

# Relaxing the I.I.D. Assumption

Adaptively Minimax Optimal Regret via Root-Entropic Regularization

---

Blair Bilodeau<sup>\*,1,2,3</sup> and Jeffrey Negrea<sup>\*,1,2,3</sup>

(Joint work with Daniel M. Roy<sup>1,2,3</sup>)

February 1, 2021

Presented to the RIKEN Center for Advanced Intelligence Project

\* Equal Contribution

<sup>1</sup>Department of Statistical Sciences, University of Toronto

<sup>2</sup>Vector Institute

<sup>3</sup>Institute for Advanced Study

# Background

---

# A Motivating Example

## Stock Market Analogy

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

**What assumptions should we make?**

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

Is assuming adversarial data too pessimistic?

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

Is assuming adversarial data too pessimistic?

Is the departure from I.I.D.-ness benign? How can we quantify that?

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

Is assuming adversarial data too pessimistic?

Is the departure from I.I.D.-ness benign? How can we quantify that?

Influence of non-stochastic forces “small”  $\Rightarrow$  maybe.

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

Is assuming adversarial data too pessimistic?

Is the departure from I.I.D.-ness benign? How can we quantify that?

Influence of non-stochastic forces “small”  $\Rightarrow$  maybe.

Meaning of “small” TBD.

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

Is assuming adversarial data too pessimistic?

Is the departure from I.I.D.-ness benign? How can we quantify that?

Influence of non-stochastic forces “small”  $\Rightarrow$  maybe.

Meaning of “small” TBD.

Want to maximize profit without having to know what drives the market.

# A Motivating Example

## Stock Market Analogy

- You need to invest your money into a stock portfolio.
- You have access to several market experts that give you advice.
- You *regret* not having always followed the *post hoc* best expert's advice

## What assumptions should we make?

A simplifying assumption is that the data are I.I.D. (e.g., Black–Scholes–Merton)

In real life, market is driven in part by non-stochastic forces.

Is assuming adversarial data too pessimistic?

Is the departure from I.I.D.-ness benign? How can we quantify that?

Influence of non-stochastic forces “small”  $\Rightarrow$  maybe.

Meaning of “small” TBD.

Want to maximize profit without having to know what drives the market.

# Sequential Prediction with Expert Advice

Sequential Prediction a.k.a. Online Learning

# Sequential Prediction with Expert Advice

## Sequential Prediction a.k.a. Online Learning

For rounds  $t = 1, \dots, T$ :

# Sequential Prediction with Expert Advice

## Sequential Prediction a.k.a. Online Learning

For rounds  $t = 1, \dots, T$ :

- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$

# Sequential Prediction with Expert Advice

## Sequential Prediction a.k.a. Online Learning

For rounds  $t = 1, \dots, T$ :

- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$
- Observe  $y(t) \in \mathcal{Y}$  from the environment

# Sequential Prediction with Expert Advice

## Sequential Prediction a.k.a. Online Learning

For rounds  $t = 1, \dots, T$ :

- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

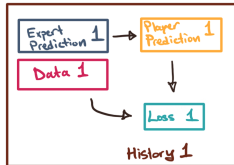


# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

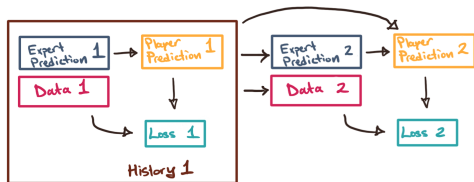


# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

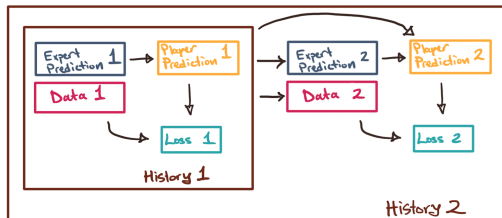


# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

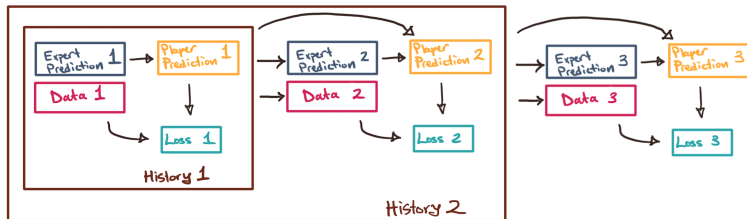


# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$

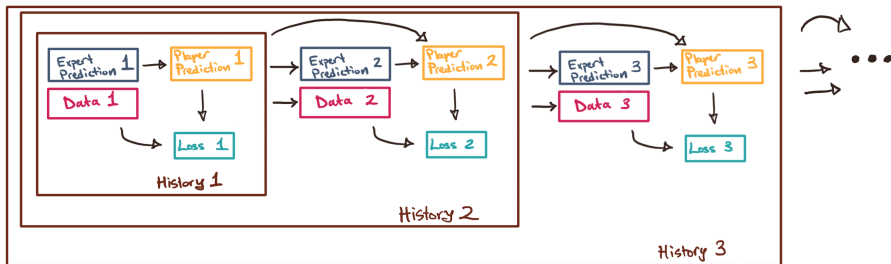


# Sequential Prediction with Expert Advice

## Sequential Prediction **with Expert Advice**

For rounds  $t = 1, \dots, T$ :

- Receive  $\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$  expert predictions
- Predict  $\hat{y}(t) \in \hat{\mathcal{Y}}$  based on historical data before time  $t$  **and expert predictions**
- Observe  $y(t) \in \mathcal{Y}$  from the environment
- Incur loss  $\ell(\hat{y}(t), y(t))$



# Measuring Performance

The measure of the **player's performance** is...

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  *reference experts*;

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  **reference experts**;
- Given by the excess cumulative **loss** of the player over the **best expert**;

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  **reference experts**;
- Given by the excess cumulative **loss** of the player over the **best expert**;

$$\text{Regret: } R(T) = \sum_{t=1}^T \ell(\hat{y}(t), y(t)) - \min_{i \in [N]} \sum_{t=1}^T \ell(x_i(t), y(t))$$

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  **reference experts**;
- Given by the excess cumulative **loss** of the player over the **best expert**;

$$\text{Regret: } R(T) = \sum_{t=1}^T \ell(\hat{y}(t), y(t)) - \min_{i \in [N]} \sum_{t=1}^T \ell(x_i(t), y(t))$$

The prediction problem is *online learnable* if a player can incur sub-linear regret:

$$\mathbb{E}R(T) \in o(T).$$

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  **reference experts**;
- Given by the excess cumulative **loss** of the player over the **best expert**;

$$\text{Regret: } R(T) = \sum_{t=1}^T \ell(\hat{y}(t), y(t)) - \min_{i \in [N]} \sum_{t=1}^T \ell(x_i(t), y(t))$$

The prediction problem is *online learnable* if a player can incur sub-linear regret:

$$\mathbb{E}R(T) \in o(T).$$

Where the  $\mathbb{E}$  is taken with respect to the randomness in the player's and expert's predictions, and the data-generating mechanism for  $(y(t))_{t \in \mathbb{N}}$ .

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  **reference experts**;
- Given by the excess cumulative **loss** of the player over the **best expert**;

$$\text{Regret: } R(T) = \sum_{t=1}^T \ell(\hat{y}(t), y(t)) - \min_{i \in [N]} \sum_{t=1}^T \ell(x_i(t), y(t))$$

The prediction problem is *online learnable* if a player can incur sub-linear regret:

$$\mathbb{E}R(T) \in o(T).$$

Where the  $\mathbb{E}$  is taken with respect to the randomness in the player's and expert's predictions, and the data-generating mechanism for  $(y(t))_{t \in \mathbb{N}}$ .

(The  $\mathbb{E}$  may be under a complicated, non-I.I.D. measure.)

# Measuring Performance

The measure of the **player's performance** is...

- Relative to the class of  $N$  **reference experts**;
- Given by the excess cumulative **loss** of the player over the **best expert**;

$$\text{Regret: } R(T) = \sum_{t=1}^T \ell(\hat{y}(t), y(t)) - \min_{i \in [N]} \sum_{t=1}^T \ell(x_i(t), y(t))$$

The prediction problem is *online learnable* if a player can incur sub-linear regret:

$$\mathbb{E}R(T) \in o(T).$$

Where the  $\mathbb{E}$  is taken with respect to the randomness in the player's and expert's predictions, and the data-generating mechanism for  $(y(t))_{t \in \mathbb{N}}$ .

(The  $\mathbb{E}$  may be under a complicated, non-I.I.D. measure.)

