Relaxing the I.I.D. Assumption

Adaptively Minimax Optimal Regret via Root-Entropic Regularization

Blair Bilodeau^{*,†,1}, Jeffrey Negrea^{*,1}, Daniel M. Roy^{†,1}

March 30, 2021 UCLA Computer Science Weekly Seminar

*Equal Contribution [†]Presenting ¹University of Toronto and Vector Institute

Background

Trading Stocks

Trading Stocks

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

Trading Stocks

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

Trading Stocks

- 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

Trading Stocks

- 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?

Trading Stocks

- 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)

Trading Stocks

- 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.

Trading Stocks

- 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?

Trading Stocks

- 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?

Trading Stocks

- 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.

Trading Stocks

- 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.

Trading Stocks

- 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.

Trading Stocks

- 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 a.k.a. Online Learning

Sequential Prediction a.k.a. Online Learning

Sequential Prediction a.k.a. Online Learning For rounds t = 1, ..., T:

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

Sequential Prediction a.k.a. Online Learning For rounds t = 1, ..., 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 a.k.a. Online Learning For rounds t = 1, ..., 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 $\ell(\hat{y}(t), y(t))$

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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

- Receive $x(t) = (x_1(t), \dots, x_N(t)) \subseteq \hat{\mathcal{Y}}$ 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 $\ell(\hat{y}(t), y(t))$



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

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

• Relative to the class of *N* reference *experts*;

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;

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;

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 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;

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:

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

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;

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 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;

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.)

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;

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

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

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

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

A constructive algorithm achieves the minimax regret: $\mathbb{E}R(T) \simeq \frac{\log N}{\Lambda}$, uniformly bounded in T.

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

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

A constructive algorithm achieves the minimax regret: $\mathbb{E}R(T) \simeq \frac{\log N}{\Delta}$, uniformly bounded in T.

Adversarial

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

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

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

A constructive algorithm achieves the minimax regret: $\mathbb{E}R(T) \simeq \frac{\log N}{\Delta}$, 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: $\mathbb{E}R(T) \simeq \sqrt{T \log N}$ for all T.

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

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

A constructive algorithm achieves the minimax regret: $\mathbb{E}R(T) \simeq \frac{\log N}{\Delta}$, 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?

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

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

A constructive algorithm achieves the minimax regret: $\mathbb{E}R(T) \simeq \frac{\log N}{\Delta}$, 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!



Simultaneous Optimality of Hedge

Can a single algorithm be optimal in both settings simultaneously? Yes!



Proposition: [MG19]

Hedge is optimal in both the stochastic and adversarial settings.

Stochastic: $\mathbb{E}R(T) \simeq (\log N)/\Delta$ uniformly in *T*.

Adversarial: $\mathbb{E}R(T) \simeq \sqrt{T \log N}$.

Simultaneous Optimality of Hedge

Can a single algorithm be optimal in both settings simultaneously? Yes!



Proposition: [MG19]

Hedge is optimal in both the stochastic and adversarial settings.

Stochastic: $\mathbb{E}R(T) \simeq (\log N)/\Delta$ uniformly in *T*.

Adversarial: $\mathbb{E}R(T) \simeq \sqrt{T \log N}$.

Real data $\not\equiv$ stochastic.

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial.

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial. \leftarrow Too pessimistic.

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial. \leftarrow Too pessimistic.

Building upon [RST11], we study a spectrum between stochastic and adversarial.

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial. \leftarrow Too pessimistic.

Building upon [RST11], we study a spectrum between stochastic and adversarial. Intuitively, fix a "neighbourhood" of distributions.

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial. \leftarrow Too pessimistic.

Building upon [RST11], we study a spectrum between stochastic and adversarial. Intuitively, fix a "neighbourhood" of distributions. Each data point drawn from an arbitrary distribution in "neighbourhood".

Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial. \leftarrow Too pessimistic.

Building upon [RST11], we study a spectrum between stochastic and adversarial. Intuitively, fix a "neighbourhood" of distributions.

Each data point drawn from an arbitrary distribution in "neighbourhood".



Real data $\not\equiv$ stochastic. \leftarrow Too optimistic.

Real data $\not\equiv$ adversarial. \leftarrow Too pessimistic.

