

Regret, Games, and Boosting

a Reverie of Gambles and Bounds



Nicolò Cesa-Bianchi

Università degli Studi di Milano
Politecnico di Milano

Joint work with...



Marco Bressan
UNIMI



Nataly Brukhim
IAS Princeton



Emmanuel Esposito
UNIMI & IIT



Yishay Mansour
Tel Aviv & Google



Shay Moran
Technion & Google



Max Thiessen
TU Wien

Learning a binary classifier

- Finite sample \mathcal{X} of datapoints with binary labels $f : \mathcal{X} \rightarrow \{-1, 1\}$



Learning a binary classifier

- ▶ Finite sample \mathcal{X} of datapoints with binary labels $f : \mathcal{X} \rightarrow \{-1, 1\}$
- ▶ We know that *simple* explanations of (\mathcal{X}, f) have good predictive power



Learning a binary classifier

- ▶ Finite sample \mathcal{X} of datapoints with binary labels $f : \mathcal{X} \rightarrow \{-1, 1\}$
- ▶ We know that *simple* explanations of (\mathcal{X}, f) have good predictive power
- ▶ Fix a simple (e.g., low VC-dimension) class of $\{-1, 1\}$ -valued functions



Learning a binary classifier

- ▶ Finite sample \mathcal{X} of datapoints with binary labels $f : \mathcal{X} \rightarrow \{-1, 1\}$
- ▶ We know that *simple* explanations of (\mathcal{X}, f) have good predictive power
- ▶ Fix a simple (e.g., low VC-dimension) class of $\{-1, 1\}$ -valued functions
- ▶ **Simple explanations:** (convex) combination of functions in the class whose sign correlates well with f on the sample \mathcal{X}

Weak Learning

- ▶ Let \mathcal{H} be the projection on \mathcal{X} of the functions in our VC class
- ▶ Each $h \in \mathcal{H}$ has the form $h : \mathcal{X} \rightarrow \{-1, 1\}$



Weak Learning

- ▶ Let \mathcal{H} be the projection on \mathcal{X} of the functions in our VC class
- ▶ Each $h \in \mathcal{H}$ has the form $h : \mathcal{X} \rightarrow \{-1, 1\}$
- ▶ **WL (weak learning) assumption:** There exists $\gamma > 0$ such that

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad \mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) \leq \frac{1}{2} - \gamma$$

Weak Learning

- ▶ Let \mathcal{H} be the projection on \mathcal{X} of the functions in our VC class
- ▶ Each $h \in \mathcal{H}$ has the form $h : \mathcal{X} \rightarrow \{-1, 1\}$
- ▶ **WL (weak learning) assumption:** There exists $\gamma > 0$ such that

$$\forall q \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad \mathbb{P}_{X \sim q}(h(X) \neq f(X)) < \frac{1}{2}$$

Connections with minimax

- ▶ Let M be the $|\mathcal{H}| \times |\mathcal{X}|$ boolean matrix of elements $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$



Connections with minimax

- ▶ Let M be the $|\mathcal{H}| \times |\mathcal{X}|$ boolean matrix of elements $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$
- ▶ Rewrite WL assumption using M :

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad \mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) < \frac{1}{2} \iff \max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) < \frac{1}{2}$$

Connections with minimax

- ▶ Let M be the $|\mathcal{H}| \times |\mathcal{X}|$ boolean matrix of elements $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$
- ▶ Rewrite WL assumption using M :

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad \mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) < \frac{1}{2} \iff \max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) < \frac{1}{2}$$

- ▶ By the **minimax theorem**, we know that

$$\max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) = \min_{\mathbf{p} \in \Delta_{\mathcal{H}}} \max_{x \in \mathcal{X}} M(\mathbf{p}, x)$$

Connections with minimax

- ▶ Let M be the $|\mathcal{H}| \times |\mathcal{X}|$ boolean matrix of elements $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$
- ▶ Rewrite WL assumption using M :

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad \mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) < \frac{1}{2} \iff \max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) < \frac{1}{2}$$

- ▶ By the **minimax theorem**, we know that

$$\max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) = \min_{\mathbf{p} \in \Delta_{\mathcal{H}}} \max_{x \in \mathcal{X}} M(\mathbf{p}, x)$$

- ▶ Therefore there exists a **distribution** p^* over \mathcal{H} such that

$$\max_{x \in \mathcal{X}} M(\mathbf{p}^*, x) < \frac{1}{2} \iff \mathbb{P}_{H \sim p^*}(H(x) \neq f(x)) < \frac{1}{2} \text{ for all } x \in \mathcal{X}$$

Connections with minimax

- ▶ Let M be the $|\mathcal{H}| \times |\mathcal{X}|$ boolean matrix of elements $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$
- ▶ Rewrite WL assumption using M :

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad \mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) < \frac{1}{2} \iff \max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) < \frac{1}{2}$$

- ▶ By the **minimax theorem**, we know that

$$\max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M(h, \mathbf{q}) = \min_{\mathbf{p} \in \Delta_{\mathcal{H}}} \max_{x \in \mathcal{X}} M(\mathbf{p}, x)$$