# Optimality in the Stochastic and Adversarial Regimes

# Optimality in the Stochastic and Adversarial Regimes

Stochastic-with-a-Gap

# Optimality in the Stochastic and Adversarial Regimes

## Stochastic-with-a-Gap

- Expert predictions and data are I.I.D. over time from some distribution.
- There is an expert whose mean loss is  $\Delta$  smaller than the others.

# Optimality in the Stochastic and Adversarial Regimes

## Stochastic-with-a-Gap

- Expert predictions and data are I.I.D. over time from some distribution.
- There is an expert whose mean loss is  $\Delta$  smaller than the others.

### Theorem (Gaillard et al. 2014 + Mourtada and Gaïffas 2019)

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \frac{\log N}{\Delta}, \text{ uniformly bounded in } T.$$

# Optimality in the Stochastic and Adversarial Regimes

## Stochastic-with-a-Gap

- Expert predictions and data are I.I.D. over time from some distribution.
- There is an expert whose mean loss is  $\Delta$  smaller than the others.

### Theorem (Gaillard et al. 2014 + Mourtada and Gaïffas 2019)

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \frac{\log N}{\Delta}, \text{ uniformly bounded in } T.$$

## Adversarial

- Compete against expert predictions and data that maximize  $R(T)$ .

# Optimality in the Stochastic and Adversarial Regimes

## Stochastic-with-a-Gap

- Expert predictions and data are I.I.D. over time from some distribution.
- There is an expert whose mean loss is  $\Delta$  smaller than the others.

### Theorem (Gaillard et al. 2014 + Mourtada and Gaïffas 2019)

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \frac{\log N}{\Delta}, \text{ uniformly bounded in } T.$$

## Adversarial

- Compete against expert predictions and data that maximize  $R(T)$ .

### Theorem (Vovk 1998, see also [FS97; CL06])

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \sqrt{T \log N} \text{ for all } T.$$

# Optimality in the Stochastic and Adversarial Regimes

## Stochastic-with-a-Gap

- Expert predictions and data are I.I.D. over time from some distribution.
- There is an expert whose mean loss is  $\Delta$  smaller than the others.

### Theorem (Gaillard et al. 2014 + Mourtada and Gaïffas 2019)

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \frac{\log N}{\Delta}, \text{ uniformly bounded in } T.$$

## Adversarial

- Compete against expert predictions and data that maximize  $R(T)$ .

### Theorem (Vovk 1998, see also [FS97; CL06])

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \sqrt{T \log N} \text{ for all } T.$$

Can a single algorithm be optimal in both settings simultaneously?

# Optimality in the Stochastic and Adversarial Regimes

## Stochastic-with-a-Gap

- Expert predictions and data are I.I.D. over time from some distribution.
- There is an expert whose mean loss is  $\Delta$  smaller than the others.

### Theorem (Gaillard et al. 2014 + Mourtada and Gaïffas 2019)

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \frac{\log N}{\Delta}, \text{ uniformly bounded in } T.$$

## Adversarial

- Compete against expert predictions and data that maximize  $R(T)$ .

### Theorem (Vovk 1998, see also [FS97; CL06])

A constructive algorithm achieves the minimax regret:

$$\mathbb{E}R(T) \asymp \sqrt{T \log N} \text{ for all } T.$$

Can a single algorithm be optimal in both settings simultaneously?

# Simultaneous Optimality of Hedge

Can a single algorithm be optimal in both settings simultaneously? **Yes!** [MG19]

# Simultaneous Optimality of Hedge

Can a single algorithm be optimal in both settings simultaneously? **Yes!** [MG19]

**Stochastic-with-a-gap:**  $\mathbb{E}R(T) \asymp (\log N)/\Delta$  uniformly in  $T$ .

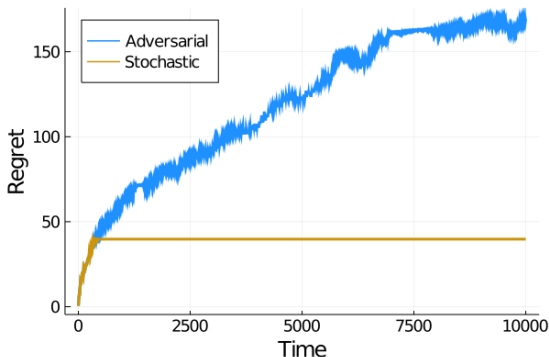
**Adversarial:**  $\mathbb{E}R(T) \asymp \sqrt{T \log N}$

# Simultaneous Optimality of Hedge

Can a single algorithm be optimal in both settings simultaneously? **Yes!** [MG19]

**Stochastic-with-a-gap:**  $\mathbb{E}R(T) \asymp (\log N)/\Delta$  uniformly in  $T$ .

**Adversarial:**  $\mathbb{E}R(T) \asymp \sqrt{T \log N}$



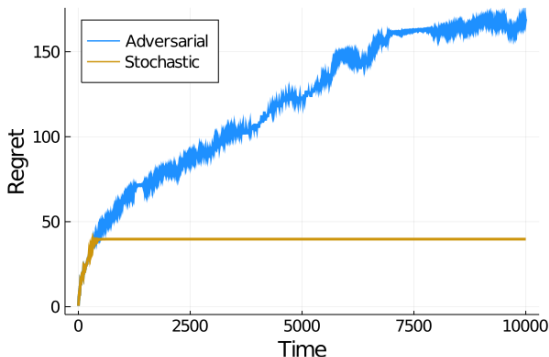
**The same algorithm, Hedge, was used in both cases!**

# Simultaneous Optimality of Hedge

Can a single algorithm be optimal in both settings simultaneously? **Yes!** [MG19]

**Stochastic-with-a-gap:**  $\mathbb{E}R(T) \asymp (\log N)/\Delta$  uniformly in  $T$ .

**Adversarial:**  $\mathbb{E}R(T) \asymp \sqrt{T \log N}$



**The same algorithm, Hedge, was used in both cases!**

# Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.

# Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

# Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.

## Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.  $\leftarrow$  Too pessimistic.

# Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.  $\leftarrow$  Too pessimistic.

We provide a spectrum between stochastic and adversarial;

# Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.  $\leftarrow$  Too pessimistic.

We provide a spectrum between stochastic and adversarial;

Intuitively, fix a “neighbourhood” of distributions;

# Beyond Stochastic and Adversarial

Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.  $\leftarrow$  Too pessimistic.

We provide a spectrum between stochastic and adversarial;

Intuitively, fix a “neighbourhood” of distributions;

Each data point drawn from an arbitrary distribution in “neighbourhood”.

# Beyond Stochastic and Adversarial

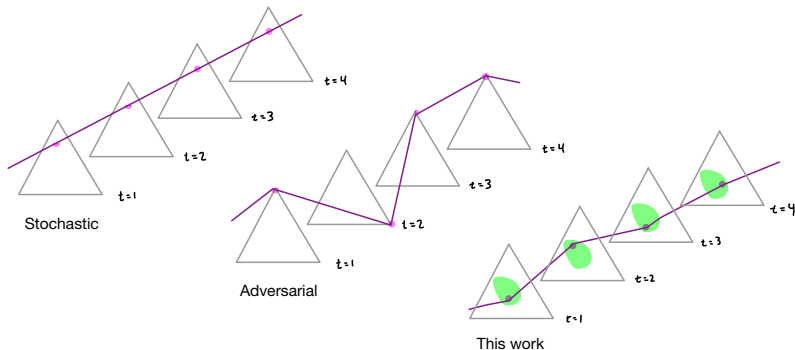
Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.  $\leftarrow$  Too pessimistic.

We provide a spectrum between stochastic and adversarial;

Intuitively, fix a “neighbourhood” of distributions;

Each data point drawn from an arbitrary distribution in “neighbourhood”.



# Beyond Stochastic and Adversarial

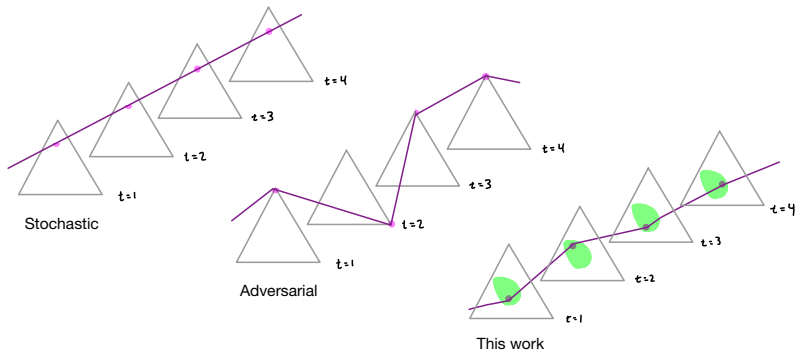
Real data  $\neq$  stochastic.  $\leftarrow$  Too optimistic.

Real data  $\neq$  adversarial.  $\leftarrow$  Too pessimistic.

