \cdot \inf_{M,s} P(M, s)$.

Theorem: At the recommended strategy profile s^\star , the mechanism is Nash incentive-compatible, individually rational, and approximately efficient with $P(M, s^\star) \leq 2 \cdot \inf_{M,s} P(M, s)$.

Theorem (High-dimensional distributions with bounded variance): The recommended strategy profile s^\star is an $\mathcal{O}(1/m)$ -approximate Nash equilibrium. Moreover, the mechanism is approximately efficient with $P(M, s^\star) \leq (2 + \mathcal{O}(1/m)) \cdot \inf_{M,s} P(M, s)$.

We need to show that $s^{\star} = \{(n_i^{\star}, f_i^{\star}, h_i^{\star})\}_i$ is a Nash equilibrium, i.e

$$p_i(M, (s_i^{\star}, s_{-i}^{\star})) \leq p_i(M, (s_i, s_{-i}^{\star})) \quad \text{for all agents } i \text{ and all deviations } s_i$$

We need to show that $s^\star = \{(n_i^\star, f_i^\star, h_i^\star)\}_i$ is a Nash equilibrium, i.e

$$p_i(M, (s_i^\star, s_{-i}^\star)) \leq p_i(M, (s_i, s_{-i}^\star)) \quad \text{for all agents } i \text{ and all deviations } s_i$$

Step 1: First, we will show that for any amount of data collected n_i , submitting it truthfully and using the recommended estimator minimizes the penalty, i.e

$$p_i\left(M, \left((n_i, f_i^\star, h_i^\star), s_{-i}^\star\right)\right) \leq p_i\left(M, \left((n_i, f_i, h_i), s_{-i}^\star\right)\right) \quad \text{for all } (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$$

We need to show that $s^{\star} = \{(n_i^{\star}, f_i^{\star}, h_i^{\star})\}_i$ is a Nash equilibrium, i.e

$$p_i(M, (s_i^{\star}, s_{-i}^{\star})) \leq p_i(M, (s_i, s_{-i}^{\star})) \quad \text{for all agents } i \text{ and all deviations } s_i$$

Step 1: First, we will show that for any amount of data collected n_i , submitting it truthfully and using the recommended estimator minimizes the penalty, i.e

$$p_i\left(M, ((n_i, f_i^{\star}, h_i^{\star}), s_{-i}^{\star})\right) \leq p_i\left(M, ((n_i, f_i, h_i), s_{-i}^{\star})\right) \quad \text{for all } (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$$

Step 2: Then, we will show the agent's penalty is minimized when she collects n_i^{\star} samples under $(f_i^{\star}, h_i^{\star})$, i.e

$$p_i\left(M, ((n_i^{\star}, f_i^{\star}, h_i^{\star}), s_{-i}^{\star})\right) \leq p_i\left(M, ((n_i, f_i, h_i), s_{-i}^{\star})\right) \quad \text{for all } n_i \in \mathbb{N}$$

Step 1: First, we will show that for any amount of data collected n_i , submitting it truthfully and using the recommended estimator minimizes the penalty, i.e

$$p_i \left(M, \left((n_i, f_i^*, h_i^*), s_{-i}^* \right) \right) \leq p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^* \right) \right) \quad \text{for all } (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$$

Step 1: First, we will show that for any amount of data collected n_i , submitting it truthfully and using the recommended estimator minimizes the penalty, i.e

$$p_i \left(M, \left((n_i, f_i^*, h_i^*), s_{-i}^* \right) \right) \leq p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^* \right) \right) \quad \text{for all } (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$$

We need to show, for all $(n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$,

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[\left(h_i^* (X_i, f_i^*(X_i), A_i) - \mu \right)^2 \right] + cn_i \leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[\left(h_i (X_i, f_i(X_i), A_i) - \mu \right)^2 \right] + cn_i$$

Step 1: First, we will show that for any amount of data collected n_i , submitting it truthfully and using the recommended estimator minimizes the penalty, i.e

$$p_i \left(M, ((n_i, f_i^*, h_i^*), s_{-i}^*) \right) \leq p_i \left(M, ((n_i, f_i, h_i), s_{-i}^*) \right) \quad \text{for all } (n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$$

We need to show, for all $(n_i, f_i, h_i) \in \mathbb{N} \times \mathcal{F} \times \mathcal{H}$,

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^* (X_i, f_i^*(X_i), A_i) - \mu)^2 \right] + cn_i \leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i (X_i, f_i(X_i), A_i) - \mu)^2 \right] + cn_i$$

Or equivalently,

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^* (X_i, f_i^*(X_i), A_i) - \mu)^2 \right] = \inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i (X_i, f_i(X_i), A_i) - \mu)^2 \right]$$

We are given $X_1^n = \{X_1, \dots, X_n\}$, drawn from $\mathcal{N}(\mu, \sigma^2)$ where σ^2 is known. Let $\hat{\mu}(X)$ be an estimator for μ . We wish to show

We are given $X_1^n = \{X_1, \dots, X_n\}$, drawn from $\mathcal{N}(\mu, \sigma^2)$ where σ^2 is known. Let $\hat{\mu}(X)$ be an estimator for μ . We wish to show

$$\text{minimax risk} = \inf_{\hat{\mu}} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} \left[(\mu - \hat{\mu}(X_1^n))^2 \right] = \frac{\sigma^2}{n}$$

We are given $X_1^n = \{X_1, \dots, X_n\}$, drawn from $\mathcal{N}(\mu, \sigma^2)$ where σ^2 is known. Let $\hat{\mu}(X)$ be an estimator for μ . We wish to show

$$\text{minimax risk} = \inf_{\hat{\mu}} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} \left[(\mu - \hat{\mu}(X_1^n))^2 \right] = \frac{\sigma^2}{n}$$

Upper bound via an estimator: We can use the sample mean

$$\hat{\mu}_{\text{sm}}(X) = (X_1 + \dots + X_n)/n.$$

We are given $X_1^n = \{X_1, \dots, X_n\}$, drawn from $\mathcal{N}(\mu, \sigma^2)$ where σ^2 is known. Let $\hat{\mu}(X)$ be an estimator for μ . We wish to show

$$\text{minimax risk} = \inf_{\hat{\mu}} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} \left[(\mu - \hat{\mu}(X_1^n))^2 \right] = \frac{\sigma^2}{n}$$

Upper bound via an estimator: We can use the sample mean

$$\hat{\mu}_{\text{sm}}(X) = (X_1 + \dots + X_n)/n.$$

$$\text{minimax risk} \leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} \left[(\mu - \hat{\mu}_{\text{sm}}(X_1^n))^2 \right] = \frac{\sigma^2}{n}$$

Lower bound via Bayes' risk: Choose a prior Λ for μ . Then lower bound via the Bayes' risk under Λ .

Lower bound via Bayes' risk: Choose a prior Λ for μ . Then lower bound via the Bayes' risk under Λ .

We will use $\Lambda = \mathcal{N}(0, \tau^2)$. For any estimator $\hat{\mu}$,

Lower bound via Bayes' risk: Choose a prior Λ for μ . Then lower bound via the Bayes' risk under Λ .

We will use $\Lambda = \mathcal{N}(0, \tau^2)$. For any estimator $\hat{\mu}$,

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2] \geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2 \mid \mu] \right] \longleftarrow \text{sup} \geq \text{avg}$$

Lower bound via Bayes' risk: Choose a prior Λ for μ . Then lower bound via the Bayes' risk under Λ .

We will use $\Lambda = \mathcal{N}(0, \tau^2)$. For any estimator $\hat{\mu}$,

$$\begin{aligned} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2] &\geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2 \mid \mu] \right] && \leftarrow \text{sup} \geq \text{avg} \\ &= \mathbb{E}_{X_1^n} \left[\mathbb{E}_{\mu \sim \Lambda} [(\mu - \hat{\mu}(X_1^n))^2 \mid X_1^n] \right] && \leftarrow \text{Swap order of expectation} \end{aligned}$$

Lower bound via Bayes' risk: Choose a prior Λ for μ . Then lower bound via the Bayes' risk under Λ .

We will use $\Lambda = \mathcal{N}(0, \tau^2)$. For any estimator $\hat{\mu}$,

$$\begin{aligned}
 \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2] &\geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2 \mid \mu] \right] && \leftarrow \text{sup} \geq \text{avg} \\
 &= \mathbb{E}_{X_1^n} \left[\mathbb{E}_{\mu \sim \Lambda} [(\mu - \hat{\mu}(X_1^n))^2 \mid X_1^n] \right] && \leftarrow \text{Swap order of expectation} \\
 &\geq \mathbb{E}_{X_1^n} \left[\frac{\sigma^2}{n + \sigma^2/\tau^2} \right] && \leftarrow \text{By normal conjugacy and choosing } \hat{\mu} = \text{posterior mean}
 \end{aligned}$$

Lower bound via Bayes' risk: Choose a prior Λ for μ . Then lower bound via the Bayes' risk under Λ .