Building upon [RST11], we study a spectrum between stochastic and adversarial. Intuitively, fix a "neighbourhood" of distributions.

Each data point drawn from an arbitrary distribution in "neighbourhood".



Algorithms should be robust to a spectrum of data-generating mechanisms.

Algorithms should be robust to a spectrum 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.

Algorithms should be robust to a spectrum 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?

Algorithms should be robust to a spectrum 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?

• Given a spectrum of settings Λ ...

Algorithms should be robust to a spectrum 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?

- Given a spectrum of settings Λ ...
- If we know $\theta \in \Lambda$ in advance, the best achievable performance is $R^*_{\theta}(T)$.

Algorithms should be robust to a spectrum 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?

- Given a spectrum of settings Λ...
- If we know $\theta \in \Lambda$ in advance, the best achievable performance is $R^*_{\theta}(T)$.
- We want an algorithm that does as well as possible without knowing θ :

 $R_{\theta}(T) \leq C R_{\theta}^{*}(T)$ uniformly in θ for large enough T.

Algorithms should be robust to a spectrum 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?

- Given a spectrum of settings Λ...
- If we know $\theta \in \Lambda$ in advance, the best achievable performance is $R^*_{\theta}(T)$.
- We want an algorithm that does as well as possible without knowing θ :

 $R_{\theta}(T) \leq C R_{\theta}^{*}(T)$ uniformly in θ for large enough T.

We show Hedge is suboptimal between Stochastic and Adversarial.

We show Hedge is suboptimal between Stochastic and Adversarial. This was surprising for us.

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.

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 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...

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!

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.

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.
Motivating Intuition

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

- 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, ...

- 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"

- 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)

- 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 Δ_0 , but not which it is...

- 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 Δ_0 , but not which it is...
 - then we are *almost* in the stochastic-with-a-gap case

- 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 Δ_0 , but not which it is...
 - then we are *almost* in the stochastic-with-a-gap case
 - so we might hope for regret Θ((log N)/Δ₀)

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 Δ_0 , but not which it is...
 - then we are *almost* in the stochastic-with-a-gap case
 - so we might hope for regret Θ((log N)/Δ₀)

Theorem BNR20

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

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

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 Δ_0 , but not which it is...
 - then we are *almost* in the stochastic-with-a-gap case
 - so we might hope for regret Θ((log N)/Δ₀)

Theorem BNR20

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

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

Relaxing the I.I.D. Assumption

Intuition

Experts and environment may collude.

Intuition

Experts and environment may collude.

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

Intuition

Experts and environment may collude.

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

Formal Framework

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

Experts and environment 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.
 - The choice of distribution is made based on outcomes of the previous rounds.

Experts and environment 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.
 - The choice of distribution is made based on outcomes of the previous rounds.

More Details

• Time-Homogeneous: \mathcal{D} does not depend on t.

Experts and environment 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.
 - The choice of distribution is made based on outcomes of the previous rounds.

More Details

- Time-Homogeneous: \mathcal{D} does not depend on t.
- Convex: environment can flip a coin to select between basic elements of \mathcal{D} .

Experts and environment 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.
 - The choice of distribution is made based on outcomes of the previous rounds.

More Details

- Time-Homogeneous: \mathcal{D} does not depend on t.
- \bullet Convex: environment can flip a coin to select between basic elements of $\mathcal{D}.$
- Environment may aim to maximize regret subject to the constraint.

Experts and environment 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.
 - The choice of distribution is made based on outcomes of the previous rounds.

More Details

- Time-Homogeneous: \mathcal{D} does not depend on t.
- \bullet Convex: environment can flip a coin to select between basic elements of $\mathcal{D}.$
- Environment may aim to maximize regret subject to the constraint.

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

• simplify the abstract complexity of the constraint;

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- yield matching lower and upper bounds on regret.

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- yield matching lower and upper bounds on regret.

Effective Experts

The set ${\mathcal D}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- yield matching lower and upper bounds on regret.

Effective Experts

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

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- 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.

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- 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

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- 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

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

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- 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.

The set $\ensuremath{\mathcal{D}}$ may be complex and difficult to describe for a particular application.

We want to characterize the constraint using quantities that:

- simplify the abstract complexity of the constraint;
- differentiate whether the data is "easy" or not, independent of algorithms;
- 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.

- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$

- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \left\{ \mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \right\}$

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

- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.



- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.


- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.



- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.



- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.



- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.



- $\mathcal{I}_0 = \{ \text{experts that are optimal for some } \mu \in \mathcal{D} \}$ $N_0 = |\mathcal{I}_0|$
- $\Delta_0 = \inf_{\mu \in \mathcal{D}} \{\mu \text{-expected difference in loss of best expert and best expert not in } \mathcal{I}_0 \}$
- **Setting:** the means for each expert are jointly defined by a parameter α , N = 5, $\mathcal{I}_0 = \{1, 3, 5\}$, $N_0 = 3$.



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

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

• $N_0 = 1$,

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

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

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

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

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

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

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

• $N_0 = N$,

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

• $N_0 = N$, $\Delta_0 = +\infty$

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

- 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 μ_0 , and any radius, $r \ge 0$. $\mathcal{D} = \mathsf{Ball}(\mu_0, r)$

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

- Pick any distribution μ_0 , and any radius, $r \ge 0$. $\mathcal{D} = \text{Ball}(\mu_0, r)$
- Suppose that μ_0 has a gap between each of the mean losses.

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

- Pick any distribution μ_0 , and any radius, $r \ge 0$. $\mathcal{D} = \mathsf{Ball}(\mu_0, r)$
- Suppose that μ_0 has a gap between each of the mean losses.
- N_0 , Δ_0 depend on the radius of the ball...

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

• $N_0 = N$, $\Delta_0 = +\infty$

Adversarial-with-an-E-gap (Mourtada and Gaïffas 2019)

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

- Pick any distribution μ_0 , and any radius, $r \ge 0$. $\mathcal{D} = \mathsf{Ball}(\mu_0, r)$
- Suppose that μ_0 has a gap between each of the mean losses.
- N_0 , Δ_0 depend on the radius of the ball...



radius

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



• N₀ non-decreasing with radius

radius



- N₀ non-decreasing with radius
- N₀ increases discretely

radius



- N₀ non-decreasing with radius
- N₀ increases discretely

radius



- N₀ non-decreasing with radius
- N₀ increases discretely
- Δ_0^{-1} increases between N_0 jumps

radius



- N₀ non-decreasing with radius
- N₀ increases discretely
- Δ_0^{-1} increases between N_0 jumps
- Δ_0^{-1} resets at each jump

radius

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



- N₀ non-decreasing with radius
- N₀ increases discretely
- Δ_0^{-1} increases between N_0 jumps
- Δ_0^{-1} resets at each jump

Lexicographical order on (N_0, Δ_0^{-1})

radius

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



- N₀ non-decreasing with radius
- N₀ increases discretely
- Δ_0^{-1} increases between N_0 jumps
- Δ_0^{-1} resets at each jump

Lexicographical order on (N_0, Δ_0^{-1})

• For nested \mathcal{D} s, larger one is "harder" to learn.

radius

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



- N₀ non-decreasing with radius
- N_0 increases discretely
- Δ_0^{-1} increases between N_0 jumps
- Δ_0^{-1} resets at each jump

Lexicographical order on (N_0, Δ_0^{-1})

- For nested \mathcal{D} s, larger one is "harder" to learn.
- (N₀, Δ₀⁻¹) quantifies the difficulty.

radius

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



- N₀ non-decreasing with radius
- N_0 increases discretely
- Δ_0^{-1} increases between N_0 jumps
- Δ_0^{-1} resets at each jump

Lexicographical order on (N_0, Δ_0^{-1})

- For nested \mathcal{D} s, larger one is "harder" to learn.
- (N₀, Δ₀⁻¹) quantifies the difficulty.
- How does regret change with Δ_0 ?

radius

Impact of (N_0, Δ_0^{-1}) on Regret



Impact of (N_0, Δ_0^{-1}) on Regret


Impact of (N_0, Δ_0^{-1}) on Regret



Impact of (N_0, Δ_0^{-1}) on Regret



Performance of Hedge

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

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:

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.

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*.

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

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

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 $\eta : \mathbb{N} \to \mathbb{R}$; initialize the weights as uniform; define

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

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 $\eta : \mathbb{N} \to \mathbb{R}$; initialize the weights as uniform; define

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

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

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

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 $\eta : \mathbb{N} \to \mathbb{R}$; initialize the weights as uniform; define

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

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

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

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

Recall:

- N_0 is the number of of effective experts,
- Δ_0 is the effective stochastic gap.