We provide a spectrum between stochastic and adversarial;

Intuitively, fix a “neighbourhood” of distributions;

Each data point drawn from an arbitrary distribution in “neighbourhood”.



# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

## Definition BNR20

An algorithm is **adaptively minimax optimal** for a spectrum of settings if:

- it achieves the minimax optimal performance in each setting; and
- it does not require knowledge of the true setting in advance.

# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

## Definition BNR20

An algorithm is **adaptively minimax optimal** for a spectrum of settings if:

- it achieves the minimax optimal performance in each setting; and
- it does not require knowledge of the true setting in advance.

How to formalize this?

# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

## Definition BNR20

An algorithm is **adaptively minimax optimal** for a spectrum of settings if:

- it achieves the minimax optimal performance in each setting; and
- it does not require knowledge of the true setting in advance.

How to formalize this?

- For an abstract range of settings  $\Theta$ ...

# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

## Definition BNR20

An algorithm is **adaptively minimax optimal** for a spectrum of settings if:

- it achieves the minimax optimal performance in each setting; and
- it does not require knowledge of the true setting in advance.

## How to formalize this?

- For an abstract range of settings  $\Theta$ ...
- Parameterize the minimax regret in each setting:  $(R_{\theta}^*(T))_{\theta \in \Theta}$ .

# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

## Definition BNR20

An algorithm is **adaptively minimax optimal** for a spectrum of settings if:

- it achieves the minimax optimal performance in each setting; and
- it does not require knowledge of the true setting in advance.

## How to formalize this?

- For an abstract range of settings  $\Theta$ ...
- Parameterize the minimax regret in each setting:  $(R_\theta^*(T))_{\theta \in \Theta}$ .
- Algorithm satisfies  $R_\theta(T) \leq C R_\theta^*(T)$  uniformly in  $\theta$  for large enough  $T$ .

# Adaptively Minimax Optimal Algorithms

Prediction algorithms should be robust to a range of data generating mechanisms.

## Definition BNR20

An algorithm is **adaptively minimax optimal** for a spectrum of settings if:

- it achieves the minimax optimal performance in each setting; and
- it does not require knowledge of the true setting in advance.

## How to formalize this?

- For an abstract range of settings  $\Theta$ ...
- Parameterize the minimax regret in each setting:  $(R_\theta^*(T))_{\theta \in \Theta}$ .
- Algorithm satisfies  $R_\theta(T) \leq C R_\theta^*(T)$  uniformly in  $\theta$  for large enough  $T$ .

# Overview of Our Work

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

We initially set out hoping to prove that Hedge was adaptive to all scenarios.

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

We initially set out hoping to prove that Hedge was adaptive to all scenarios.

## **Theorem BNR20**

Without oracle knowledge to tune the learning rate, Hedge is not simultaneously minimax optimal at all settings between stochastic-with-a-gap and adversarial.

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

We initially set out hoping to prove that Hedge was adaptive to all scenarios.

## **Theorem BNR20**

Without oracle knowledge to tune the learning rate, Hedge is not simultaneously minimax optimal at all settings between stochastic-with-a-gap and adversarial.

We provide a new algorithm that achieves the minimax rate in all settings...

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

We initially set out hoping to prove that Hedge was adaptive to all scenarios.

## **Theorem BNR20**

Without oracle knowledge to tune the learning rate, Hedge is not simultaneously minimax optimal at all settings between stochastic-with-a-gap and adversarial.

We provide a new algorithm that achieves the minimax rate in all settings...  
...without knowledge of which setting prevails!

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

We initially set out hoping to prove that Hedge was adaptive to all scenarios.

## Theorem BNR20

Without oracle knowledge to tune the learning rate, Hedge is not simultaneously minimax optimal at all settings between stochastic-with-a-gap and adversarial.

We provide a new algorithm that achieves the minimax rate in all settings...  
...without knowledge of which setting prevails!

## Theorem BNR20

There is an adaptively minimax optimal algorithm: Meta-CARE.

# Overview of Our Work

We show Hedge is suboptimal between Stochastic and Adversarial.

This was surprising for us.

We initially set out hoping to prove that Hedge was adaptive to all scenarios.

## Theorem BNR20

Without oracle knowledge to tune the learning rate, Hedge is not simultaneously minimax optimal at all settings between stochastic-with-a-gap and adversarial.

We provide a new algorithm that achieves the minimax rate in all settings...  
...without knowledge of which setting prevails!

## Theorem BNR20

There is an adaptively minimax optimal algorithm: Meta-CARE.

# Main Result

## Motivating Intuition

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know *only*  $N_0$  of the experts can ever be “the best”, and which ones, ...

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know *only*  $N_0$  of the experts can ever be “the best”, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know **only  $N_0$  of the experts can ever be “the best”**, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”
  - so we might strive to have regret  $\Theta(\sqrt{T \log N_0})$  in  $(T, N_0)$

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know **only  $N_0$  of the experts can ever be “the best”**, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”
  - so we might strive to have regret  $\Theta(\sqrt{T \log N_0})$  in  $(T, N_0)$
- If we know **one expert is better than the rest by  $\Delta_0$** , but not which it is...

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know *only*  $N_0$  of the experts can ever be “the best”, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”
  - so we might strive to have regret  $\Theta(\sqrt{T \log N_0})$  in  $(T, N_0)$
- If we know *one expert is better than the rest by*  $\Delta_0$ , but not which it is...
  - then we are *almost* in the stochastic-with-a-gap case

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know *only*  $N_0$  of the experts can ever be “the best”, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”
  - so we might strive to have regret  $\Theta(\sqrt{T \log N_0})$  in  $(T, N_0)$
- If we know *one expert is better than the rest by*  $\Delta_0$ , but not which it is...
  - then we are *almost* in the stochastic-with-a-gap case
  - so we might hope for regret  $\Theta((\log N)/\Delta_0)$

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know **only  $N_0$  of the experts can ever be “the best”**, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”
  - so we might strive to have regret  $\Theta(\sqrt{T \log N_0})$  in  $(T, N_0)$
- If we know **one expert is better than the rest by  $\Delta_0$** , but not which it is...
  - then we are *almost* in the stochastic-with-a-gap case
  - so we might hope for regret  $\Theta((\log N)/\Delta_0)$

## Theorem BNR20

The adaptively minimax optimal rate of regret, which Meta-CARE achieves, is

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & N_0 \geq 2 \\ (\log N)/\Delta_0 & N_0 = 1. \end{cases}$$

# Main Result

## Motivating Intuition

- In the adversarial case the minimax optimal regret is  $\Theta(\sqrt{T \log N})$
- If we know **only  $N_0$  of the experts can ever be “the best”**, and which ones, ...
  - we could restrict an adversarially optimal algorithm to the “best experts”
  - so we might strive to have regret  $\Theta(\sqrt{T \log N_0})$  in  $(T, N_0)$
- If we know **one expert is better than the rest by  $\Delta_0$** , but not which it is...
  - then we are *almost* in the stochastic-with-a-gap case
  - so we might hope for regret  $\Theta((\log N)/\Delta_0)$

## Theorem BNR20

The adaptively minimax optimal rate of regret, which Meta-CARE achieves, is

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & N_0 \geq 2 \\ (\log N)/\Delta_0 & N_0 = 1. \end{cases}$$

# Hedge Algorithm

---

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

A proper algorithm assigns probability  $w_i(t)$  to expert  $i$  at time  $t$ .

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

A proper algorithm assigns probability  $w_i(t)$  to expert  $i$  at time  $t$ .

## Hedge Algorithm

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

A proper algorithm assigns probability  $w_i(t)$  to expert  $i$  at time  $t$ .

## Hedge Algorithm

- Fix learning rate schedule  $\eta : \mathbb{N} \rightarrow \mathbb{R}$ ; initialize the weights as uniform; define

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

A proper algorithm assigns probability  $w_i(t)$  to expert  $i$  at time  $t$ .

## Hedge Algorithm

- Fix learning rate schedule  $\eta : \mathbb{N} \rightarrow \mathbb{R}$ ; initialize the weights as uniform; define

$$\ell_i(t) = \ell(x_i(t), y(t)), \quad L_i(t) = \sum_{s=1}^t \ell_i(s).$$

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

A proper algorithm assigns probability  $w_i(t)$  to expert  $i$  at time  $t$ .

## Hedge Algorithm

- Fix learning rate schedule  $\eta : \mathbb{N} \rightarrow \mathbb{R}$ ; initialize the weights as uniform; define

$$\ell_i(t) = \ell(x_i(t), y(t)), \quad L_i(t) = \sum_{s=1}^t \ell_i(s).$$

- Update weights for each  $i \in [N]$  using

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

# Hedge Algorithm

We will consider only finite expert classes and bounded losses  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow [0, 1]$ .