We will use $\Lambda = \mathcal{N}(0, \tau^2)$. For any estimator $\hat{\mu}$,

$$\begin{aligned}
 \sup_{\mu \in \mathbb{R}} \mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2] &\geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{X_1^n} [(\mu - \hat{\mu}(X_1^n))^2 \mid \mu] \right] && \leftarrow \text{sup} \geq \text{avg} \\
 &= \mathbb{E}_{X_1^n} \left[\mathbb{E}_{\mu \sim \Lambda} [(\mu - \hat{\mu}(X_1^n))^2 \mid X_1^n] \right] && \leftarrow \text{Swap order of expectation} \\
 &\geq \mathbb{E}_{X_1^n} \left[\frac{\sigma^2}{n + \sigma^2/\tau^2} \right] && \leftarrow \text{By normal conjugacy and choosing } \hat{\mu} = \text{posterior mean} \\
 &= \frac{\sigma^2}{n + \sigma^2/\tau^2} \rightarrow \frac{\sigma^2}{n} && \text{as } \tau^2 \rightarrow \infty
 \end{aligned}$$

We will apply the same recipe to show that f_i^\star, h_i^\star are the minimax-optimal submission functions and estimators for the agent.

$$\inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_\mu \left[\left(h_i (X_i, f_i(X_i), A_i) - \mu \right)^2 \right] = \sup_{\mu \in \mathbb{R}} \mathbb{E}_\mu \left[\left(h_i^\star (X_i, f_i^\star(X_i), A_i) - \mu \right)^2 \right]$$

We will apply the same recipe to show that f_i^\star, h_i^\star are the minimax-optimal submission functions and estimators for the agent.

$$\inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_\mu \left[\left(h_i(X_i, f_i(X_i), A_i) - \mu \right)^2 \right] = \sup_{\mu \in \mathbb{R}} \mathbb{E}_\mu \left[\left(h_i^\star(X_i, f_i^\star(X_i), A_i) - \mu \right)^2 \right]$$

But the data available to the agent is not i.i.d!

- ▶ The corruption is data-dependent.

We will apply the same recipe to show that f_i^\star, h_i^\star are the minimax-optimal submission functions and estimators for the agent.

$$\inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_\mu \left[\left(h_i(X_i, f_i(X_i), A_i) - \mu \right)^2 \right] = \sup_{\mu \in \mathbb{R}} \mathbb{E}_\mu \left[\left(h_i^\star(X_i, f_i^\star(X_i), A_i) - \mu \right)^2 \right]$$

But the data available to the agent is not i.i.d!

- ▶ The corruption is data-dependent.
- ▶ In fact, X_i, Z_i, Z'_i is not even jointly Gaussian.

We show

$$\inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] \leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^*(X_i, f_i^*(X_i), A_i) - \mu)^2 \right]$$

We show

$$\begin{aligned} \inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] &\leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^*(X_i, f_i^*(X_i), A_i) - \mu)^2 \right] \\ &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\left(\frac{(m-2)n_i^*}{\left(\sigma^2 + \alpha^2 (\sigma^2/n_i + \sigma^2/n_i^*) Z^2 \right)} + \frac{n_i + n_i^*}{\sigma^{-2}} \right)^{-1} \right] \end{aligned}$$

We show

$$\begin{aligned}
 \inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] &\leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^*(X_i, f_i^*(X_i), A_i) - \mu)^2 \right] \\
 &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\left(\frac{(m-2)n_i^*}{\left(\sigma^2 + \alpha^2 (\sigma^2/n_i + \sigma^2/n_i^*) Z^2 \right)} + \frac{n_i + n_i^*}{\sigma^{-2}} \right)^{-1} \right] =: R_{\infty}(n_i) \quad (\text{say})
 \end{aligned}$$

We show

$$\begin{aligned} \inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] &\leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^{\star}(X_i, f_i^{\star}(X_i), A_i) - \mu)^2 \right] \\ &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\left(\frac{(m-2)n_i^{\star}}{\left(\sigma^2 + \alpha^2 (\sigma^2/n_i + \sigma^2/n_i^{\star}) Z^2 \right)} + \frac{n_i + n_i^{\star}}{\sigma^{-2}} \right)^{-1} \right] =: R_{\infty}(n_i) \quad (\text{say}) \end{aligned}$$

Key ingredients

- ▶ When $f_i^{\star} = \text{identity}$, first condition on X_i, Z_i , then $Z_i' \sim \mathcal{N}(0, \sigma^2 + \eta^2)$.

We show

$$\begin{aligned} \inf_{f_i, h_i} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] &\leq \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\mu} \left[(h_i^*(X_i, f_i^*(X_i), A_i) - \mu)^2 \right] \\ &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\left(\frac{(m-2)n_i^*}{\left(\sigma^2 + \alpha^2 (\sigma^2/n_i + \sigma^2/n_i^*) Z^2 \right)} + \frac{n_i + n_i^*}{\sigma^{-2}} \right)^{-1} \right] =: R_{\infty}(n_i) \quad (\text{say}) \end{aligned}$$

Key ingredients

- ▶ When $f_i^* = \text{identity}$, first condition on X_i, Z_i , then $Z_i' \sim \mathcal{N}(0, \sigma^2 + \eta^2)$.
- ▶ Properties of Gaussians
- ▶ Lots of algebra

Choose prior $\Lambda = \mathcal{N}(0, \tau^2)$ for μ . Then for any f_i, h_i , we have

Choose prior $\Lambda = \mathcal{N}(0, \tau^2)$ for μ . Then for any f_i, h_i , we have

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] \geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \mu \right] \right] \longleftarrow \text{sup} \geq \text{avg}$$

Choose prior $\Lambda = \mathcal{N}(0, \tau^2)$ for μ . Then for any f_i, h_i , we have

$$\begin{aligned} \sup_{\mu \in \mathbb{R}} \mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] &\geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \mu \right] \right] \longleftarrow \text{sup} \geq \text{avg} \\ &= \mathbb{E}_{\text{data} \sim \mu} \left[\mathbb{E}_{\mu \sim \Lambda} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \text{data} \right] \right] \longleftarrow \text{Swap order of expectation} \end{aligned}$$

Choose prior $\Lambda = \mathcal{N}(0, \tau^2)$ for μ . Then for any f_i, h_i , we have

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] \geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \mu \right] \right] \leftarrow \text{sup} \geq \text{avg}$$

$$= \mathbb{E}_{\text{data} \sim \mu} \left[\mathbb{E}_{\mu \sim \Lambda} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \text{data} \right] \right] \leftarrow \text{Swap order of expectation}$$

$$\geq \mathbb{E}_{\text{data}} \left[\inf_{f_i} \left(|Z'_i| \left(\sigma^2 + \alpha^2 \left(\frac{1}{|f_i(X_i)|} \sum_{y \in f_i(X_i)} y - \frac{1}{|Z_i|} \sum_{z \in Z_i} z \right)^2 \right)^{-1} + \frac{|X_i| + |Z_i|}{\sigma^2} + \frac{1}{\tau^2} \right)^{-1} \right]$$

(X_i, Z_i, Z'_i, μ) is not jointly Gaussian, but $Z'_i, \mu \mid X_i, Z_i$ is Gaussian.

Minimized by choosing $h_i =$ posterior mean.

Choose prior $\Lambda = \mathcal{N}(0, \tau^2)$ for μ . Then for any f_i, h_i , we have

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] \geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \mu \right] \right] \leftarrow \text{sup} \geq \text{avg}$$

$$= \mathbb{E}_{\text{data} \sim \mu} \left[\mathbb{E}_{\mu \sim \Lambda} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \text{data} \right] \right] \leftarrow \text{Swap order of expectation}$$

$$\geq \mathbb{E}_{\text{data}} \left[\inf_{f_i} \left(|Z'_i| \left(\sigma^2 + \alpha^2 \left(\frac{1}{|f_i(X_i)|} \sum_{y \in f_i(X_i)} y - \frac{1}{|Z_i|} \sum_{z \in Z_i} z \right)^2 \right)^{-1} + \frac{|X_i| + |Z_i|}{\sigma^2} + \frac{1}{\tau^2} \right)^{-1} \right]$$

(X_i, Z_i, Z'_i, μ) is not jointly Gaussian, but $Z'_i, \mu \mid X_i, Z_i$ is Gaussian.

Minimized by choosing $h_i =$ posterior mean.

$$= \dots = R_\tau(n_i) \quad (\text{say})$$

Minimized by applying the Hardy-Littlewood inequality and choosing $f_i(X_i) = \left\{ \left(1 + \sigma^2 / (|X| \ell^2) \right)^{-1} x, \forall x \in X_i \right\}$.