- ▶ Therefore there exists a **distribution** p^* over \mathcal{H} such that

$$\max_{x \in \mathcal{X}} M(\mathbf{p}^*, x) < \frac{1}{2} \iff \mathbb{P}_{H \sim p^*}(H(x) \neq f(x)) < \frac{1}{2} \text{ for all } x \in \mathcal{X}$$

- ▶ This \mathbf{p}^* *explains* (f, \mathcal{X}) : for all $x \in \mathcal{X}$ $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H(x)]) = f(x)$

where $\text{sgn}(0) = \perp$

Weak learning and simple explanations

- ▶ WL assumption is equivalent to $\text{sgn}(\mathbb{E}_{p^*}[H]) = f$, existence of a simple explanation for (\mathcal{X}, f)



Weak learning and simple explanations

- ▶ WL assumption is equivalent to $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H]) = f$, existence of a simple explanation for (\mathcal{X}, f)
- ▶ We can find \mathbf{p}^* using LP, but **boosting** is a very simple and efficient alternative



Weak learning and simple explanations

- ▶ WL assumption is equivalent to $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H]) = f$, existence of a simple explanation for (\mathcal{X}, f)
- ▶ We can find \mathbf{p}^* using LP, but **boosting** is a very simple and efficient alternative
- ▶ WL oracle for $z \in [0, 1]$: Given $\mathbf{q} \in \Delta_{\mathcal{X}}$ the oracle returns $h \in \mathcal{H}$ such that

$$\mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) \leq z$$

Weak learning and simple explanations

- ▶ WL assumption is equivalent to $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H]) = f$, existence of a simple explanation for (\mathcal{X}, f)
- ▶ We can find \mathbf{p}^* using LP, but **boosting** is a very simple and efficient alternative
- ▶ WL oracle for $z \in [0, 1]$: Given $\mathbf{q} \in \Delta_{\mathcal{X}}$ the oracle returns $h \in \mathcal{H}$ such that

$$\mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) \leq z$$

Definition

$z \in [0, 1]$ is **boostable** if we can find \mathbf{p}^* such that $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H]) = f$ using a WL oracle for z

Weak learning and simple explanations

- ▶ WL assumption is equivalent to $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H]) = f$, existence of a simple explanation for (\mathcal{X}, f)
- ▶ We can find \mathbf{p}^* using LP, but **boosting** is a very simple and efficient alternative
- ▶ WL oracle for $z \in [0, 1]$: Given $\mathbf{q} \in \Delta_{\mathcal{X}}$ the oracle returns $h \in \mathcal{H}$ such that

$$\mathbb{P}_{X \sim \mathbf{q}}(h(X) \neq f(X)) \leq z$$

Definition

$z \in [0, 1]$ is **boostable** if we can find \mathbf{p}^* such that $\text{sgn}(\mathbb{E}_{\mathbf{p}^*}[H]) = f$ using a WL oracle for z

We now show how **to boost a WL oracle using online learning**

Online learning

Recall $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$

The online learning protocol

For each round $t \geq 1$:

1. The learner chooses $\mathbf{p}_t \in \Delta_{\mathcal{H}}$
2. The adversary reveals $x_t \in \mathcal{X}$
3. The learner suffers loss $M(\mathbf{p}_t, x_t) = \mathbb{P}_{H \sim \mathbf{p}_t}(H(x_t) \neq f(x_t))$



Online learning

Recall $M(h, x) = \mathbb{I}\{h(x) \neq f(x)\}$

The online learning protocol

For each round $t \geq 1$:

1. The learner chooses $\mathbf{p}_t \in \Delta_{\mathcal{H}}$
2. The adversary reveals $x_t \in \mathcal{X}$
3. The learner suffers loss $M(\mathbf{p}_t, x_t) = \mathbb{P}_{H \sim \mathbf{p}_t}(H(x_t) \neq f(x_t))$

For all T and for all $x_1, \dots, x_T \in \mathcal{X}$, if the learner runs the **Hedge algorithm**, then

$$\frac{1}{T} \sum_{t=1}^T M(\mathbf{p}_t, x_t) \leq \min_{h \in \mathcal{H}} \frac{1}{T} \sum_{t=1}^T M(h, x_t) + \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

Boosting via online learning

We run Hedge over the **dual game** $M' = \mathbf{1}\mathbf{1}^\top - M^\top$ against a **WL oracle for z** as adversary

Boosting algorithm

For each round $t \geq 1$:

1. Hedge chooses $\mathbf{p}_t \in \Delta_{\mathcal{X}}$ (Hedge learns distributions over \mathcal{X})
2. The WL oracle returns $h_t \in \mathcal{H}$ satisfying $M(h_t, \mathbf{p}_t) \leq z$
3. Hedge gets loss $M'(\mathbf{p}_t, h_t) = 1 - M(h_t, \mathbf{p}_t)$

Analysis

- The WL oracle for z satisfies $M'(\mathbf{p}_t, h_t) = 1 - M(h_t, \mathbf{p}_t) \geq 1 - z \quad t = 1, 2, \dots$



