# Random Classification Noise
# does not Defeat All Convex Potential Boosters
# Irrespective of Model Choice

**Yishay Mansour**

Tel Aviv U.
& Google Research
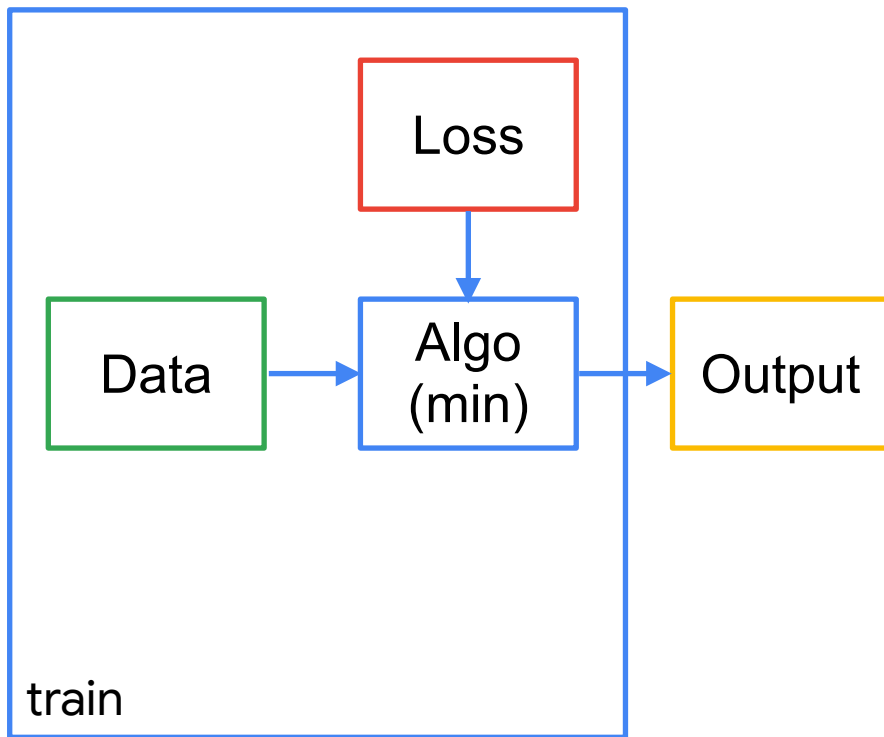
**Richard Nock**