Choose prior $\Lambda = \mathcal{N}(0, \tau^2)$ for μ . Then for any f_i, h_i , we have

$$\sup_{\mu \in \mathbb{R}} \mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \right] \geq \mathbb{E}_{\mu \sim \Lambda} \left[\mathbb{E}_{\text{data} \sim \mu} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \mu \right] \right] \leftarrow \text{sup} \geq \text{avg}$$

$$= \mathbb{E}_{\text{data} \sim \mu} \left[\mathbb{E}_{\mu \sim \Lambda} \left[(h_i(X_i, f_i(X_i), A_i) - \mu)^2 \mid \text{data} \right] \right] \leftarrow \text{Swap order of expectation}$$

$$\geq \mathbb{E}_{\text{data}} \left[\inf_{f_i} \left(|Z'_i| \left(\sigma^2 + \alpha^2 \left(\frac{1}{|f_i(X_i)|} \sum_{y \in f_i(X_i)} y - \frac{1}{|Z_i|} \sum_{z \in Z_i} z \right)^2 \right)^{-1} + \frac{|X_i| + |Z_i|}{\sigma^2} + \frac{1}{\tau^2} \right)^{-1} \right]$$

(X_i, Z_i, Z'_i, μ) is not jointly Gaussian, but $Z'_i, \mu \mid X_i, Z_i$ is Gaussian.

Minimized by choosing $h_i =$ posterior mean.

$$= \dots = R_\tau(n_i) \quad (\text{say}) \leftarrow \text{Minimized by applying the Hardy-Littlewood inequality and choosing } f_i(X_i) = \left\{ \left(1 + \sigma^2 / (|X| \ell^2) \right)^{-1} x, \forall x \in X_i \right\}.$$

$$\rightarrow R_\infty(n_i)$$

Step 2: Then, we will show the agent's penalty is minimized when she collects n_i samples under (f_i^\star, h_i^\star) , i.e

$$p_i \left(M, \left((n_i^\star, f_i^\star, h_i^\star), s_{-i}^\star \right) \right) \leq p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right) \quad \text{for all } n_i \in \mathbb{N}$$

Step 2: Then, we will show the agent's penalty is minimized when she collects n_i samples under (f_i^\star, h_i^\star) , i.e

$$p_i \left(M, \left((n_i^\star, f_i^\star, h_i^\star), s_{-i}^\star \right) \right) \leq p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right) \quad \text{for all } n_i \in \mathbb{N}$$

From Step 1 we have,

$$q_i(n_i) := \inf_{f_i, h_i} p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right)$$

Step 2: Then, we will show the agent's penalty is minimized when she collects n_i samples under (f_i^\star, h_i^\star) , i.e

$$p_i \left(M, \left((n_i^\star, f_i^\star, h_i^\star), s_{-i}^\star \right) \right) \leq p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right) \quad \text{for all } n_i \in \mathbb{N}$$

From Step 1 we have,

$$q_i(n_i) := \inf_{f_i, h_i} p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right) = \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\left(\frac{(m-2)n_i^\star}{\left(\sigma^2 + \alpha^2 \left(\sigma^2/n_i + \sigma^2/n_i^\star \right) Z^2 \right)} + \frac{n_i + n_i^\star}{\sigma^{-2}} \right)^{-1} \right] + cn_i$$

Step 2: Then, we will show the agent's penalty is minimized when she collects n_i samples under (f_i^\star, h_i^\star) , i.e

$$p_i \left(M, \left((n_i^\star, f_i^\star, h_i^\star), s_{-i}^\star \right) \right) \leq p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right) \quad \text{for all } n_i \in \mathbb{N}$$

From Step 1 we have,

$$q_i(n_i) := \inf_{f_i, h_i} p_i \left(M, \left((n_i, f_i, h_i), s_{-i}^\star \right) \right) = \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\left(\frac{(m-2)n_i^\star}{\left(\sigma^2 + \alpha^2 \left(\sigma^2/n_i + \sigma^2/n_i^\star \right) Z^2 \right)} + \frac{n_i + n_i^\star}{\sigma^{-2}} \right)^{-1} \right] + cn_i$$

- The term inside \mathbb{E} is convex. Hence so is $q(n_i)$.
- $q(n_i)$ is minimized at $n_i = n_i^\star$ (by our choice of α).

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' subn
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' subn
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$

We set α to be the smallest number larger than $\sqrt{n_i^*}$ such that $G(\alpha) = 0$, where,

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' subn
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$

We set α to be the smallest number larger than $\sqrt{n_i^*}$ such that $G(\alpha) = 0$, where,

$$G(\alpha) := \left(\frac{m-4}{m-2} \frac{4\alpha^2}{\sigma/\sqrt{cm}} - 1 \right) \frac{4\alpha}{\sqrt{\sigma}(m/c)^{1/4}} - \left(4(m+1) \frac{\alpha^2}{\sigma\sqrt{m/c}} - 1 \right) \sqrt{2\pi} \exp\left(\frac{\sigma\sqrt{m/c}}{8\alpha^2}\right) \text{Erfc}\left(\frac{\sqrt{\sigma}(m/c)^{1/4}}{2\sqrt{2}\alpha}\right)$$

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' subn
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$

We set α to be the smallest number larger than $\sqrt{n_i^*}$ such that $G(\alpha) = 0$, where,

$$G(\alpha) := \left(\frac{m-4}{m-2} \frac{4\alpha^2}{\sigma/\sqrt{cm}} - 1 \right) \frac{4\alpha}{\sqrt{\sigma}(m/c)^{1/4}} - \left(4(m+1) \frac{\alpha^2}{\sigma\sqrt{m/c}} - 1 \right) \sqrt{2\pi} \exp\left(\frac{\sigma\sqrt{m/c}}{8\alpha^2}\right) \text{Erfc}\left(\frac{\sqrt{\sigma}(m/c)^{1/4}}{2\sqrt{2}\alpha}\right)$$

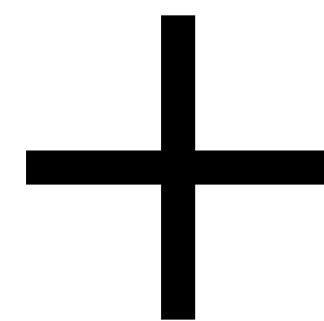
- ▶ $\alpha^2 \geq n_i^*$: step 1 of NIC (sufficiently penalise untruthful agents).
- ▶ "smallest number larger than": for efficiency (don't over-penalize truthful agents).
- ▶ $G(\alpha) = 0$: step 2 of NIC (collect a sufficient amount of data).

- ▶ For each agent i :
 - ▶ $Z_i \leftarrow$ sample $n^* = \sigma/\sqrt{cm}$ points from others' subn
 - ▶ Set noise variance $\eta_i^2 = \alpha^2 (\text{mean}(Y_i) - \text{mean}(Z_i))^2$

We set α to be the smallest number larger than $\sqrt{n_i^*}$ such that $G(\alpha) = 0$, where,

$$G(\alpha) := \left(\frac{m-4}{m-2} \frac{4\alpha^2}{\sigma/\sqrt{cm}} - 1 \right) \frac{4\alpha}{\sqrt{\sigma}(m/c)^{1/4}} - \left(4(m+1) \frac{\alpha^2}{\sigma\sqrt{m/c}} - 1 \right) \sqrt{2\pi} \exp\left(\frac{\sigma\sqrt{m/c}}{8\alpha^2}\right) \text{Erfc}\left(\frac{\sqrt{\sigma}(m/c)^{1/4}}{2\sqrt{2}\alpha}\right)$$

- ▶ $\alpha^2 \geq n_i^*$: step 1 of NIC (sufficiently penalise untruthful agents).
- ▶ "smallest number larger than": for efficiency (don't over-penalize truthful agents).
- ▶ $G(\alpha) = 0$: step 2 of NIC (collect a sufficient amount of data).