Analysis

- ▶ The WL oracle for z satisfies $M'(\mathbf{p}_t, h_t) = 1 - M(h_t, \mathbf{p}_t) \geq 1 - z \quad t = 1, 2, \dots$
- ▶ Therefore, by the properties of Hedge

$$1 - z \leq \frac{1}{T} \sum_{t=1}^T M'(\mathbf{p}_t, h_t) \leq \min_{x \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T M'(x, h_t) + \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

Analysis

- ▶ The WL oracle for z satisfies $M'(\mathbf{p}_t, h_t) = 1 - M(h_t, \mathbf{p}_t) \geq 1 - z \quad t = 1, 2, \dots$
- ▶ Therefore, by the properties of Hedge

$$1 - z \leq \frac{1}{T} \sum_{t=1}^T M'(\mathbf{p}_t, h_t) \leq \min_{x \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T M'(x, h_t) + \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

- ▶ Hence, if $z < \frac{1}{2}$ and T is large enough,

$$\min_{x \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T M'(x, h_t) > \frac{1}{2}$$

Analysis

- ▶ The WL oracle for z satisfies $M'(\mathbf{p}_t, h_t) = 1 - M(h_t, \mathbf{p}_t) \geq 1 - z \quad t = 1, 2, \dots$
- ▶ Therefore, by the properties of Hedge

$$1 - z \leq \frac{1}{T} \sum_{t=1}^T M'(\mathbf{p}_t, h_t) \leq \min_{x \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T M'(x, h_t) + \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

- ▶ Hence, if $z < \frac{1}{2}$ and T is large enough,

$$\min_{x \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T M'(x, h_t) > \frac{1}{2}$$

- ▶ So, $h_t(x) = f(x)$ for more than half of the h_t on each $x \in \mathcal{X}$, which implies

$$\operatorname{sgn}\left(\frac{1}{T} \sum_{t=1}^T h_t\right) = f$$

A game-theoretic view

- ▶ Game matrix $B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ for binary classification



A game-theoretic view

- ▶ Game matrix $B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ for binary classification
- ▶ $\frac{1}{2} = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} B(\mathbf{u}, \mathbf{v}) = V(B)$ is the value of game B $(\mathcal{Y} = \{-1, 1\})$



A game-theoretic view

- ▶ Game matrix $B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ for binary classification
- ▶ $\frac{1}{2} = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} B(\mathbf{u}, \mathbf{v}) = V(B)$ is the value of game B $(\mathcal{Y} = \{-1, 1\})$

Theorem (Boosting Theorem)

Any $z < V(B)$ is boostable



A game-theoretic view

- ▶ Game matrix $B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ for binary classification
- ▶ $\frac{1}{2} = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} B(\mathbf{u}, \mathbf{v}) = V(B)$ is the value of game B $(\mathcal{Y} = \{-1, 1\})$

Theorem (Boosting Theorem)

Any $z < V(B)$ is boostable

$$V(B) = \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} B(\mathbf{u}, \mathbf{v})$$

is the smallest error attainable with a biased coin \mathbf{u} against the worst-possible distribution \mathbf{v}

A game-theoretic view

- ▶ Game matrix $B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ for binary classification
- ▶ $\frac{1}{2} = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} B(\mathbf{u}, \mathbf{v}) = V(B)$ is the value of game B $(\mathcal{Y} = \{-1, 1\})$

Theorem (Boosting Theorem)

Any $z < V(B)$ is boostable

$$V(B) = \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} B(\mathbf{u}, \mathbf{v})$$

is the smallest error attainable with a biased coin \mathbf{u} against the worst-possible distribution \mathbf{v}

Any $z \in [0, 1]$ is either boostable or coin attainable

Cost-sensitive binary classification

- Game matrix $C = \begin{pmatrix} 0 & c^+ \\ c^- & 0 \end{pmatrix}$ for cost-sensitive classification
where $0 < c^+, c^- \leq 1$ are costs for **false positive** and **false negative** mistakes



Cost-sensitive binary classification

- ▶ Game matrix $C = \begin{pmatrix} 0 & c^+ \\ c^- & 0 \end{pmatrix}$ for cost-sensitive classification
where $0 < c^+, c^- \leq 1$ are costs for **false positive** and **false negative** mistakes
- ▶ $\min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = \frac{c^- c^+}{c^- + c^+} = V(C)$ is the value of game C



Cost-sensitive binary classification

- Game matrix $C = \begin{pmatrix} 0 & c^+ \\ c^- & 0 \end{pmatrix}$ for cost-sensitive classification where $0 < c^+, c^- \leq 1$ are costs for **false positive** and **false negative** mistakes
- $\min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = \frac{c^- c^+}{c^- + c^+} = V(C)$ is the value of game C
- **$(\mathcal{H} \times \mathcal{X})$ -matrix:**

$$M_C(h, x) = c^+ \underbrace{\mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{false positive}} + c^- \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{false negative}}$$