All explicit algorithms we will consider are *proper*:

the player randomly selects an expert to emulate at each time.

A proper algorithm assigns probability  $w_i(t)$  to expert  $i$  at time  $t$ .

## Hedge Algorithm

- Fix learning rate schedule  $\eta : \mathbb{N} \rightarrow \mathbb{R}$ ; initialize the weights as uniform; define

$$\ell_i(t) = \ell(x_i(t), y(t)), \quad L_i(t) = \sum_{s=1}^t \ell_i(s).$$

- Update weights for each  $i \in [N]$  using

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[N]$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[N]$

- with a flat prior, and

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[N]$

- with a flat prior, and
- model likelihood  $\exp\{-\eta \ell_i(t)\}$  for the  $t$ -th observation under parameter  $i$ .

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[N]$

- with a flat prior, and
- model likelihood  $\exp\{-\eta \ell_i(t)\}$  for the  $t$ -th observation under parameter  $i$ .

## Variational Formulation

$$w(t) = \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t-1) \rangle - \frac{1}{\eta(t)} H(w) \right)$$

where

$$H(w) = - \sum_{i \in [N]} w_i \log(w_i).$$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[N]$

- with a flat prior, and
- model likelihood  $\exp\{-\eta \ell_i(t)\}$  for the  $t$ -th observation under parameter  $i$ .

## Variational Formulation

$$w(t) = \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t-1) \rangle - \frac{1}{\eta(t)} H(w) + \frac{1}{\eta(t)} \sum_{i=1}^N w_i \log(N) \right)$$

where

$$H(w) = - \sum_{i \in [N]} w_i \log(w_i).$$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[N]$

- with a flat prior, and
- model likelihood  $\exp\{-\eta \ell_i(t)\}$  for the  $t$ -th observation under parameter  $i$ .

## Variational Formulation

$$w(t) = \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t-1) \rangle + \frac{1}{\eta(t)} \text{KL}(w \| \text{Unif}([N])) \right)$$

where

$$\text{KL}(w \| p) = \sum_{i \in [N]} w_i \log(w_i / p_i).$$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[M]$

- with a flat prior, and
- model likelihood  $\exp\{-\eta \ell_i(t)\}$  for the  $t$ -th observation under parameter  $i$ .

## Variational Formulation

$$w(t) = \arg \min_{w \in \text{simp}([M])} \left( \langle w, L(t-1) \rangle + \frac{1}{\eta(t)} \text{KL}(w \| \text{Unif}([M])) \right)$$

## Gibbs Posterior

$$\hat{\pi}_t(\theta) = \arg \min_{\hat{\pi} \in \mathcal{M}(\Theta)} \left( \mathbb{E}_{\theta \sim \hat{\pi}} L_\theta(t-1) + \frac{1}{\eta(t)} \text{KL}(\hat{\pi} \| \pi) \right)$$

# Hedge as Bayesian Inference

$$w_i(t) \propto \exp \{ -\eta(t) L_i(t-1) \}.$$

For  $\eta(t) = \eta$  constant in  $t$ , this looks like Bayes rule for a parameter in  $[M]$

- with a flat prior, and
- model likelihood  $\exp\{-\eta \ell_i(t)\}$  for the  $t$ -th observation under parameter  $i$ .

## Variational Formulation

$$w(t) = \arg \min_{w \in \text{simp}([M])} \left( \langle w, L(t-1) \rangle + \frac{1}{\eta(t)} \text{KL}(w \| \text{Unif}([M])) \right)$$

## Gibbs Posterior

$$\hat{\pi}_t(\theta) = \arg \min_{\hat{\pi} \in \mathcal{M}(\Theta)} \left( \mathbb{E}_{\theta \sim \hat{\pi}} L_\theta(t-1) + \frac{1}{\eta(t)} \text{KL}(\hat{\pi} \| \pi) \right)$$

$$\hat{\pi}_t(\theta) \propto \pi(\theta) \exp\{-\eta(t) L_\theta(t-1)\}$$

## Relaxing the I.I.D. Assumption

---

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .
  - **Time-Homogeneous:**  $\mathcal{D}$  does not depend on  $t$

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .
  - Time-Homogeneous:  $\mathcal{D}$  does not depend on  $t$
  - Convexity  $\Leftrightarrow$  environment can flip a coin to select between basic elements of  $\mathcal{D}$

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .
  - Time-Homogeneous:  $\mathcal{D}$  does not depend on  $t$
  - Convexity  $\Leftrightarrow$  environment can flip a coin to select between basic elements of  $\mathcal{D}$
  - Environment may aim to maximize regret subject to the constraint

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .
  - Time-Homogeneous:  $\mathcal{D}$  does not depend on  $t$
  - Convexity  $\Leftrightarrow$  environment can flip a coin to select between basic elements of  $\mathcal{D}$
  - Environment may aim to maximize regret subject to the constraint
- The choice of distribution is made based on outcomes of the previous rounds.

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .
  - Time-Homogeneous:  $\mathcal{D}$  does not depend on  $t$
  - Convexity  $\Leftrightarrow$  environment can flip a coin to select between basic elements of  $\mathcal{D}$
  - Environment may aim to maximize regret subject to the constraint
- The choice of distribution is made based on outcomes of the previous rounds.

# Our Setting: Time-Homogeneous Convex Constraints

## Intuition

Experts and observations may collude.

Realizations  $(x(t), y(t))$  are sampled from an adversarial conditional distribution.

## Formal Framework

- Fix a convex set of distributions  $\mathcal{D} \subseteq \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$ .
- $(x(t), y(t))$  drawn from an element of  $\mathcal{D}$  given the history prior to  $t$ .
  - Time-Homogeneous:  $\mathcal{D}$  does not depend on  $t$
  - Convexity  $\Leftrightarrow$  environment can flip a coin to select between basic elements of  $\mathcal{D}$
  - Environment may aim to maximize regret subject to the constraint
- The choice of distribution is made based on outcomes of the previous rounds.

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\},$

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\},$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\},$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ ,

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ , e.g.  $\text{Ball}(\mu_0, r)$  for  $r > 0$

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ , e.g.  $\text{Ball}(\mu_0, r)$  for  $r > 0$
- $r \rightarrow 0$  gives the stochastic case, specifically I.I.D.  $\mu_0$ .

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ , e.g.  $\text{Ball}(\mu_0, r)$  for  $r > 0$
- $r \rightarrow 0$  gives the stochastic case, specifically I.I.D.  $\mu_0$ .
- $r \rightarrow \infty$  gives adversarial case.

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ , e.g.  $\text{Ball}(\mu_0, r)$  for  $r > 0$
- $r \rightarrow 0$  gives the stochastic case, specifically I.I.D.  $\mu_0$ .
- $r \rightarrow \infty$  gives adversarial case. Smoothly transitions in between as  $r$  varies.

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ , e.g.  $\text{Ball}(\mu_0, r)$  for  $r > 0$
- $r \rightarrow 0$  gives the stochastic case, specifically I.I.D.  $\mu_0$ .
- $r \rightarrow \infty$  gives adversarial case. Smoothly transitions in between as  $r$  varies.
- A small neighborhood leads to a slight relaxation of I.I.D.-ness.

# Examples

**Stochastic:**  $\mathcal{D} = \{\mu_0\}$ ,

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y}) \leftarrow$  contains point masses!

**Adversarial-with-an- $\mathbb{E}$ -gap** (Mourtada and Gaïffas 2019)

- One expert has at least  $\Delta > 0$  less  $\mathbb{E}$  loss than the rest on every round.

**Neighborhood-of-I.I.D.**

- Fix a metric on the space of distributions over  $\hat{\mathcal{Y}}^N \times \mathcal{Y}$
- Pick any  $\mu_0$ , and let  $\mathcal{D}$  be a neighborhood of  $\mu_0$ , e.g.  $\text{Ball}(\mu_0, r)$  for  $r > 0$
- $r \rightarrow 0$  gives the stochastic case, specifically I.I.D.  $\mu_0$ .
- $r \rightarrow \infty$  gives adversarial case. Smoothly transitions in between as  $r$  varies.
- A small neighborhood leads to a slight relaxation of I.I.D.-ness.

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

**Effective Experts**

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

## Effective Experts

$$\mathcal{I}_0 = \{\text{experts that are optimal in } \mathbb{E} \text{ for some } \mu \in \mathcal{D}\}$$

$$N_0 = |\mathcal{I}_0|$$

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

## Effective Experts

$$\mathcal{I}_0 = \{\text{experts that are optimal in } \mathbb{E} \text{ for some } \mu \in \mathcal{D}\}$$

$$N_0 = |\mathcal{I}_0|$$

Analogous to the single best expert in the stochastic-with-a-gap setting.