Google Research

**Robert C. Williamson**

Tübingen U.
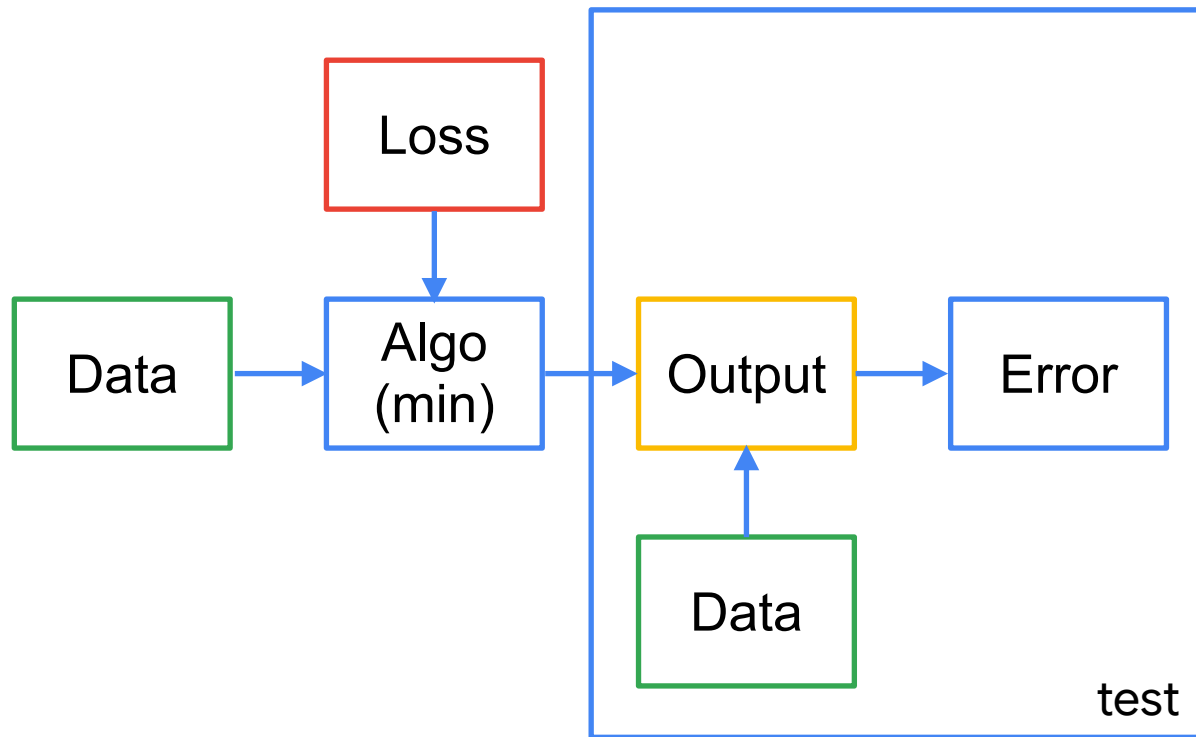& Tübingen AI Center

# Why this work ?

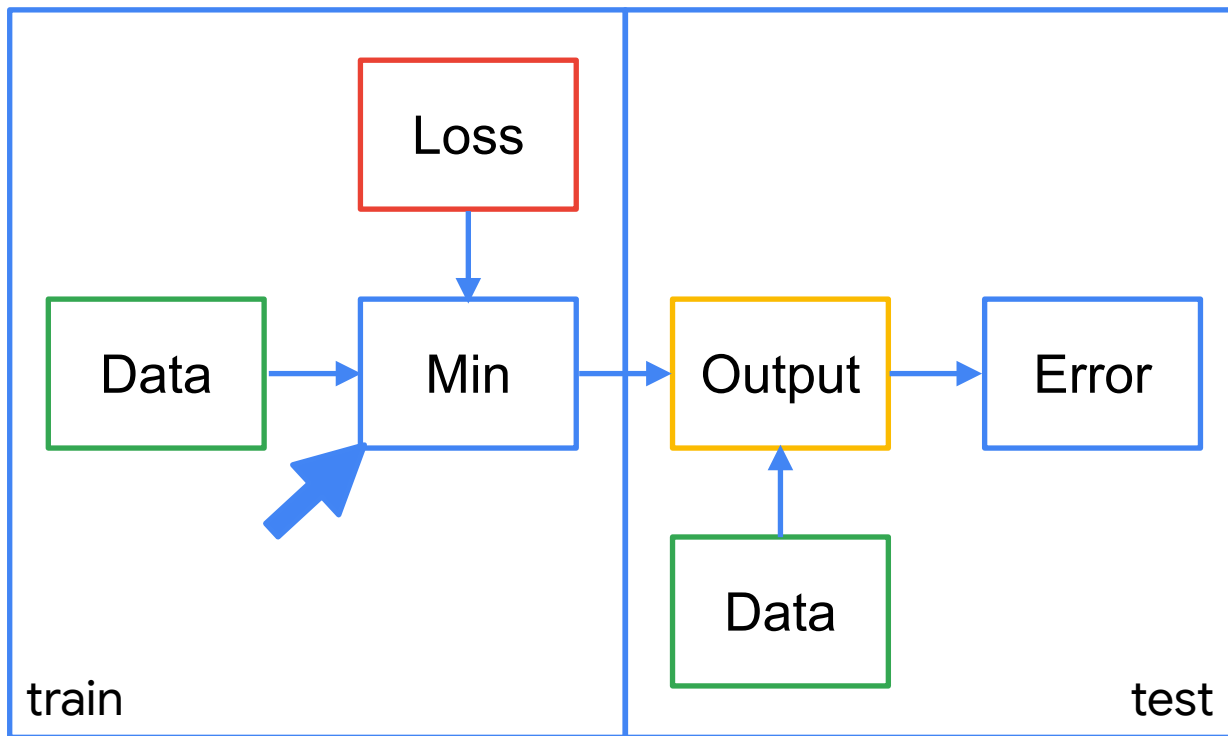Mansour, Nock & Williamson — ICML'23

# Long & Servedio (L&S) - Setting II

Long, P.-M. and Servedio, R.-A. Random classification noise defeats all convex potential boosters. In $25^{th}$ ICML, pp. 608–615, 2008b.