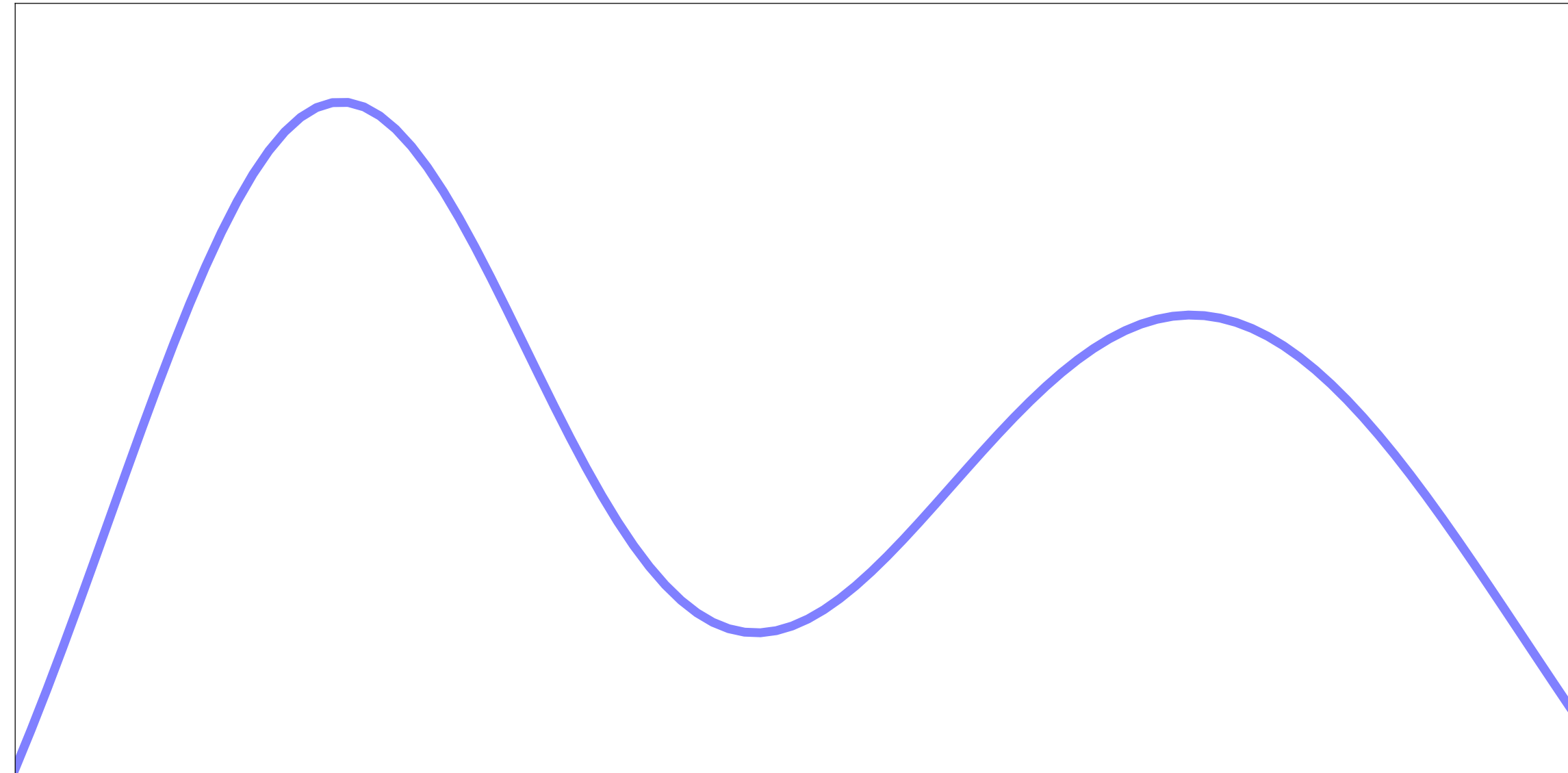
1. Mechanism design for collaborative normal mean estimation

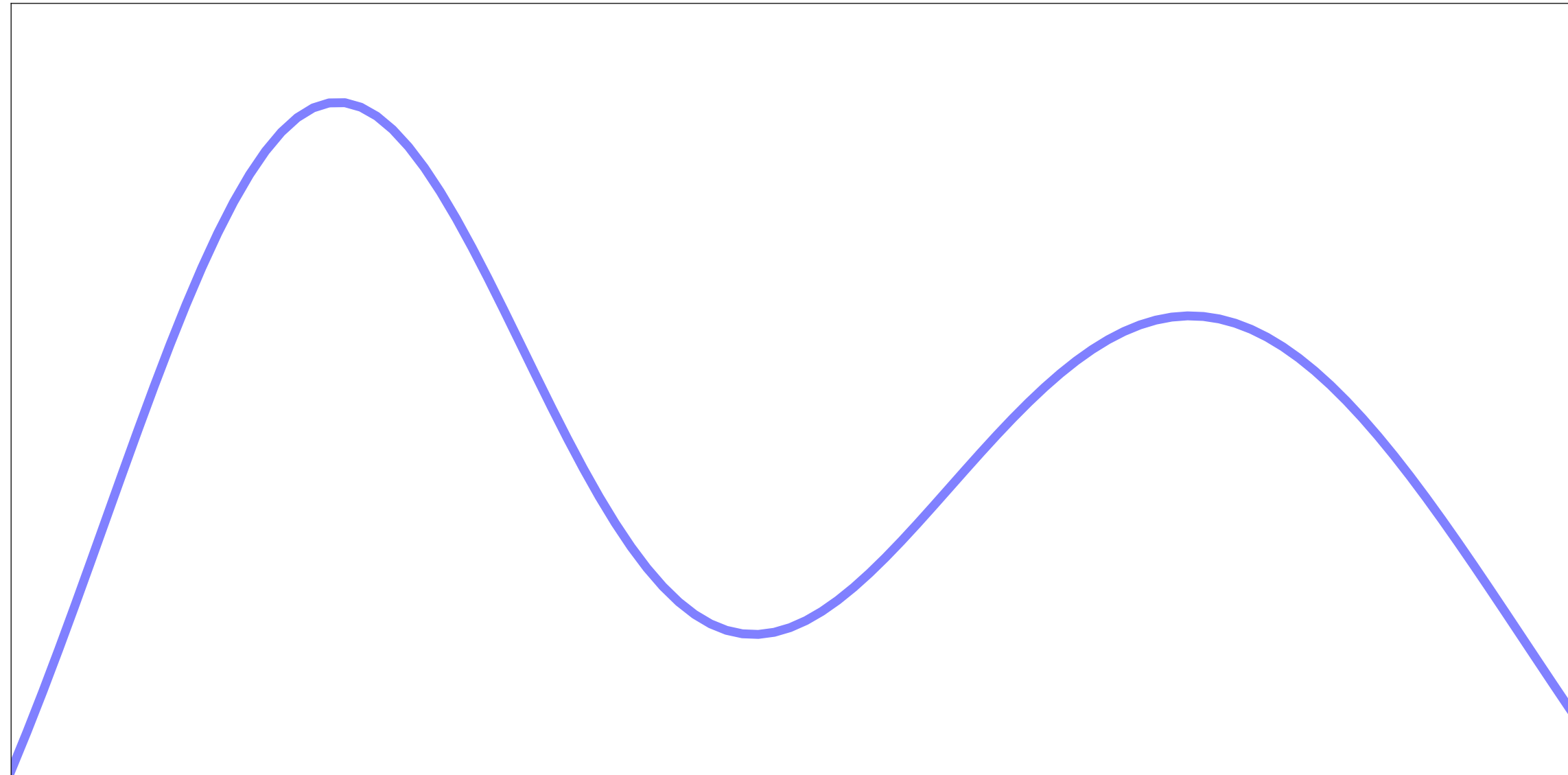
(Chen, Zhu, Kandasamy, *Neurips 2023*)

- ▶ Intuitions and Challenges
- ▶ Problem formalism
- ▶ Mechanism and theoretical analysis

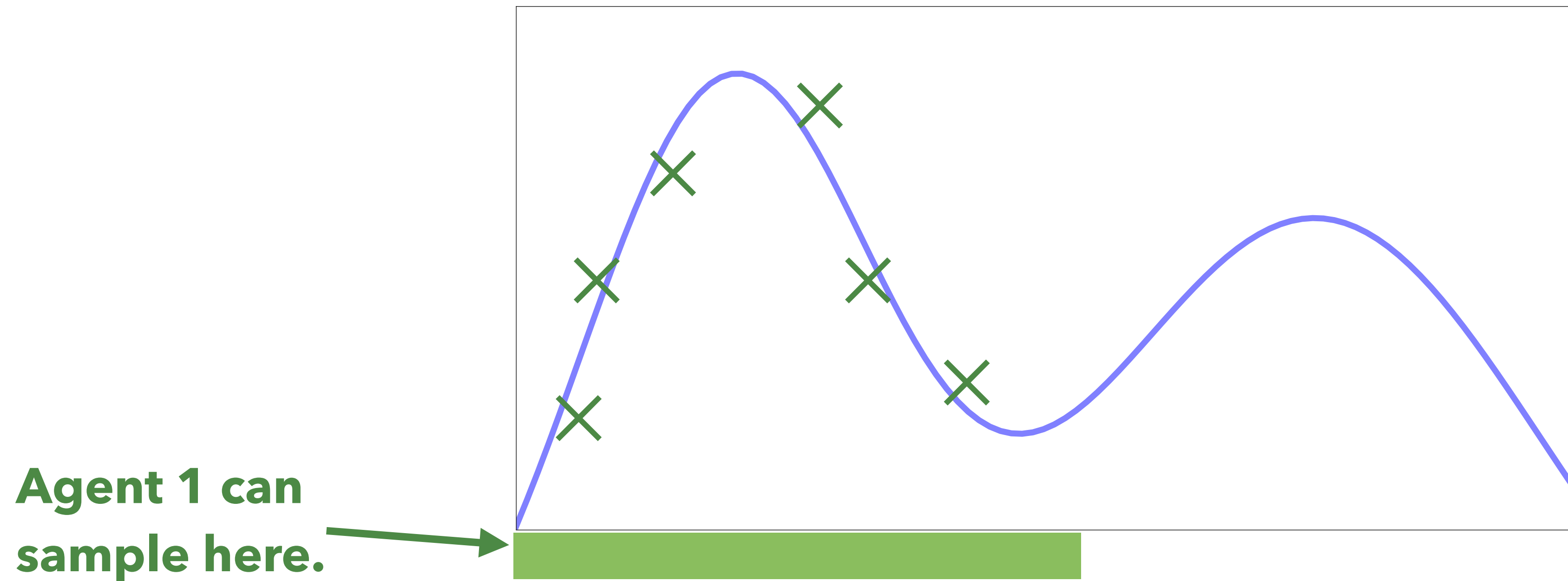
2. Extensions & Future work

- ▶ **Collaborative supervised learning, design of experiments**
- ▶ **Data marketplaces**