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

## Effective Experts

$$\mathcal{I}_0 = \{\text{experts that are optimal in } \mathbb{E} \text{ for some } \mu \in \mathcal{D}\}$$

$$N_0 = |\mathcal{I}_0|$$

Analogous to the single best expert in the stochastic-with-a-gap setting.

## Effective Stochastic Gap

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

## Effective Experts

$$\mathcal{I}_0 = \{\text{experts that are optimal in } \mathbb{E} \text{ for some } \mu \in \mathcal{D}\}$$

$$N_0 = |\mathcal{I}_0|$$

Analogous to the single best expert in the stochastic-with-a-gap setting.

## Effective Stochastic Gap

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

## Effective Experts

$$\mathcal{I}_0 = \{\text{experts that are optimal in } \mathbb{E} \text{ for some } \mu \in \mathcal{D}\}$$

$$N_0 = |\mathcal{I}_0|$$

Analogous to the single best expert in the stochastic-with-a-gap setting.

## Effective Stochastic Gap

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

Analogous to the gap in the stochastic-with-a-gap setting.

# Constraint-Characterizing Quantities

We use quantities to characterize the constraint that:

- are representative of whether the data is “easy” or not;
- yield matching lower and upper bounds on regret.

## Effective Experts

$$\mathcal{I}_0 = \{\text{experts that are optimal in } \mathbb{E} \text{ for some } \mu \in \mathcal{D}\}$$

$$N_0 = |\mathcal{I}_0|$$

Analogous to the single best expert in the stochastic-with-a-gap setting.

## Effective Stochastic Gap

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

Analogous to the gap in the stochastic-with-a-gap setting.

## Constraint Example

$$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\} \quad N_0 = |\mathcal{I}_0|$$

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,

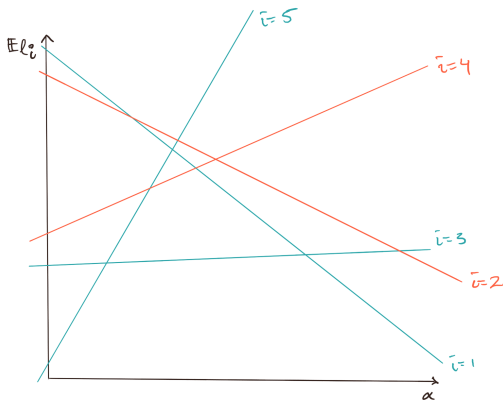
$$N = 5, \mathcal{I}_0 = \{1, 3, 5\}, N_0 = 3.$$

# Constraint Example

$$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\} \quad N_0 = |\mathcal{I}_0|$$

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,  
 $N = 5$ ,  $\mathcal{I}_0 = \{1, 3, 5\}$ ,  $N_0 = 3$ .



# Constraint Example

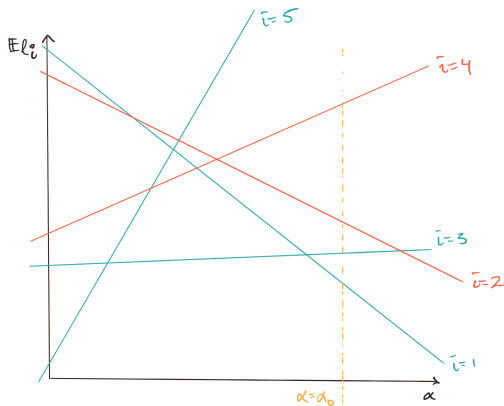
$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\}$

$$N_0 = |\mathcal{I}_0|$$

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,

$$N = 5, \mathcal{I}_0 = \{1, 3, 5\}, N_0 = 3.$$



# Constraint Example

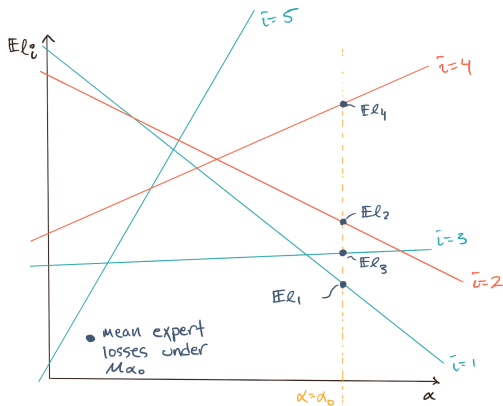
$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\}$

$$N_0 = |\mathcal{I}_0|$$

$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,

$$N = 5, \mathcal{I}_0 = \{1, 3, 5\}, N_0 = 3.$$



# Constraint Example

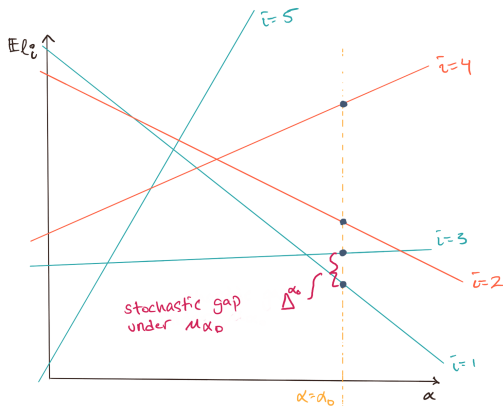
$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\}$

$$N_0 = |\mathcal{I}_0|$$

$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,

$$N = 5, \mathcal{I}_0 = \{1, 3, 5\}, N_0 = 3.$$



# Constraint Example

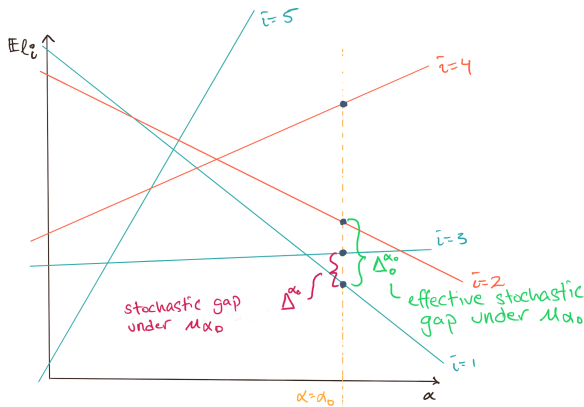
$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\}$

$$N_0 = |\mathcal{I}_0|$$

$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,

$$N = 5, \mathcal{I}_0 = \{1, 3, 5\}, N_0 = 3.$$



# Constraint Example

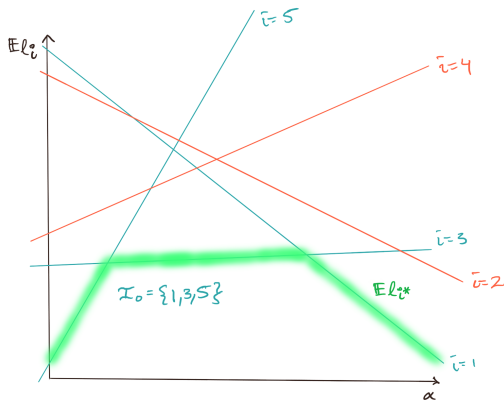
$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\}$

$$N_0 = |\mathcal{I}_0|$$

$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,

$$N = 5, \mathcal{I}_0 = \{1, 3, 5\}, N_0 = 3.$$

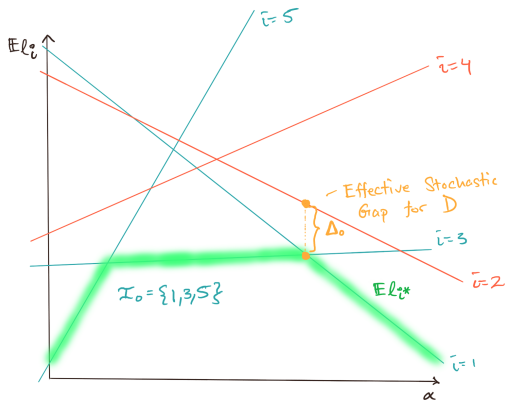


# Constraint Example

$$\mathcal{I}_0 = \{\text{experts that are optimal for some } \mu \in \mathcal{D}\} \quad N_0 = |\mathcal{I}_0|$$

$$\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu\text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0\}$$

**Setting:** the means for each expert are jointly defined by a parameter  $\alpha$ ,  
 $N = 5$ ,  $\mathcal{I}_0 = \{1, 3, 5\}$ ,  $N_0 = 3$ .



## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\},$

## $(N_0, \Delta_0)$ for Examples

Stochastic-with-a-gap:  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [N]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [N]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [N]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [N]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [N]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [N]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .
- By design,  $N_0 = 1$  and  $\Delta_0 = \Delta$ .

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .
- By design,  $N_0 = 1$  and  $\Delta_0 = \Delta$ .