Google Research
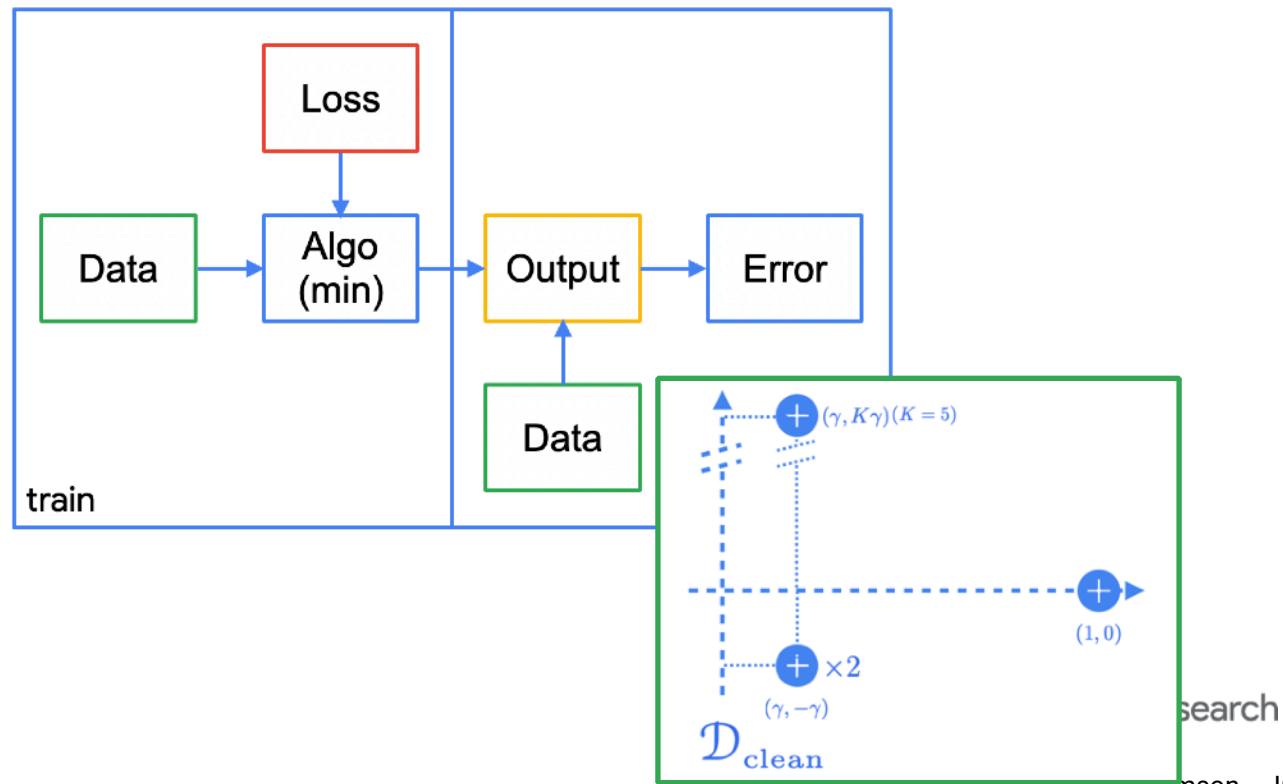
# Setting II

# Setting I

# Setting II - test data

# Setting II - training data



$$\eta_{\scriptscriptstyle Y} = \frac{1}{N+1}$$

$\mathcal{D}_{\mathbf{noisy}}$

$(N > 1)$

Loss

Data → Algo (min) → Output → Error

Data

train

$\mathcal{D}_{clean}$

$(\gamma, K\gamma)(K = 5)$

$(1,0)$

$(\gamma, -\gamma)$

# Setting II - data summary



$$\eta_Y = \frac{1}{N+1}$$

$\mathcal{D}_{noisy}$

Loss

Data

Algo (min)