Cost-sensitive binary classification

- Game matrix $C = \begin{pmatrix} 0 & c^+ \\ c^- & 0 \end{pmatrix}$ for cost-sensitive classification where $0 < c^+, c^- \leq 1$ are costs for **false positive** and **false negative** mistakes
- $\min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = \frac{c^- c^+}{c^- + c^+} = V(C)$ is the value of game C
- **($\mathcal{H} \times \mathcal{X}$)-matrix:**

$$M_C(h, x) = c^+ \underbrace{\mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{false positive}} + c^- \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{false negative}}$$

- WL oracle for z : $\max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M_C(h, \mathbf{q}) \leq z$

Cost-sensitive binary classification

- Game matrix $C = \begin{pmatrix} 0 & c^+ \\ c^- & 0 \end{pmatrix}$ for cost-sensitive classification where $0 < c^+, c^- \leq 1$ are costs for **false positive** and **false negative** mistakes
- $\min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = \frac{c^- c^+}{c^- + c^+} = V(C)$ is the value of game C
- **$(\mathcal{H} \times \mathcal{X})$ -matrix:**

$$M_C(h, x) = c^+ \underbrace{\mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{false positive}} + c^- \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{false negative}}$$

- WL oracle for z : $\max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M_C(h, \mathbf{q}) \leq z$
- Which values of z are boostable?

Cost-sensitive binary classification

- Game matrix $C = \begin{pmatrix} 0 & c^+ \\ c^- & 0 \end{pmatrix}$ for cost-sensitive classification where $0 < c^+, c^- \leq 1$ are costs for **false positive** and **false negative** mistakes
- $\min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = \frac{c^- c^+}{c^- + c^+} = V(C)$ is the value of game C
- **$(\mathcal{H} \times \mathcal{X})$ -matrix:**

$$M_C(h, x) = c^+ \underbrace{\mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{false positive}} + c^- \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{false negative}}$$

- WL oracle for z : $\max_{\mathbf{q} \in \Delta_{\mathcal{X}}} \min_{h \in \mathcal{H}} M_C(h, \mathbf{q}) \leq z$
- Which values of z are boostable?

Theorem (Cost-sensitive Boosting Theorem)

Any $z < V(C)$ is boostable

(boostable vs. coin attainable dichotomy)

Relationship to Bayes optimal prediction

- ▶ Bayes optimal prediction minimizes the conditional risk:

$$f^*(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$



Relationship to Bayes optimal prediction

- ▶ Bayes optimal prediction minimizes the conditional risk:

$$f^*(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$

- ▶ Worst-case conditional Bayes risk:

$$\max_{p(Y|X=x)} \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$

Relationship to Bayes optimal prediction

- ▶ Bayes optimal prediction minimizes the conditional risk:

$$f^*(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$

- ▶ Worst-case conditional Bayes risk:

$$\max_{p(Y|X=x)} \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$

- ▶ Cost-sensitive binary classification

$$\max_{p(Y|X=x)} \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[C(\hat{y}, Y) \mid X = x] = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = V(C)$$

Relationship to Bayes optimal prediction

- ▶ Bayes optimal prediction minimizes the conditional risk:

$$f^*(x) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$

- ▶ Worst-case conditional Bayes risk:

$$\max_{p(Y|X=x)} \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(\hat{y}, Y) \mid X = x]$$

- ▶ Cost-sensitive binary classification

$$\max_{p(Y|X=x)} \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[C(\hat{y}, Y) \mid X = x] = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_{\mathcal{Y}}} C(\mathbf{u}, \mathbf{v}) = V(C)$$

- ▶ The worst-case conditional Bayes risk for a binary prediction game defines the threshold between boostability and coin attainability

Simultaneous false positive and false negative guarantees

- The FP game $W^+ = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ and the FN game $W^- = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}$



Simultaneous false positive and false negative guarantees

- ▶ The FP game $W^+ = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ and the FN game $W^- = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}$
- ▶ The associated $(\mathcal{H} \times \mathcal{X})$ -matrices

$$M^+(h, x) = \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\} \quad M^-(h, x) = \mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}$$

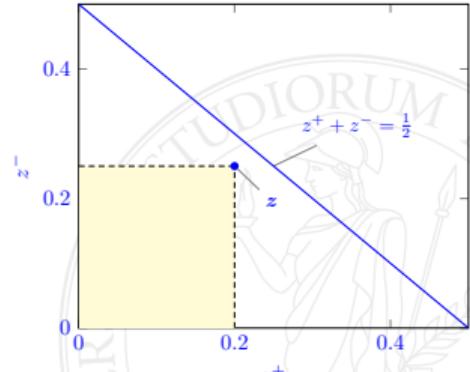
Simultaneous false positive and false negative guarantees

- The FP game $W^+ = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ and the FN game $W^- = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}$
- The associated $(\mathcal{H} \times \mathcal{X})$ -matrices

$$M^+(h, x) = \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\} \quad M^-(h, x) = \mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}$$

- WL oracle for $\mathbf{z} = (z^+, z^-)$
(simultaneous guarantees for FP and FN mistakes):