**Neighborhood-of-I.I.D.**

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .
- By design,  $N_0 = 1$  and  $\Delta_0 = \Delta$ .

**Neighborhood-of-I.I.D.**

- Pick any distribution  $\mu_0$ , and any radius,  $r \geq 0$ .  $\mathcal{D} = \text{Ball}(\mu_0, r)$

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .
- By design,  $N_0 = 1$  and  $\Delta_0 = \Delta$ .

**Neighborhood-of-I.I.D.**

- Pick any distribution  $\mu_0$ , and any radius,  $r \geq 0$ .  $\mathcal{D} = \text{Ball}(\mu_0, r)$
- Suppose that  $\mu_0$  has a gaps between all the mean losses.

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .
- By design,  $N_0 = 1$  and  $\Delta_0 = \Delta$ .

**Neighborhood-of-I.I.D.**

- Pick any distribution  $\mu_0$ , and any radius,  $r \geq 0$ .  $\mathcal{D} = \text{Ball}(\mu_0, r)$
- Suppose that  $\mu_0$  has a gaps between all the mean losses.
- $N_0$ ,  $\Delta_0$  depend on the radius of the ball...

## $(N_0, \Delta_0)$ for Examples

**Stochastic-with-a-gap:**  $\mathcal{D} = \{\mu_0\}$ ,

- $N_0 = 1$ ,  $\mathcal{I}_0 = \left\{ i^* = \arg \min_{i \in [M]} \mathbb{E}_{\mu_0}[\ell_i] \right\}$ ,  $\Delta_0 = \min_{i \neq i^*} \mathbb{E}_{\mu_0}[\ell_i - \ell_{i^*}]$

**Adversarial:**  $\mathcal{D} = \mathcal{M}(\hat{\mathcal{Y}}^N \times \mathcal{Y})$

- $N_0 = N$ ,  $\Delta_0 = +\infty$

**Adversarial-with-an- $\mathbb{E}$ -gap**

- All measures where a common expert is better than others in  $\mathbb{E}$  by  $\Delta > 0$ .
- By design,  $N_0 = 1$  and  $\Delta_0 = \Delta$ .

**Neighborhood-of-I.I.D.**

- Pick any distribution  $\mu_0$ , and any radius,  $r \geq 0$ .  $\mathcal{D} = \text{Ball}(\mu_0, r)$
- Suppose that  $\mu_0$  has a gaps between all the mean losses.
- $N_0$ ,  $\Delta_0$  depend on the radius of the ball...

## Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

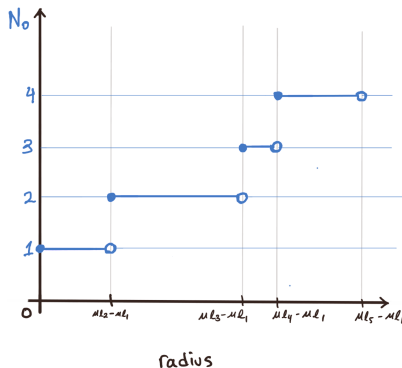
$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$



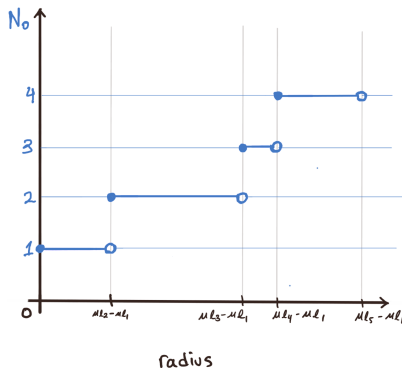
# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius



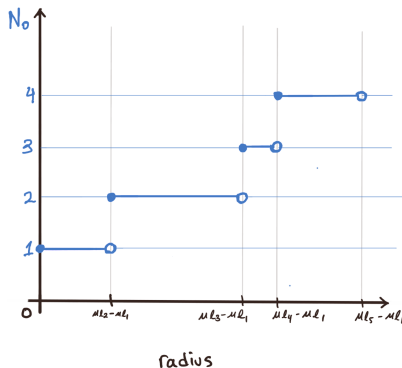
# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely



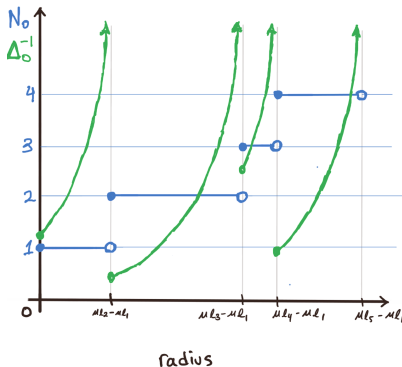
# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely



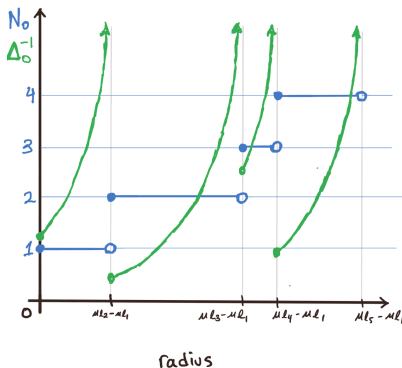
# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely
- $\Delta_0^{-1}$  increases between jumps in  $N_0$



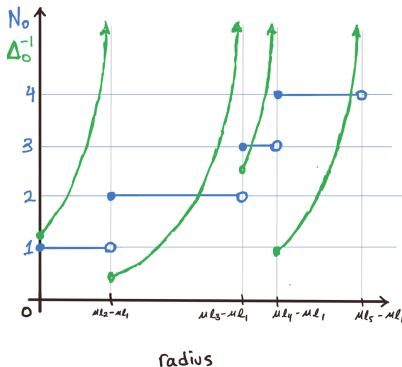
# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely
- $\Delta_0^{-1}$  increases between jumps in  $N_0$
- $\Delta_0^{-1}$  resets each time  $N_0$  increases



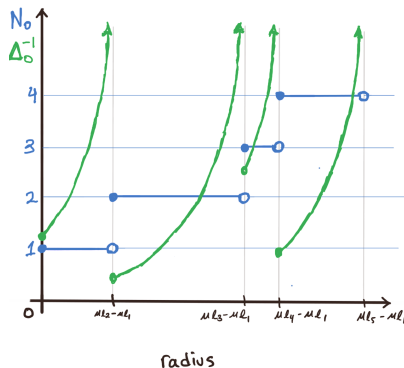
# Interpreting $(N_0, \Delta_0^{-1})$

Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely
- $\Delta_0^{-1}$  increases between jumps in  $N_0$
- $\Delta_0^{-1}$  resets each time  $N_0$  increases



Lexicographical order on  $(N_0, \Delta_0^{-1})$  respects " $\subseteq$ " for nested  $\mathcal{D}$ s.

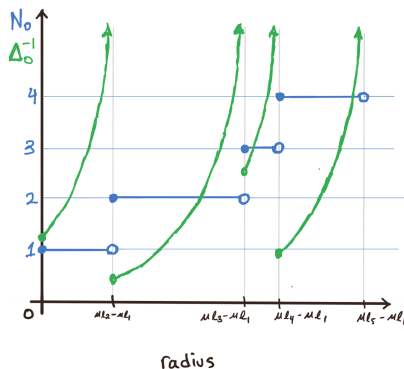
# Interpreting $(N_0, \Delta_0^{-1})$

## Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely
- $\Delta_0^{-1}$  increases between jumps in  $N_0$
- $\Delta_0^{-1}$  resets each time  $N_0$  increases



Lexicographical order on  $(N_0, \Delta_0^{-1})$  respects " $\subseteq$ " for nested  $\mathcal{D}$ s.

- For nested  $\mathcal{D}$ s, larger one is harder.

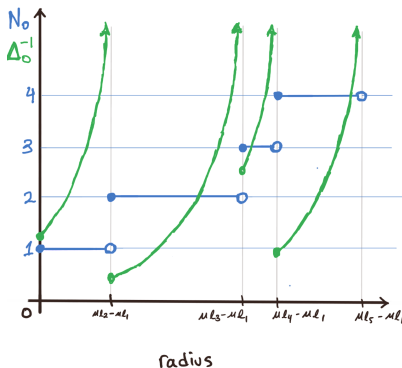
# Interpreting $(N_0, \Delta_0^{-1})$

## Minimax Regret

$$\mathbb{E}R(T) \asymp \begin{cases} \sqrt{T \log N_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

$\mathcal{D} = \text{Ball}(\mu, \text{radius})$  w/  $\mathbb{E}_\mu \ell_1 < \mathbb{E}_\mu \ell_2 < \dots$

- $N_0$  non-decreasing with radius
- $N_0$  increases discretely
- $\Delta_0^{-1}$  increases between jumps in  $N_0$
- $\Delta_0^{-1}$  resets each time  $N_0$  increases



Lexicographical order on  $(N_0, \Delta_0^{-1})$  respects " $\subseteq$ " for nested  $\mathcal{D}$ s.