Output

Error

Data

train

- 2 classes (-1,+1)
- whole domain **known**

$(\gamma, K\gamma)(K=5)$

$(1,0)$

$(\gamma, -\gamma)$

$\mathcal{D}_{clean}$

search

# Setting II - loss



$$\phi \in \mathcal{C}_{\mathrm{sur}}$$

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

$$\eta_v = \frac{1}{N+1}$$

$\mathcal{D}_{\mathrm{noisy}}$

Loss

Data → Algo (min) → Output → Error

Data

train

$(\gamma, K\gamma)(K=5)$

$(1,0)$

$(\gamma, -\gamma)$

$\mathcal{D}_{\mathrm{clean}}$

search

# Setting II - outputs

$$\phi \in \mathcal{C}_{\mathrm{sur}}$$

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

$$\eta_{\scriptscriptstyle N} = \frac{1}{N+1}$$

$$(\quad \times 2) \times N +$$

$\mathcal{D}_{\mathrm{noisy}}$

Loss

Data

Algo (min)

Real-valued (sign = class)

Output

$h$

Error

Data

train

$(\gamma, K\gamma)(K = 5)$

$(i, 0)$

$\times 2$

$(\gamma, -\gamma)$

$\mathcal{D}_{\mathrm{clean}}$

search

# Setting II - algorithm



Loss

Data → Algo (min) → Output → Error

Boosting "*à-la* Adaboost"

Data

$h_{\mathrm{boost}}$

$\eta_{\scriptscriptstyle Y} = \dfrac{1}{N+1}$

$\mathcal{D}_{\mathrm{noisy}}$

$\mathcal{D}_{\mathrm{clean}}$

# Setting II - Key result



$$\mathrm{Err}(h_{\mathrm{boost}}, \mathcal{D}_{\mathrm{clean}}) = 0.5$$

already after few iterations

Loss

Data → Algo (min) → Output → Error

$h_{\mathrm{boost}}$

Boosting "à-la Adaboost"

Data

$\eta_k = \dfrac{1}{N+1}$

$\mathcal{D}_{\mathrm{noisy}}$

$\mathcal{D}_{\mathrm{clean}}$

# Setting II - Key result



Loss

Data → Algo (min) → Output → Error

$h_{\text{boost}}$

Boosting "à-la Adaboost"

Data

$\eta_k = \frac{1}{N+1}$

$\mathcal{D}_{\text{noisy}}$

$\mathcal{D}_{\text{clean}}$

$\text{Err}(h_{\text{boost}}, \mathcal{D}_{\text{clean}}) = 0.5$

already after few iterations

# Setting I - Key result



$\phi \in \mathcal{C}_{\mathrm{sur}}$

"margin loss"

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

$\mathcal{D}_{\mathrm{noisy}}$

$\eta_x = \dfrac{1}{N+1}$

Loss

Data → Min → Output → Error

$h^*$

Data

train

$\mathcal{D}_{\mathrm{clean}}$

# Setting I - Key result



$\phi \in \mathcal{C}_{\mathrm{sur}}$

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

$\mathrm{Err}(h^*, \mathcal{D}_{\mathrm{clean}}) = 0.5$

$\mathcal{D}_{\mathrm{noisy}}$

$\eta_k = \dfrac{1}{N+1}$

Loss

Data → Min → Output → Error

$h^*$

Data

train

$\mathcal{D}_{\mathrm{clean}}$

# Setting I - Key result



$\phi \in \mathcal{C}_{\mathrm{sur}}$

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

$\mathrm{Err}(h^*, \mathcal{D}_{\mathrm{clean}}) = 0.5$

$\eta_v = \dfrac{1}{N+1}$

$\mathcal{D}_{\mathrm{noisy}}$

Loss

Data → Min → Output → Error

$h^*$

Data

train

$\mathcal{D}_{\mathrm{clean}}$

search

the "simplest" form of corruption
defeats two praised ML components:
convex [losses | boosters]...

( x100s)

Google Research

Mansour, Nock & Williamson — ICML'23

the "simplest" form of corruption
defeats two praised ML components:
convex [losses | boosters]...

or does it ?

why this work

Google Research

Mansour, Nock & Williamson — ICML'23

( x100s)

# Enters Savage

Savage, L.-J. Elicitation of personal probabilities and ex-pectations. *J. of the Am. Stat. Assoc.*, pp. 783–801, 1971.