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$



Simultaneous false positive and false negative guarantees

- ▶ The FP game $W^+ = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ and the FN game $W^- = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}$
- ▶ The associated $(\mathcal{H} \times \mathcal{X})$ -matrices

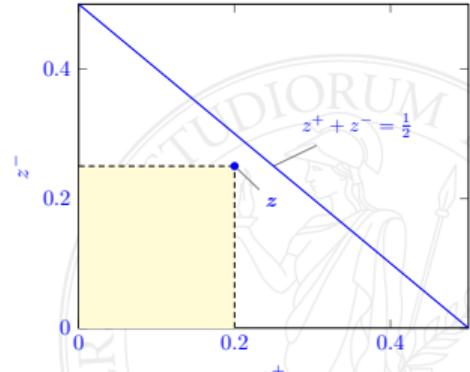
$$M^+(h, x) = \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\} \quad M^-(h, x) = \mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}$$

- ▶ WL oracle for $\mathbf{z} = (z^+, z^-)$

(simultaneous guarantees for FP and FN mistakes):

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ Which values of $\mathbf{z} = (z^+, z^-)$ are boostable?



Coin attainable guarantees

- Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

Coin attainable guarantees

- ▶ Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ This oracle returns $h : \mathcal{X} \rightarrow \mathcal{Y}$ attaining the desired guarantees against a given distribution over \mathcal{X}

Coin attainable guarantees

- ▶ Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ This oracle returns $h : \mathcal{X} \rightarrow \mathcal{Y}$ attaining the desired guarantees against a given distribution over \mathcal{X}
- ▶ **Coin-attainable region:** guarantees that are attainable with a **biased coin**

$$\mathcal{K} \equiv \left\{ (z^+, z^-) : \forall \mathbf{v} \in \Delta_{\mathcal{Y}} \quad \exists \mathbf{u} \in \Delta_{\mathcal{Y}} \quad W^+(\mathbf{u}, \mathbf{v}) \leq z^+ \wedge W^-(\mathbf{u}, \mathbf{v}) \leq z^- \right\}$$

Coin attainable guarantees

- ▶ Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ This oracle returns $h : \mathcal{X} \rightarrow \mathcal{Y}$ attaining the desired guarantees against a given distribution over \mathcal{X}
- ▶ **Coin-attainable region:** guarantees that are attainable with a **biased coin**

$$\mathcal{K} \equiv \left\{ (z^+, z^-) : \forall \mathbf{v} \in \Delta_{\mathcal{Y}} \quad \exists \mathbf{u} \in \Delta_{\mathcal{Y}} \quad W^+(\mathbf{u}, \mathbf{v}) \leq z^+ \wedge W^-(\mathbf{u}, \mathbf{v}) \leq z^- \right\}$$

- ▶ Can an oracle providing guarantees that are attainable with a coin be used for boosting?

Coin attainable guarantees

- ▶ Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ This oracle returns $h : \mathcal{X} \rightarrow \mathcal{Y}$ attaining the desired guarantees against a given distribution over \mathcal{X}
- ▶ **Coin-attainable region:** guarantees that are attainable with a **biased coin**

$$\mathcal{K} \equiv \left\{ (z^+, z^-) : \forall \mathbf{v} \in \Delta_{\mathcal{Y}} \quad \exists \mathbf{u} \in \Delta_{\mathcal{Y}} \quad W^+(\mathbf{u}, \mathbf{v}) \leq z^+ \wedge W^-(\mathbf{u}, \mathbf{v}) \leq z^- \right\}$$

- ▶ Can an oracle providing guarantees that are attainable with a coin be used for boosting?
- ▶ **Scalar case:** any $z \geq V(C)$ is coin attainable by definition and therefore not boostable by the cost-sensitive Boosting theorem

Coin attainable guarantees

- ▶ Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ This oracle returns $h : \mathcal{X} \rightarrow \mathcal{Y}$ attaining the desired guarantees against a given distribution over \mathcal{X}
- ▶ **Coin-attainable region:** guarantees that are attainable with a **biased coin**

$$\mathcal{K} \equiv \left\{ (z^+, z^-) : \forall \mathbf{v} \in \Delta_{\mathcal{Y}} \quad \exists \mathbf{u} \in \Delta_{\mathcal{Y}} \quad W^+(\mathbf{u}, \mathbf{v}) \leq z^+ \wedge W^-(\mathbf{u}, \mathbf{v}) \leq z^- \right\}$$

- ▶ Can an oracle providing guarantees that are attainable with a coin be used for boosting?
- ▶ **Scalar case:** any $z \geq V(C)$ is coin attainable by definition and therefore not boostable by the cost-sensitive Boosting theorem

Theorem (Multi-Objective Boosting Theorem)

Any $\mathbf{z} \notin \mathcal{K}$ is boostable

Coin attainable guarantees

- ▶ Recall WL oracle for $\mathbf{z} = (z^+, z^-)$:

$$\forall \mathbf{q} \in \Delta_{\mathcal{X}} \quad \exists h \in \mathcal{H} \quad M^+(h, \mathbf{q}) < z^+ \wedge M^-(h, \mathbf{q}) < z^-$$