- For nested  $\mathcal{D}$ s, larger one is harder.
- $(N_0, \Delta_0^{-1})$  quantifies the difficulty of  $\mathcal{D}$

## Performance of Hedge

---

# Hedge Regret Bounds

Consider playing Hedge with  $\eta(t) = c/\sqrt{t}$  for any convex  $\mathcal{D}$ .

Recall:

- $N_0$  is the number of effective experts,
- $\Delta_0$  is the effective stochastic gap.

# Hedge Regret Bounds

Consider playing Hedge with  $\eta(t) = c/\sqrt{t}$  for any convex  $\mathcal{D}$ .

Recall:

- $N_0$  is the number of effective experts,
- $\Delta_0$  is the effective stochastic gap.

## Theorem BNR20

Taking  $c \propto \sqrt{\log N}$  gives

$$\mathbb{E}R(T) \lesssim \begin{cases} \sqrt{T \log N} + \frac{\log N}{\Delta_0} & : N_0 \geq 2 \\ (\log N)/\Delta_0 & : N_0 = 1. \end{cases}$$

Taking  $c \propto 1$  gives

$$\mathbb{E}R(T) \lesssim (\log N_0)\sqrt{T} + \frac{(\log N)^2}{\Delta_0}$$

**We also prove matching lower bounds!**

# Hedge Regret Bounds

Consider playing Hedge with  $\eta(t) = c/\sqrt{t}$  for any convex  $\mathcal{D}$ .

Recall:

- $N_0$  is the number of effective experts,
- $\Delta_0$  is the effective stochastic gap.

## Theorem BNR20

If the player has oracle knowledge of  $N_0 > 1$ , taking  $c \propto \sqrt{\log(N_0)}$  gives

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^2}{(\log N_0) \Delta_0}.$$

# Hedge Regret Bounds

Consider playing Hedge with  $\eta(t) = c/\sqrt{t}$  for any convex  $\mathcal{D}$ .

Recall:

- $N_0$  is the number of effective experts,
- $\Delta_0$  is the effective stochastic gap.

## Theorem BNR20

If the player has oracle knowledge of  $N_0 > 1$ , taking  $c \propto \sqrt{\log(N_0)}$  gives

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^2}{(\log N_0) \Delta_0}.$$

In all three cases, we interpret terms involving...

# Hedge Regret Bounds

Consider playing Hedge with  $\eta(t) = c/\sqrt{t}$  for any convex  $\mathcal{D}$ .

Recall:

- $N_0$  is the number of effective experts,
- $\Delta_0$  is the effective stochastic gap.

## Theorem BNR20

If the player has oracle knowledge of  $N_0 > 1$ , taking  $c \propto \sqrt{\log(N_0)}$  gives

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^2}{(\log N_0) \Delta_0}.$$

In all three cases, we interpret terms involving...

- $T$ : long run regret accumulation after adapting

# Hedge Regret Bounds

Consider playing Hedge with  $\eta(t) = c/\sqrt{t}$  for any convex  $\mathcal{D}$ .

Recall:

- $N_0$  is the number of effective experts,
- $\Delta_0$  is the effective stochastic gap.

## Theorem BNR20

If the player has oracle knowledge of  $N_0 > 1$ , taking  $c \propto \sqrt{\log(N_0)}$  gives

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^2}{(\log N_0) \Delta_0}.$$

In all three cases, we interpret terms involving...

- $T$ : long run regret accumulation after adapting
- $(\log N, \Delta_0)$ : adversarial regret over adaptation period of duration  $\mathcal{O}\left(\frac{(\log N)^2}{c^2 \Delta_0^2}\right)$

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

**Answer:** Yes!

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

**Answer:** Yes!

We introduce two new algorithms in order to do this...

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

**Answer:** Yes!

We introduce two new algorithms in order to do this...

- FTRL-CARE, accomplished 1 and 3, but not 2.

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

**Answer:** Yes!

We introduce two new algorithms in order to do this...

- FTRL-CARE, accomplished 1 and 3, but not 2.
  - slightly worse dependence on  $N$ .

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

**Answer:** Yes!

We introduce two new algorithms in order to do this...

- FTRL-CARE, accomplished 1 and 3, but not 2.
  - slightly worse dependence on  $N$ .
- Meta-CARE, accomplished all 3 by *boosting* FTRL-CARE with Hedge.

# Can we do better than Hedge?

**Question:** If we don't know  $N_0$ , can we learn adaptively and minimax optimally?

In particular, is there an algorithm for which...

1.  $(T, N_0)$ -dependence matches Hedge with oracle knowledge of  $N_0$ ,
2.  $(\log N, \Delta_0)$ -dependence optimal for the stochastic case when  $N_0 = 1$ , and
3. no information is needed about the true setting?

**Answer:** Yes!

We introduce two new algorithms in order to do this...

- FTRL-CARE, accomplished 1 and 3, but not 2.
  - slightly worse dependence on  $N$ .
- Meta-CARE, accomplished all 3 by *boosting* FTRL-CARE with Hedge.

# Improved Algorithms and Bounds

---

# Intuition for Improving on Hedge

Three Key Insights:

# Intuition for Improving on Hedge

## Three Key Insights:

1. From our oracle Hedge bound, we want a learning rate  $\propto \sqrt{(\log N_0)/t}$ .

# Intuition for Improving on Hedge

## Three Key Insights:

1. From our oracle Hedge bound, we want a learning rate  $\propto \sqrt{(\log N_0)/t}$ .
2. The regret of Hedge closely depends on the *entropy* of the weights:

$$H(\mathbf{w}) = - \sum_{i \in [N]} w_i \log(w_i).$$

# Intuition for Improving on Hedge

## Three Key Insights:

1. From our oracle Hedge bound, we want a learning rate  $\propto \sqrt{(\log N_0)/t}$ .
2. The regret of Hedge closely depends on the *entropy* of the weights:

$$H(\mathbf{w}) = - \sum_{i \in [N]} w_i \log(w_i).$$

3. Worst-case adversary forces weights to concentrate to  $\text{Unif}(\mathcal{I}_0)$ , so

$$H(\mathbf{w}) \approx \log N_0.$$

# Intuition for Improving on Hedge

## Three Key Insights:

1. From our oracle Hedge bound, we want a learning rate  $\propto \sqrt{(\log N_0)/t}$ .
2. The regret of Hedge closely depends on the *entropy* of the weights:

$$H(\mathbf{w}) = - \sum_{i \in [N]} w_i \log(w_i).$$

3. Worst-case adversary forces weights to concentrate to  $\text{Unif}(\mathcal{I}_0)$ , so

$$H(\mathbf{w}) \approx \log N_0.$$

What if we could have our learning rate at time  $t$ ,  $\eta(t)$ , look like

$$\eta(t) = \sqrt{\frac{H(\mathbf{w}(t))}{t}} \quad ?$$

# Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.

# Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.

Parametrized by a sequence of regularizers  $(\psi_t)_{t \in \mathbb{N}} \subseteq \text{simp}([N]) \rightarrow \mathbb{R}$ ,

# Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.

Parametrized by a sequence of regularizers  $(\psi_t)_{t \in \mathbb{N}} \subseteq \text{simp}([N]) \rightarrow \mathbb{R}$ ,

$$w(t+1) = \arg \min_{w \in \text{simp}([N])} (\langle w, L(t) \rangle + \psi_{t+1}(w)).$$

# Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.

Parametrized by a sequence of regularizers  $(\psi_t)_{t \in \mathbb{N}} \subseteq \text{simp}([N]) \rightarrow \mathbb{R}$ ,

$$w(t+1) = \arg \min_{w \in \text{simp}([N])} (\langle w, L(t) \rangle + \psi_{t+1}(w)).$$

Hedge corresponds to  $\psi_{t+1}(w) = -\frac{H(w)}{\eta(t+1)}$ .

# Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.

Parametrized by a sequence of regularizers  $(\psi_t)_{t \in \mathbb{N}} \subseteq \text{simp}([N]) \rightarrow \mathbb{R}$ ,

$$w(t+1) = \arg \min_{w \in \text{simp}([N])} (\langle w, L(t) \rangle + \psi_{t+1}(w)).$$

Hedge corresponds to  $\psi_{t+1}(w) = -\frac{H(w)}{\eta(t+1)}$ . That is,

$$\frac{\exp \{-\eta(t+1)L(t)\}}{\sum_{i \in [N]} \exp \{-\eta(t+1)L_i(t)\}} = \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t) \rangle - \frac{H(w)}{\eta(t+1)} \right)$$

# Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.