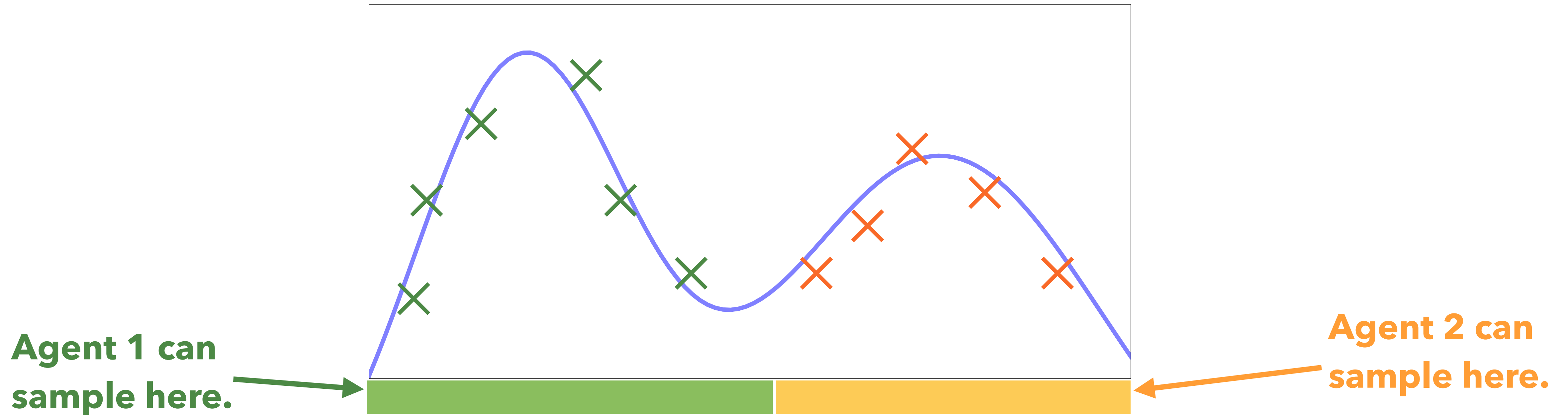




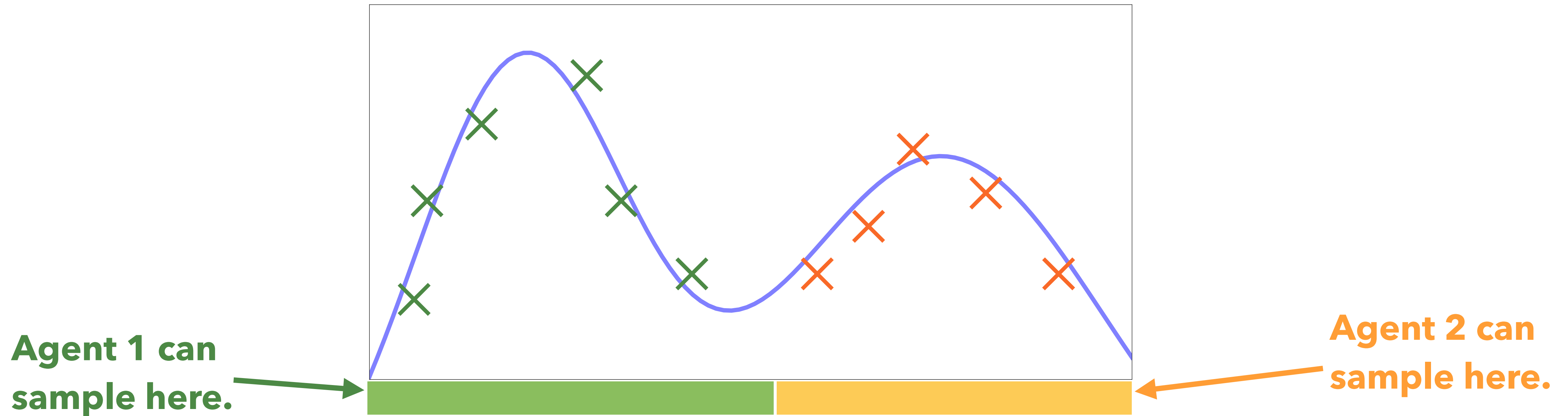
Data sharing when there is asymmetric data collection capabilities?



Data sharing when there is asymmetric data collection capabilities?

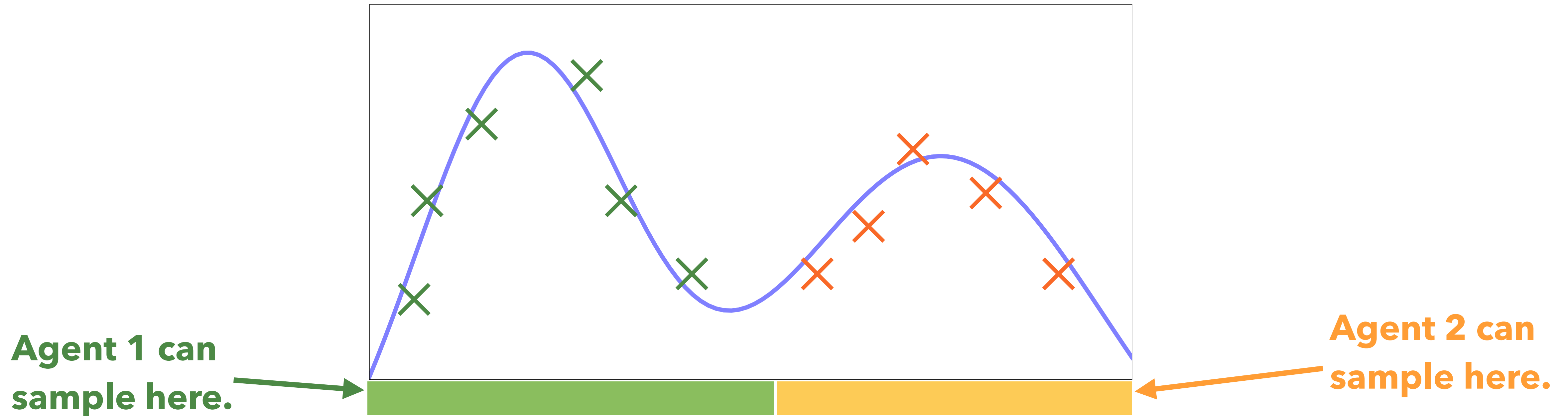


Data sharing when there is asymmetric data collection capabilities?



Data sharing when there is asymmetric data collection capabilities?