- ▶ This oracle returns $h : \mathcal{X} \rightarrow \mathcal{Y}$ attaining the desired guarantees against a given distribution over \mathcal{X}
- ▶ **Coin-attainable region:** guarantees that are attainable with a **biased coin**

$$\mathcal{K} \equiv \left\{ (z^+, z^-) : \forall \mathbf{v} \in \Delta_{\mathcal{Y}} \quad \exists \mathbf{u} \in \Delta_{\mathcal{Y}} \quad W^+(\mathbf{u}, \mathbf{v}) \leq z^+ \wedge W^-(\mathbf{u}, \mathbf{v}) \leq z^- \right\}$$

- ▶ Can an oracle providing guarantees that are attainable with a coin be used for boosting?
- ▶ **Scalar case:** any $z \geq V(C)$ is coin attainable by definition and therefore not boostable by the cost-sensitive Boosting theorem

Theorem (Multi-Objective Boosting Theorem)

Any $\mathbf{z} \notin \mathcal{K}$ is boostable

What is the geometry of \mathcal{K} ?

Scalarization

For any $0 \leq \alpha \leq 1$

$$W_\alpha = \alpha W^+ + (1 - \alpha) W^- = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$$

$$M_\alpha = \underbrace{\alpha \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{FP}} + (1 - \alpha) \underbrace{(1 - \alpha) \mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{FN}}$$

Scalarization

For any $0 \leq \alpha \leq 1$

$$W_\alpha = \alpha W^+ + (1 - \alpha) W^- = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$$

$$M_\alpha = \underbrace{\alpha \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{FP}} + (1 - \alpha) \underbrace{(1 - \alpha) \mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{FN}}$$

- ▶ This is equivalent to **convex costs** $\alpha, 1 - \alpha$ for false positive and false negative mistakes

Scalarization

For any $0 \leq \alpha \leq 1$

$$W_\alpha = \alpha W^+ + (1 - \alpha) W^- = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$$

$$M_\alpha = \underbrace{\alpha \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{FP}} + (1 - \alpha) \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{FN}}$$

- ▶ This is equivalent to **convex costs** $\alpha, 1 - \alpha$ for false positive and false negative mistakes
- ▶ In the **scalar case** $W_\alpha = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$, any $z \geq V(W_\alpha)$ is coin attainable

Scalarization

For any $0 \leq \alpha \leq 1$

$$W_\alpha = \alpha W^+ + (1 - \alpha) W^- = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$$

$$M_\alpha = \underbrace{\alpha \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{FP}} + (1 - \alpha) \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{FN}}$$

- ▶ This is equivalent to **convex costs** $\alpha, 1 - \alpha$ for false positive and false negative mistakes
- ▶ In the **scalar case** $W_\alpha = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$, any $z \geq V(W_\alpha)$ is coin attainable

Theorem (Duality)

A multi-objective guarantee \mathbf{z} is coin-attainable iff it has no boostable scalarizations:

$$\mathcal{K} \equiv \left\{ \mathbf{z} = (z^+, z^-) : \forall \boldsymbol{\alpha} = (\alpha, 1 - \alpha) \quad \langle \boldsymbol{\alpha}, \mathbf{z} \rangle \geq V(W_\alpha) \right\}$$

Scalarization

For any $0 \leq \alpha \leq 1$

$$W_\alpha = \alpha W^+ + (1 - \alpha) W^- = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$$

$$M_\alpha = \underbrace{\alpha \mathbb{I}\{h(x) = 1 \wedge f(x) = -1\}}_{\text{FP}} + (1 - \alpha) \underbrace{\mathbb{I}\{h(x) = -1 \wedge f(x) = 1\}}_{\text{FN}}$$

- ▶ This is equivalent to **convex costs** $\alpha, 1 - \alpha$ for false positive and false negative mistakes
- ▶ In the **scalar case** $W_\alpha = \begin{pmatrix} 0 & \alpha \\ 1 - \alpha & 0 \end{pmatrix}$, any $z \geq V(W_\alpha)$ is coin attainable