Google Research

# Setting I tweak (temporary)

# Class Probability Estimation

Class prediction $\rightarrow$ posterior prediction ($\hat{p}[y = 1|\boldsymbol{x}]$)

# Class Probability Estimation



Class prediction $\longrightarrow$ posterior prediction ($\hat{p}[y = 1|\boldsymbol{x}]$)

CPE loss (pointwise)

partial losses

$$\ell(y, u) \doteq [\![y = 1]\!] \cdot \boxed{\ell_1(u)} + [\![y = -1]\!] \cdot \boxed{\ell_{-1}(u)}$$

estimated posterior in [0,1]

true label / class in {-1,1}

# Class Probability Estimation

Class prediction $\rightarrow$ posterior prediction ($\hat{p}[y = 1 | \boldsymbol{x}]$)
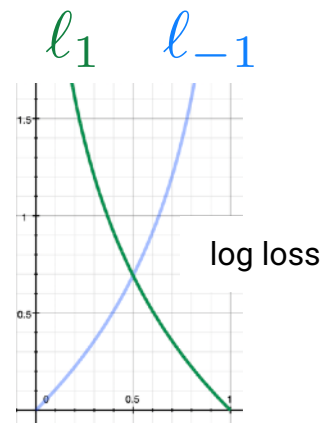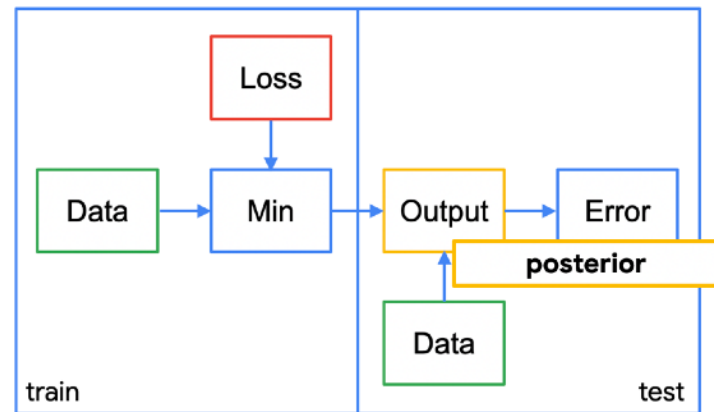
CPE loss (pointwise)

partial losses

$$\ell(y, u) \doteq [\![ y = 1 ]\!] \cdot \boxed{\ell_1(u)} + [\![ y = -1 ]\!] \cdot \boxed{\ell_{-1}(u)}$$

estimated posterior in [0,1]

true label / class in {-1,1}



$\ell_1$   $\ell_{-1}$

log loss

symmetry axis

# Class Probability Estimation

Class prediction → posterior prediction ($\hat{p}[y = 1|\boldsymbol{x}]$)

CPE loss (pointwise)

partial losses

$$\ell(y, u) \doteq [\![y = 1]\!] \cdot \boxed{\ell_1(u)} + [\![y = -1]\!] \cdot \boxed{\ell_{-1}(u)}$$

estimated posterior in [0,1]

true label / class in {-1,1}

CPE loss (population)

$$\Phi(\eta, \mathcal{D}) \doteq \mathbb{E}_{(\boldsymbol{x},y) \sim \mathcal{D}} \left[ \ell(y, \eta(\boldsymbol{x})) \right]$$





$\ell_1$  $\ell_{-1}$

log loss

Google Research

Mansour, Nock & Williamson — ICML'23

# Properness



Class prediction → posterior prediction ( $\hat{p}[y = 1|\boldsymbol{x}]$ )

CPE loss (pointwise)

partial losses

$$\ell(y, u) \doteq [\![ y = 1 ]\!] \cdot \ell_1(u) + [\![ y = -1 ]\!] \cdot \ell_{-1}(u)$$

estimated posterior in [0,1]
true label / class in {-1,1}

CPE loss (population)

$$\Phi(\eta, \mathcal{D}) \doteq \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[ \ell(y, \eta(\boldsymbol{x})) \right]$$