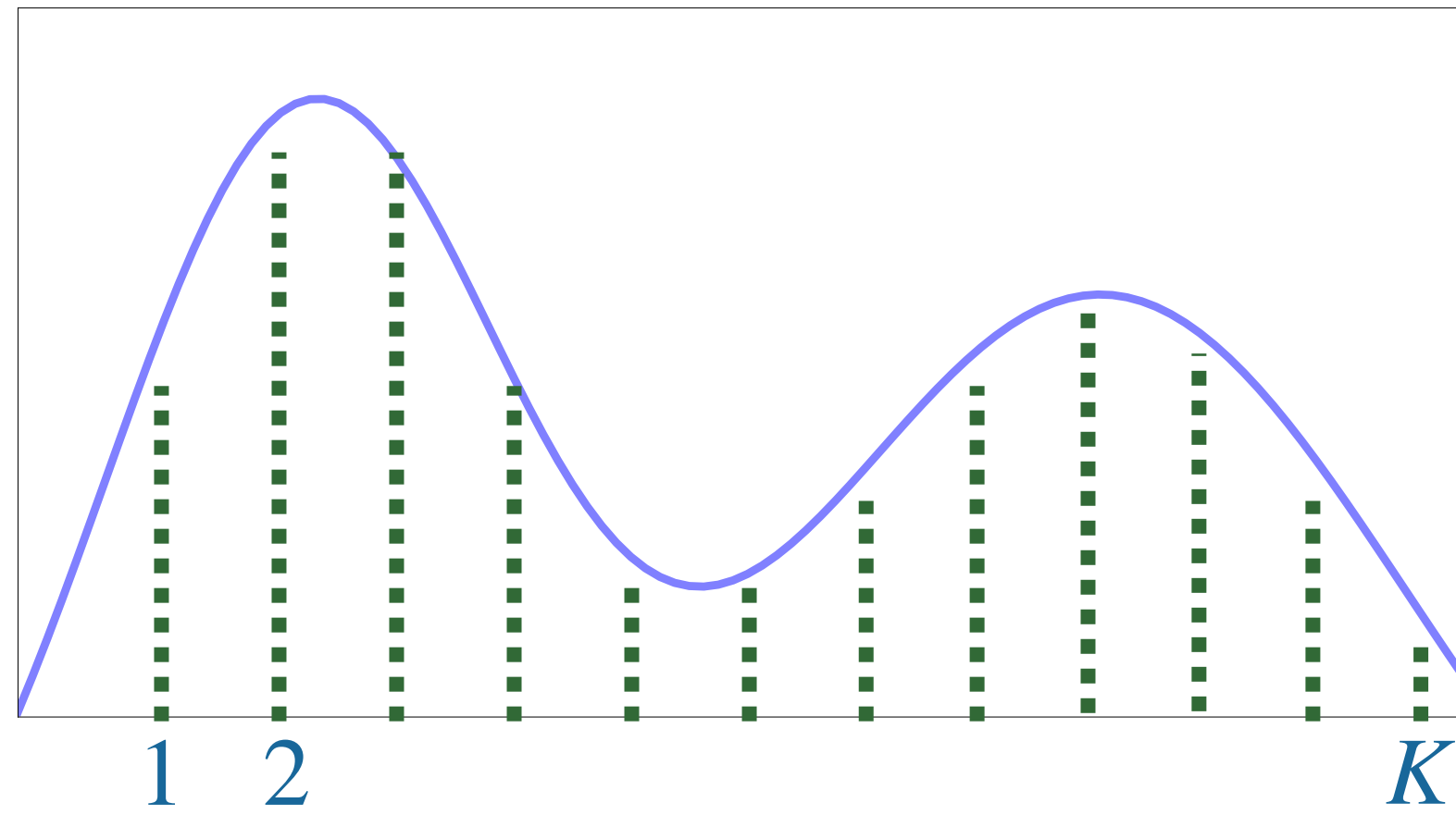
- + Agents will be more willing to collaborate due to complementarity of data.



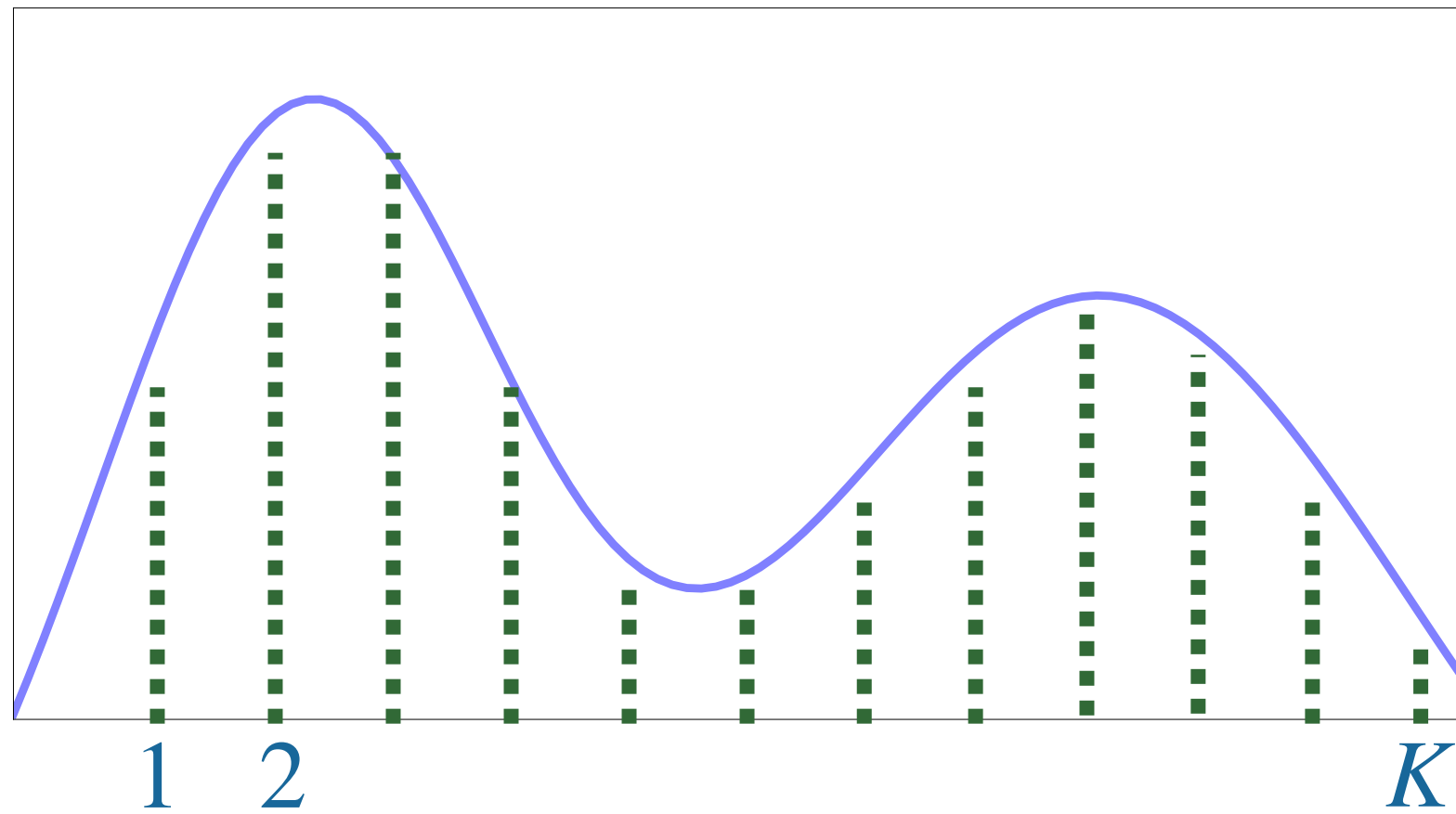
Data sharing when there is asymmetric data collection capabilities?

- + Agents will be more willing to collaborate due to complementarity of data.
- No way to validate an agent's data with other similar data.

Consider a K discretisation of the domain



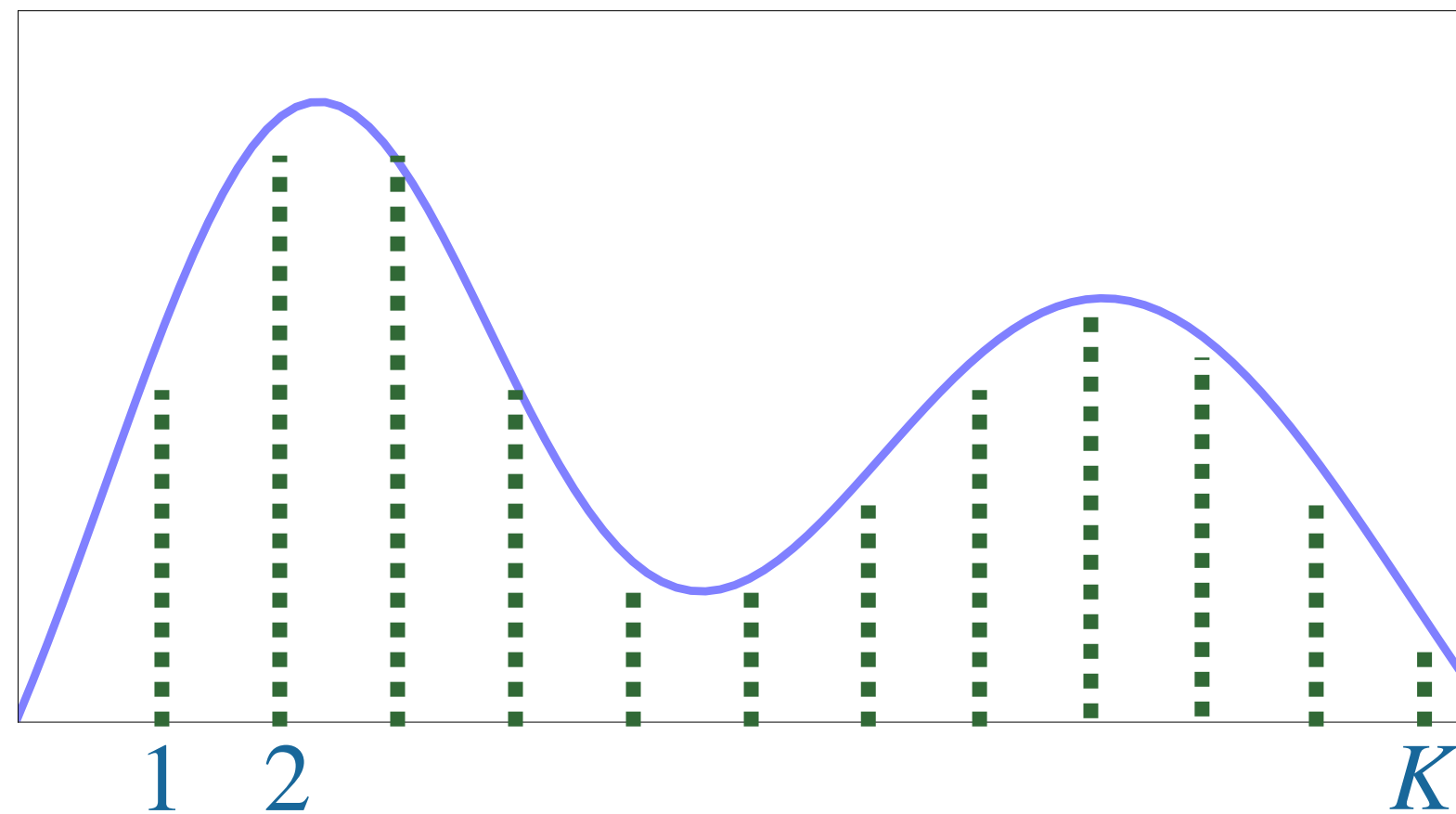
Consider a K discretisation of the domain