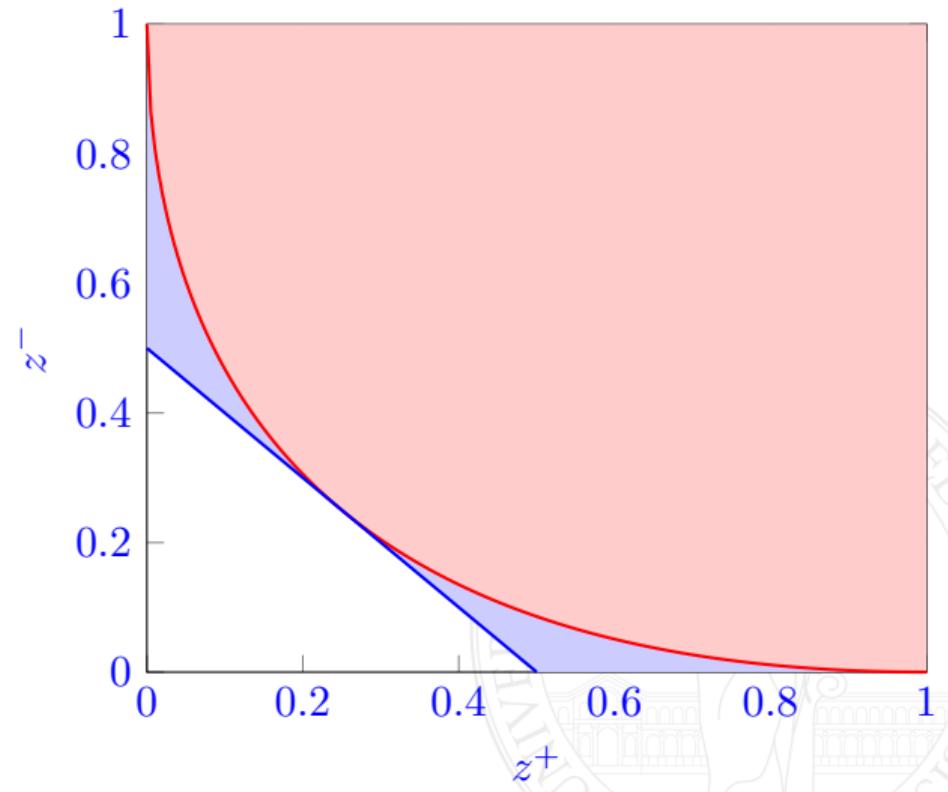
Theorem (Duality)

A multi-objective guarantee \mathbf{z} is coin-attainable iff it has no boostable scalarizations:

$$\begin{aligned} \mathcal{K} &\equiv \left\{ \mathbf{z} = (z^+, z^-) : \forall \boldsymbol{\alpha} = (\alpha, 1 - \alpha) \quad \langle \boldsymbol{\alpha}, \mathbf{z} \rangle \geq V(W_\alpha) \right\} \\ &\equiv \left\{ (z^+, z^-) \in [0, 1]^2 : \sqrt{z^+} + \sqrt{z^-} \geq 1 \right\} \end{aligned}$$

Geometry of \mathcal{K}

- ▶ Red: Coin attainable region
- ▶ Blue: Boostable only using oracle with simultaneous guarantees
- ▶ White: Boostable by both oracles



Extensions to the multiclass case

- ▶ Game matrix D for labels $\{1, \dots, n\}$ where $D(i, j)$ is cost of predicting i when true label is j



Extensions to the multiclass case

- ▶ Game matrix D for labels $\{1, \dots, n\}$ where $D(i, j)$ is cost of predicting i when true label is j
- ▶ For any subset $J \subseteq \mathcal{Y}$ of labels: $V_J(D) = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_J} D(\mathbf{u}, \mathbf{v})$



Extensions to the multiclass case

- ▶ Game matrix D for labels $\{1, \dots, n\}$ where $D(i, j)$ is cost of predicting i when true label is j
- ▶ For any subset $J \subseteq \mathcal{Y}$ of labels: $V_J(D) = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_J} D(\mathbf{u}, \mathbf{v})$
- ▶ This is the smallest loss one can ensure by predicting with a die knowing that the correct label is in J .



Extensions to the multiclass case

- ▶ Game matrix D for labels $\{1, \dots, n\}$ where $D(i, j)$ is cost of predicting i when true label is j
- ▶ For any subset $J \subseteq \mathcal{Y}$ of labels: $V_J(D) = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_J} D(\mathbf{u}, \mathbf{v})$
- ▶ This is the smallest loss one can ensure by predicting with a die knowing that the correct label is in J .
- ▶ Boosting is only possible within any two consecutive game values