Parametrized by a sequence of regularizers  $(\psi_t)_{t \in \mathbb{N}} \subseteq \text{simp}([N]) \rightarrow \mathbb{R}$ ,

$$w(t+1) = \arg \min_{w \in \text{simp}([N])} (\langle w, L(t) \rangle + \psi_{t+1}(w)).$$

Hedge corresponds to  $\psi_{t+1}(w) = -\frac{H(w)}{\eta(t+1)}$ . That is,

$$\frac{\exp \{-\eta(t+1)L(t)\}}{\sum_{i \in [N]} \exp \{-\eta(t+1)L_i(t)\}} = \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t) \rangle - \frac{H(w)}{\eta(t+1)} \right)$$

# Follow the Regularized Leader with CARE

Introducing **FTRL-CARE**

# Follow the Regularized Leader with CARE

## Introducing **FTRL-CARE**

Follow the Regularized Leader with Constraint-Adaptive Root-Entropic regularization

$$w(t+1) \in \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t) \rangle - \frac{\sqrt{t+1}}{c_1} \sqrt{H(w) + c_2} \right),$$

# Follow the Regularized Leader with CARE

## Introducing **FTRL-CARE**

Follow the Regularized Leader with Constraint-Adaptive Root-Entropic regularization

$$w(t+1) \in \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t) \rangle - \frac{\sqrt{t+1}}{c_1} \sqrt{H(w) + c_2} \right),$$

which is equivalent to solving the system of equations...

$$\eta(t+1) = c_1 \sqrt{\frac{H(w(t+1)) + c_2}{t+1}} \quad \text{and} \quad w(t+1) = \frac{\exp \{-\eta(t+1)L(t)\}}{\sum_{i \in [N]} \exp \{-\eta(t+1)L_i(t)\}}.$$

# Follow the Regularized Leader with CARE

## Introducing FTRL-CARE

Follow the Regularized Leader with Constraint-Adaptive Root-Entropic regularization

$$w(t+1) \in \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t) \rangle - \frac{\sqrt{t+1}}{c_1} \sqrt{H(w) + c_2} \right),$$

which is equivalent to solving the system of equations...

$$\eta(t+1) = c_1 \sqrt{\frac{H(w(t+1)) + c_2}{t+1}} \quad \text{and} \quad w(t+1) = \frac{\exp \{-\eta(t+1)L(t)\}}{\sum_{i \in [M]} \exp \{-\eta(t+1)L_i(t)\}}.$$

## Theorem BNR20

For any convex  $\mathcal{D}$ , FTRL-CARE achieves

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^{3/2}}{\Delta_0}.$$

# Follow the Regularized Leader with CARE

## Introducing FTRL-CARE

Follow the Regularized Leader with Constraint-Adaptive Root-Entropic regularization

$$w(t+1) \in \arg \min_{w \in \text{simp}([N])} \left( \langle w, L(t) \rangle - \frac{\sqrt{t+1}}{c_1} \sqrt{H(w) + c_2} \right),$$

which is equivalent to solving the system of equations...

$$\eta(t+1) = c_1 \sqrt{\frac{H(w(t+1)) + c_2}{t+1}} \quad \text{and} \quad w(t+1) = \frac{\exp \{-\eta(t+1)L(t)\}}{\sum_{i \in [M]} \exp \{-\eta(t+1)L_i(t)\}}.$$

## Theorem BNR20

For any convex  $\mathcal{D}$ , FTRL-CARE achieves

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^{3/2}}{\Delta_0}.$$

**CARE if you can, Hedge if you must; or, Meta-CARE for All**

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

## CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

## Meta-CARE

- Treat the predictions of Hedge and FTRL-CARE as *meta-experts*...

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

## Meta-CARE

- Treat the predictions of Hedge and FTRL-CARE as *meta-experts*...
- Use Hedge on these two meta-experts.

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

## Meta-CARE

- Treat the predictions of Hedge and FTRL-CARE as *meta-experts*...
- Use Hedge on these two meta-experts.
- Incur best regret of the two, plus some excess from meta-learning.

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

## Meta-CARE

- Treat the predictions of Hedge and FTRL-CARE as *meta-experts*...
- Use Hedge on these two meta-experts.
- Incur best regret of the two, plus some excess from meta-learning.
- Excess regret from meta-learning does not affect the order.

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

## Meta-CARE

- Treat the predictions of Hedge and FTRL-CARE as *meta-experts*...
- Use Hedge on these two meta-experts.
- Incur best regret of the two, plus some excess from meta-learning.
- Excess regret from meta-learning does not affect the order.

## Theorem BNR20

For any convex  $\mathcal{D}$ , Meta-CARE achieves

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \mathbb{I}_{[N_0=1]} \frac{\log N}{\Delta_0} + \mathbb{I}_{[N_0 \geq 2]} \frac{(\log N)^{3/2}}{\Delta_0}.$$

# CARE if you can, Hedge if you must; or, Meta-CARE for All

FTRL-CARE has adaptively minimax optimal dependence on  $(T, N_0)$ ...

... but when  $N_0 = 1$ , it incurs total regret of order  $\frac{(\log N)^{3/2}}{\Delta_0}$  instead of  $\frac{(\log N)}{\Delta_0}$ .

To be minimax optimal even when  $N_0 = 1$ , combine Hedge and FTRL-CARE.

## Meta-CARE

- Treat the predictions of Hedge and FTRL-CARE as *meta-experts*...
- Use Hedge on these two meta-experts.
- Incur best regret of the two, plus some excess from meta-learning.
- Excess regret from meta-learning does not affect the order.

## Theorem BNR20

For any convex  $\mathcal{D}$ , Meta-CARE achieves

$$\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \mathbb{I}_{[N_0=1]} \frac{\log N}{\Delta_0} + \mathbb{I}_{[N_0 \geq 2]} \frac{(\log N)^{3/2}}{\Delta_0}.$$

# Summary

---

# Our Contributions

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.
4. Provided a new algorithm, Meta-CARE, and corresponding regret bounds.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.
4. Provided a new algorithm, Meta-CARE, and corresponding regret bounds.
  - Adapts optimally to our full spectrum of relaxations of the I.I.D. assumption.

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.
4. Provided a new algorithm, Meta-CARE, and corresponding regret bounds.
  - Adapts optimally to our full spectrum of relaxations of the I.I.D. assumption.
  - ...without using oracle knowledge of  $N_0$ .

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.
4. Provided a new algorithm, Meta-CARE, and corresponding regret bounds.
  - Adapts optimally to our full spectrum of relaxations of the I.I.D. assumption.
  - ...without using oracle knowledge of  $N_0$ .

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.
4. Provided a new algorithm, Meta-CARE, and corresponding regret bounds.
  - Adapts optimally to our full spectrum of relaxations of the I.I.D. assumption.
  - ...without using oracle knowledge of  $N_0$ .

# Our Contributions

1. Introduced a spectrum of relaxations of the I.I.D. assumption.
  - Indexed by time-homogeneous convex constraints on the environment.
  - Interpolate between the pure stochastic and adversarial settings.
  - Data that we want to predict won't be purely adversarial or stochastic.
  - We want to know that we do well in intermediate scenarios as well.
2. Characterized the difficulty of learning along the spectrum using  $N_0$  and  $\Delta_0$ .
  - Defined what it means to be adaptively minimax optimal along the spectrum.
3. Derived regret bounds for Hedge along spectrum from I.I.D. to adversarial.
  - In terms of the constraint  $\mathcal{D}$  via explicit dependence on  $(N_0, \Delta_0)$ .
  - Requires oracle knowledge to get minimax optimal dependence on  $T$  and  $N_0$ .
  - Therefore Hedge is not adaptively minimax optimal.
4. Provided a new algorithm, Meta-CARE, and corresponding regret bounds.
  - Adapts optimally to our full spectrum of relaxations of the I.I.D. assumption.
  - ...without using oracle knowledge of  $N_0$ .

# References

- ▶ [CL06] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- ▶ [FS97] Y. Freund and R. Schapire. “A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting”. *Journal of Computer and System Sciences* 55 (1 1997), pp. 119–139.
- ▶ [GSE14] P. Gaillard, G. Stoltz, and T. van Erven. “A second-order bound with excess losses”. In: *Proceedings of the 27th Conference on Learning Theory*. 2014.
- ▶ [MG19] J. Mourtada and S. Gaïffas. “On the optimality of the Hedge algorithm in the stochastic regime.”. *Journal of Machine Learning Research* 20.83 (2019), pp. 1–28.
- ▶ [Vov98] V. Vovk. “A Game of Prediction with Expert Advice”. *Journal of Computer and System Sciences* 56 (2 1998), pp. 153–173.