Agent i can sample from distribution k at cost $c_{i,k}$.

$$\text{Penalty, } p_i = \sum_{k=1}^K \text{est-err}_k + \sum_{k=1}^K c_{i,k} n_{i,k}$$

Consider a K discretisation of the domain



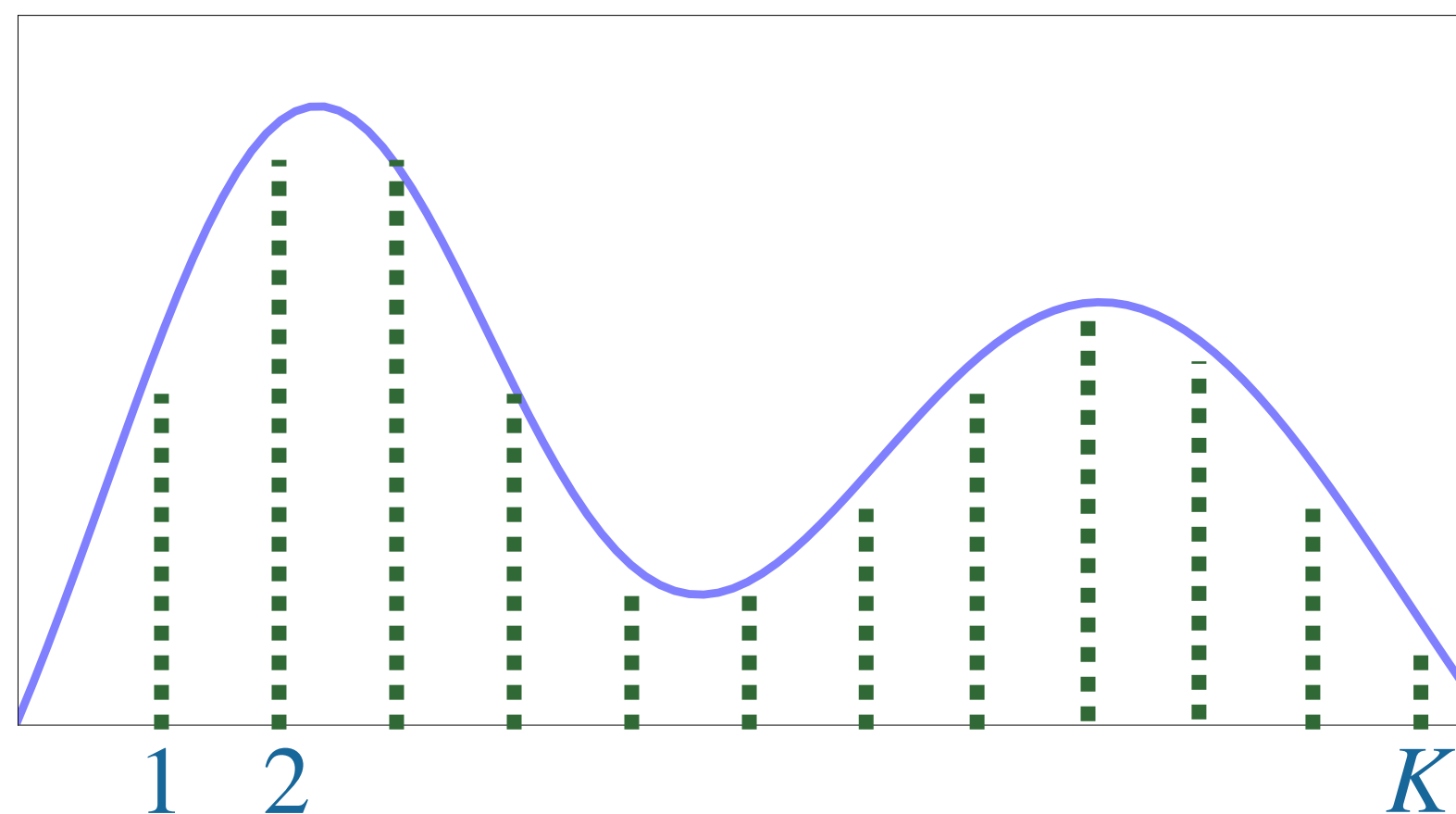
Agent i can sample from distribution k at cost $c_{i,k}$.

$$\text{Penalty, } p_i = \sum_{k=1}^K \text{est-err}_k + \sum_{k=1}^K c_{i,k} n_{i,k}$$

Theorem: There exists a NIC and IR mechanism for which,

$$P(M, s^*) \leq 8\sqrt{m} \cdot \inf_{M,s} P(M, s)$$

Consider a K discretisation of the domain



Agent i can sample from distribution k at cost $c_{i,k}$.

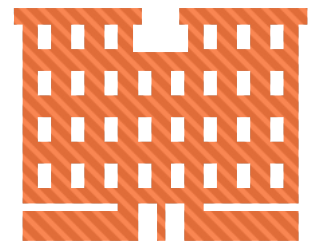
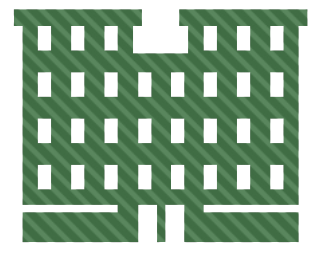
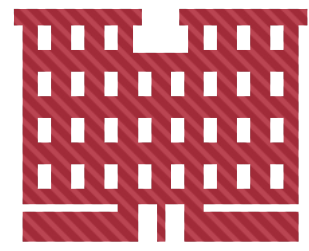
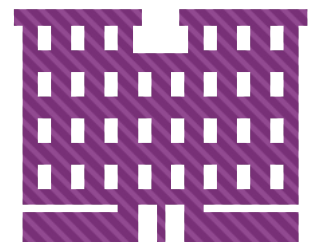
$$\text{Penalty, } p_i = \sum_{k=1}^K \text{est-err}_k + \sum_{k=1}^K c_{i,k} n_{i,k}$$

Theorem: There exists a NIC and IR mechanism for which,

$$P(M, s^*) \leq 8\sqrt{m} \cdot \inf_{M,s} P(M, s)$$

Theorem (hardness): There exists a set of costs $\{c_{i,k}\}_{i,k}$ such that for any NIC and IR mechanism, we have

$$P(M, s^*) \in \Omega\left(\sqrt{m}\right) \cdot \inf_{M,s} P(M, s)$$

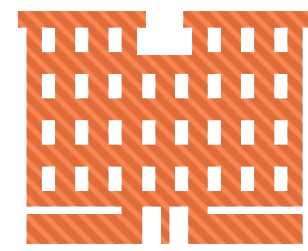
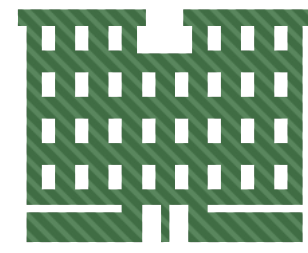
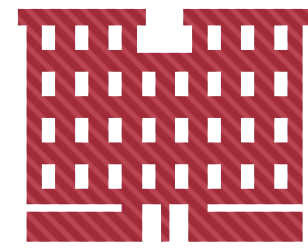
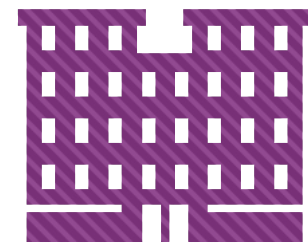


**Data
contributors**

Marketplace



**Data
consumers**



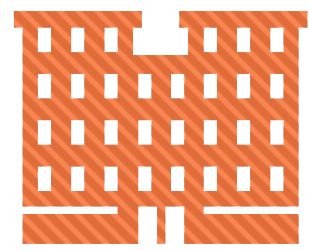
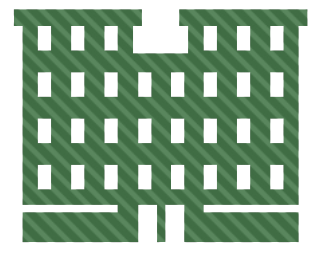
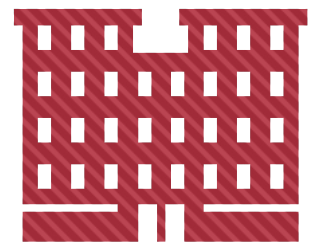
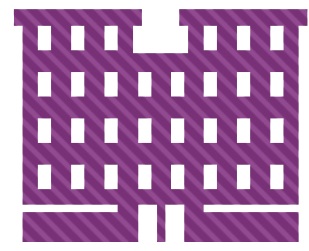
Data contributors