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

Recall:

- N_0 is the number of of effective experts,
- Δ_0 is the effective stochastic gap.

Theorem BNR20

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

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

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!

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

Recall:

- N_0 is the number of of effective experts,
- Δ_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

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

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

Recall:

- N_0 is the number of of effective experts,
- Δ_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...

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

Recall:

- N_0 is the number of of effective experts,
- Δ_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

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

Recall:

- N_0 is the number of of effective experts,
- Δ_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 $\mathcal{O}\left(\frac{(\log N)^2}{c^2\Lambda^2}\right)$

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...

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

- 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

- 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?

- 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!

- 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!

- 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.

- 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!

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

- 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!

- 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.

- 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!

- 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

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

- 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:

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

- 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 $\operatorname{Unif}(\mathcal{I}_0)$, so

 $H(w) \approx \log N_0.$

- 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 $\operatorname{Unif}(\mathcal{I}_0)$, so

 $H(w) \approx \log N_0.$

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

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

Follow the Regularized Leader

FTRL is a fundamental online linear optimization algorithm.
FTRL is a fundamental online linear optimization algorithm.

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

FTRL is a fundamental online linear optimization algorithm.

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

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

FTRL is a fundamental online linear optimization algorithm.

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

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

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

FTRL is a fundamental online linear optimization algorithm.

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

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

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

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

FTRL is a fundamental online linear optimization algorithm.

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

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

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

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

Introducing FTRL-CARE:

FTRL is a fundamental online linear optimization algorithm.

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

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

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

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

Introducing FTRL-CARE:

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

$$w(t+1) \in \operatorname*{arg\,min}_{w\in \operatorname{simp}([N])} \left(\left\langle w, \ L(t) \right\rangle - rac{\sqrt{t+1}}{c_1} \sqrt{H(w) + c_2}
ight).$$

FTRL is a fundamental online linear optimization algorithm.

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

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

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

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

Introducing FTRL-CARE:

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

$$w(t+1) \in \operatorname*{arg\,min}_{w\in \operatorname{simp}([N])} \left(\left\langle w, \ L(t) \right\rangle - rac{\sqrt{t+1}}{c_1} \sqrt{H(w) + c_2}
ight).$$

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20For any convex \mathcal{D} , FTRL-CARE achieves $\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^{3/2}}{\Delta_0}$.When $N_0 = 1$, it incurs total regret of order $\frac{(\log N)^{3/2}}{\Delta_0}$ instead of $\frac{(\log N)}{\Delta_0}$.

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20For any convex \mathcal{D} , FTRL-CARE achieves $\mathbb{E}R(T) \lesssim \sqrt{T \log N_0} + \frac{(\log N)^{3/2}}{\Delta_0}$ 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.

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20

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

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

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*.

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20

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

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

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.

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20

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

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

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 negligible excess from meta-learning.

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20

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

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

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 negligible excess from meta-learning.

Theorem BNR20

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

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

FTRL-CARE is *almost* adaptively minimax optimal.

Theorem BNR20

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

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

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 negligible excess from meta-learning.

Theorem BNR20

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

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

Proof Technique 1: FTRL-CARE looks like Hedge with oracle knowledge.

Proof Technique 1: FTRL-CARE looks like Hedge with oracle knowledge.

Lemma BNR20

FTRL-CARE is equivalent to solving the following system of equations:

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

Proof Technique 1: FTRL-CARE looks like Hedge with oracle knowledge.

Lemma BNR20

FTRL-CARE is equivalent to solving the following system of equations:

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

Proof Technique 2: Concentration of measure holds under our relaxation of i.i.d.

Proof Technique 1: FTRL-CARE looks like Hedge with oracle knowledge.

Lemma BNR20

FTRL-CARE is equivalent to solving the following system of equations:

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

Proof Technique 2: Concentration of measure holds under our relaxation of i.i.d.