$$0 = v_1 < v_2 < \dots < v_j < v_{j+1} < \dots < v_N = V_{\mathcal{Y}}(D)$$

Extensions to the multiclass case

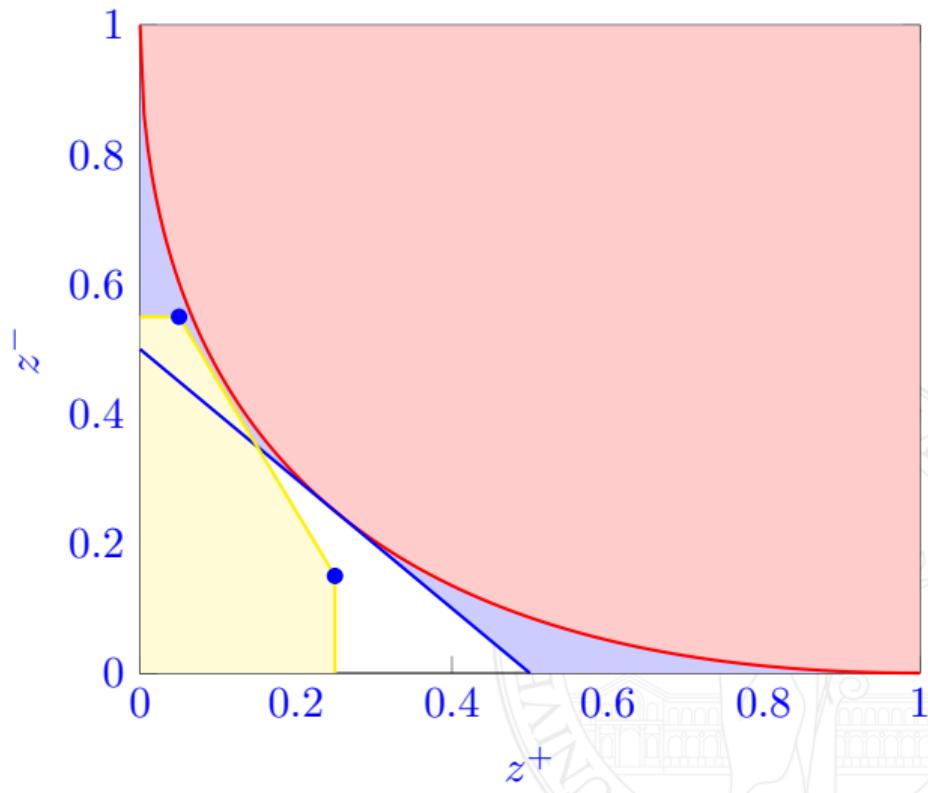
- ▶ Game matrix D for labels $\{1, \dots, n\}$ where $D(i, j)$ is cost of predicting i when true label is j
- ▶ For any subset $J \subseteq \mathcal{Y}$ of labels: $V_J(D) = \min_{\mathbf{u} \in \Delta_{\mathcal{Y}}} \max_{\mathbf{v} \in \Delta_J} D(\mathbf{u}, \mathbf{v})$
- ▶ This is the smallest loss one can ensure by predicting with a die knowing that the correct label is in J .
- ▶ Boosting is only possible within any two consecutive game values

$$0 = v_1 < v_2 < \dots < v_j < v_{j+1} < \dots < v_N = V_{\mathcal{Y}}(D)$$

- ▶ Extension to multiobjective multiclass

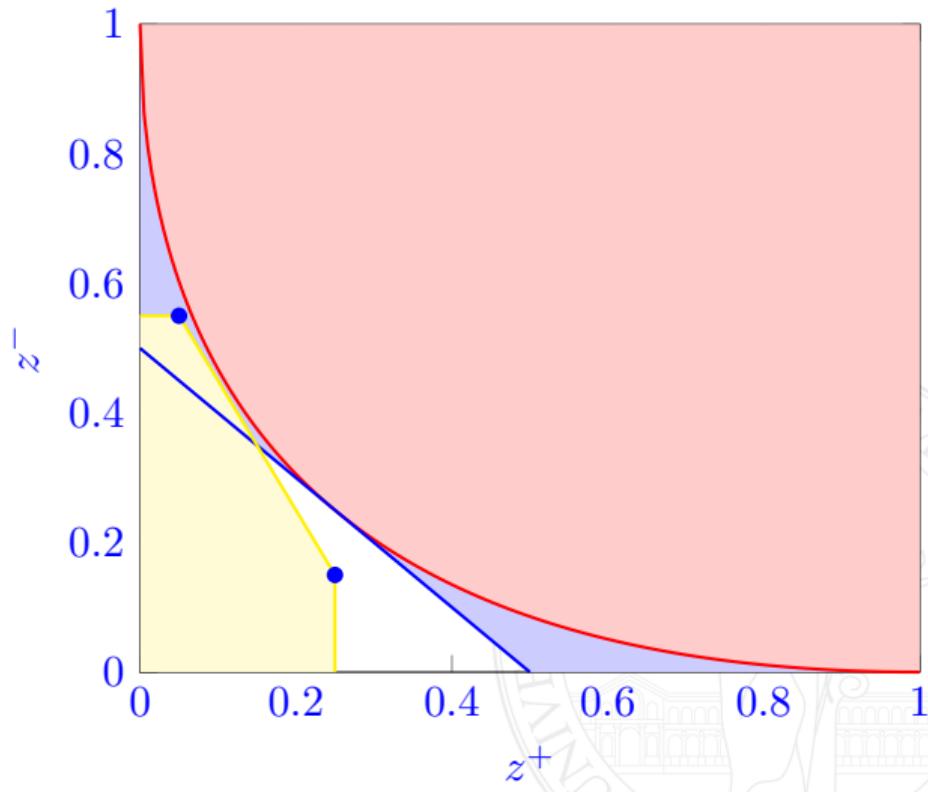
Open problems

- ▶ Suppose that the WL oracle provides guarantees in some convex region



Open problems

- ▶ Suppose that the WL oracle provides guarantees in some convex region
- ▶ Can we boost as long as the oracle region does not intersect the coin attainable region?



Open problems

- ▶ Suppose that the WL oracle provides guarantees in some convex region
- ▶ Can we boost as long as the oracle region does not intersect the coin attainable region?
- ▶ Can we do it adaptively?