Marketplace



Data consumers

Consumers purchase data from contributors via a marketplace:



Data contributors

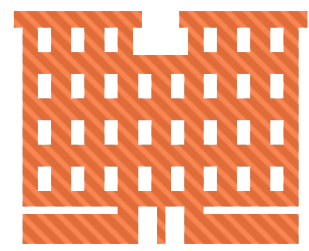
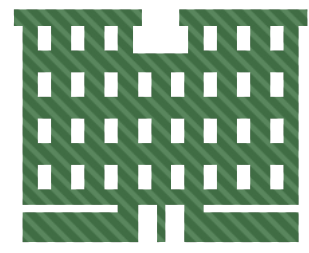
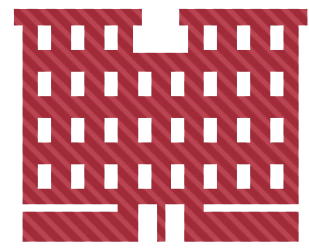
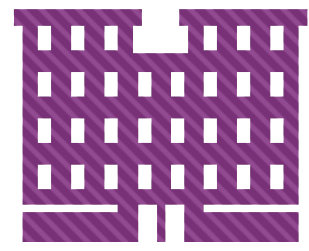
Marketplace



Data consumers

Consumers purchase data from contributors via a marketplace:

- ▶ Ensure contributors do not fabricate/poison data.



**Data
contributors**

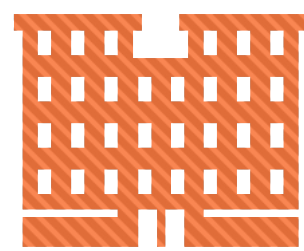
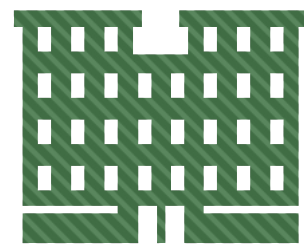
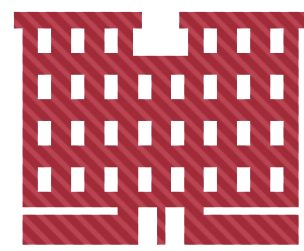
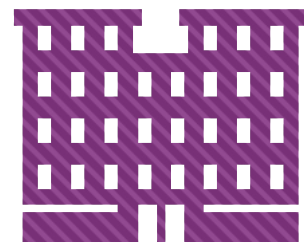
Marketplace



**Data
consumers**

Consumers purchase data from contributors via a marketplace:

- ▶ Ensure contributors do not fabricate/poison data.
- ▶ Pricing data that is being sold to consumers.



Data contributors

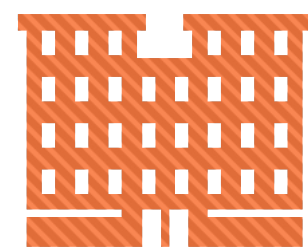
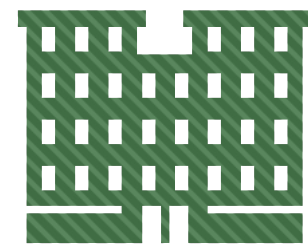
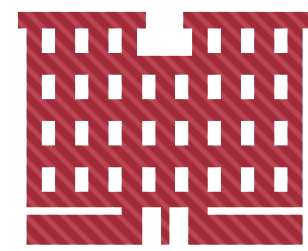
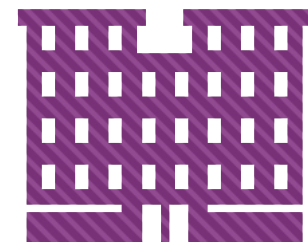
Marketplace



Data consumers

Consumers purchase data from contributors via a marketplace:

- ▶ Ensure contributors do not fabricate/poison data.
- ▶ Pricing data that is being sold to consumers.
- ▶ Re-distributing the revenue back to the contributors.



Data contributors

Marketplace



Data consumers

Consumers purchase data from contributors via a marketplace:

- ▶ Ensure contributors do not fabricate/poison data.
- ▶ Pricing data that is being sold to consumers.
- ▶ Re-distributing the revenue back to the contributors.
- ▶ Learn consumer valuation of data via online feedback.



**Keran
Chen**



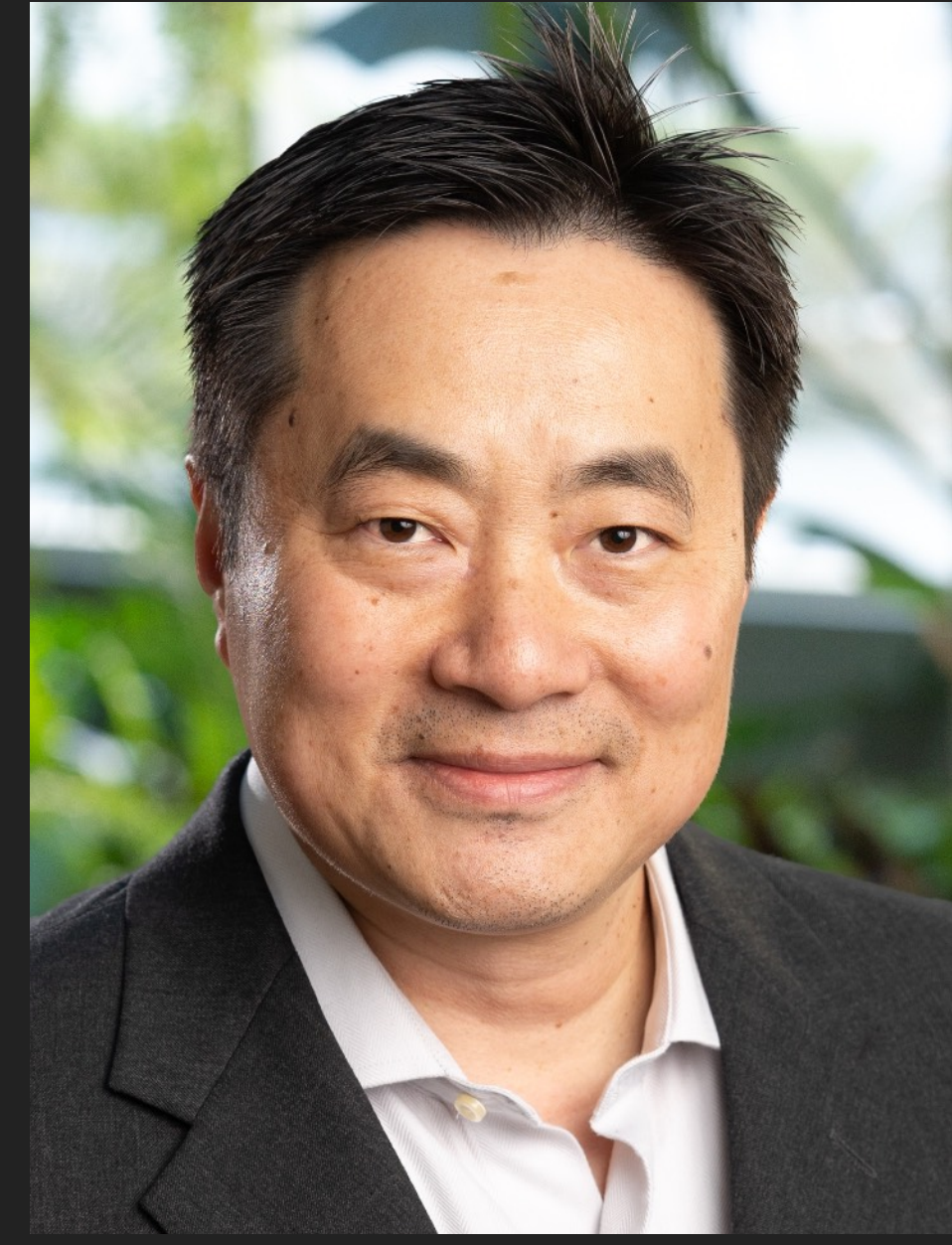
**Yiding
Chen**



**Alex
Clinton**



**Joon
Suk Huh**



**Jerry
Zhu**

THANK YOU!

kandasamy@cs.wisc.edu

- ▶ Data sharing has many benefits
 - ▶ Maximize the value created by data.
 - ▶ Democratize data
- ▶ But strategic agents can free-ride in naive mechanisms, either by not contributing data, or contributing fabricated datasets.
- ▶ Our mechanism is IR and NIC while achieving a factor 2 of the global minimum social penalty.