Lemma BNR20

For any prediction algorithm, constraint \mathcal{D}_{r} and data-generating mechanism,

$$\sup_{i \in [N] \setminus \mathcal{I}_0} \mathbb{E} \min_{i_0 \in \mathcal{I}_0} \exp \left\{ \lambda \sum_{t=0}^{T} \left[\ell_{i_0}(t) - \ell_i(t) \right] \right\} \leq \exp \left\{ T \left[\lambda^2 / 2 - \lambda \Delta_0 \right] \right\}.$$

Proof Technique 1: FTRL-CARE looks like Hedge with oracle knowledge.

Lemma BNR20

FTRL-CARE is equivalent to solving the following system of equations:

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

Proof Technique 2: Concentration of measure holds under our relaxation of i.i.d.

Lemma BNR20

For any prediction algorithm, constraint \mathcal{D}_{r} and data-generating mechanism,

$$\sup_{i \in [N] \setminus \mathcal{I}_0} \mathbb{E} \min_{i_0 \in \mathcal{I}_0} \exp \left\{ \lambda \sum_{t=0}^{T} \left[\ell_{i_0}(t) - \ell_i(t) \right] \right\} \leq \exp \left\{ T \left[\lambda^2 / 2 - \lambda \Delta_0 \right] \right\}.$$

Summary

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

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - Data that we want to predict won't be purely adversarial or stochastic.

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

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N_0 , and the effective stochastic gap, Δ_0 .

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N₀, and the effective stochastic gap, Δ₀.
- 3. Formalized the notion of adaptive minimax optimality.

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N_0 , and the effective stochastic gap, Δ_0 .
- 3. Formalized the notion of adaptive minimax optimality.
- 4. Proved Hedge **is not** adaptively minimax optimal along spectrum from I.I.D. to adversarial.

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N_0 , and the effective stochastic gap, Δ_0 .
- 3. Formalized the notion of adaptive minimax optimality.
- 4. Proved Hedge **is not** adaptively minimax optimal along spectrum from I.I.D. to adversarial.
 - Requires oracle knowledge to get minimax optimal dependence on T and N_0 .

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N_0 , and the effective stochastic gap, Δ_0 .
- 3. Formalized the notion of adaptive minimax optimality.
- 4. Proved Hedge **is not** adaptively minimax optimal along spectrum from I.I.D. to adversarial.
 - Requires oracle knowledge to get minimax optimal dependence on T and N_0 .
- 5. Provided a new algorithm, Meta-CARE that is adaptively minimax optimal.

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N_0 , and the effective stochastic gap, Δ_0 .
- 3. Formalized the notion of adaptive minimax optimality.
- 4. Proved Hedge **is not** adaptively minimax optimal along spectrum from I.I.D. to adversarial.
 - Requires oracle knowledge to get minimax optimal dependence on T and N_0 .
- 5. Provided a new algorithm, Meta-CARE that is adaptively minimax optimal.
 - Performs as well as possible relative to the constraint on the adversary,

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N₀, and the effective stochastic gap, Δ₀.
- 3. Formalized the notion of adaptive minimax optimality.
- 4. Proved Hedge **is not** adaptively minimax optimal along spectrum from I.I.D. to adversarial.
 - Requires oracle knowledge to get minimax optimal dependence on T and N_0 .
- 5. Provided a new algorithm, Meta-CARE that is adaptively minimax optimal.
 - Performs as well as possible relative to the constraint on the adversary, without knowledge of the constraint.

- 1. Introduced a spectrum of relaxations of the I.I.D. assumption.
 - 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 minimax regret under time-homogeneous convex constraints.
 - Depends on the number of effective experts, N₀, and the effective stochastic gap, Δ₀.
- 3. Formalized the notion of adaptive minimax optimality.
- 4. Proved Hedge **is not** adaptively minimax optimal along spectrum from I.I.D. to adversarial.
 - Requires oracle knowledge to get minimax optimal dependence on T and N_0 .
- 5. Provided a new algorithm, Meta-CARE that is adaptively minimax optimal.
 - Performs as well as possible relative to the constraint on the adversary, without knowledge of the constraint.

References

- ► [CL06] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games.* Cambirdge 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.
- [RST11] A. Rakhlin, K. Sridharan, and A. Tewari. "Online learning: Stochastic, constrained, and smoothed adversaries". In: Advances in Neural Information Processing Systems 25. 2011.
References (cont.)

[Vov98] V. Vovk. "A Game of Prediction with Expert Advice". Journal of Computer and System Sciences 56 (2 1998), pp. 153–173.