Quality: **strict properness** (strict optimum = Bayes prediction)

$$\eta_{\mathrm{Bayes}} = \arg\min_{\eta} \Phi(\eta, \mathcal{D})$$



$\ell_1 \quad \ell_{-1}$

log loss

Mansour, Nock & Williamson — ICML'23

# Properness



Class prediction → posterior prediction ($\hat{p}[y = 1 | \boldsymbol{x}]$)

CPE loss (pointwise)

partial losses

$$\ell(y, u) \doteq [\![y = 1]\!] \cdot \boxed{\ell_1(u)} + [\![y = -1]\!] \cdot \boxed{\ell_{-1}(u)}$$

estimated posterior in [0,1]

true label / class in {-1,1}

CPE loss (population)

$$\Phi(\eta, \mathcal{D}) \doteq \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[ \ell(y, \eta(\boldsymbol{x})) \right]$$

Quality:     **properness** (          optima $\supseteq$ Bayes prediction)



$$\eta_{\text{Bayes}} \in \arg\min_{\eta} \Phi(\eta, \mathcal{D})$$
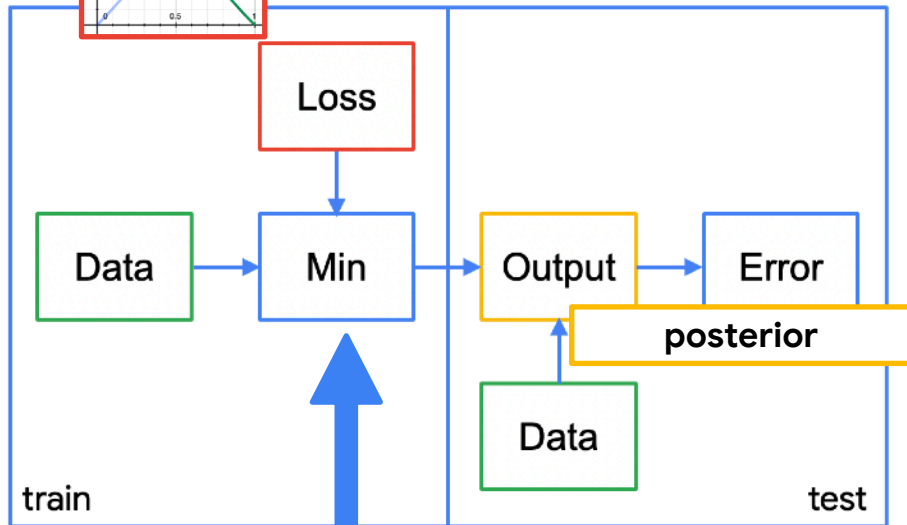


0/1 loss

symmetry axis

# Back to L&S (Setting I)

# Savage on L&S (Setting I)



- strictly proper,
- symmetric,
- differentiable

Loss

Data → Min → Output → Error
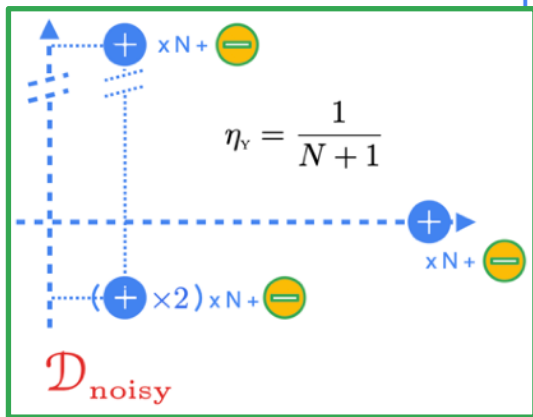
posterior

Data

train

test

⚠️ Because we train on the full domain, min sought =

# Savage on L&S (Setting I)



- strictly proper,
- symmetric,
- differentiable

$$\eta_Y = \frac{1}{N+1}$$

$$\mathcal{D}_{\mathbf{noisy}}$$

$$\eta_{\mathrm{Bayes}} = \frac{N}{N+1}$$
$$(\forall \boldsymbol{x}\,!)$$

Loss

Data → Min → Output → Error

posterior

Data

train                    test

Mansour, Nock & Williamson — ICML'23

# Savage on L&S (Setting I)



- strictly proper,
- symmetric,
- differentiable

$$\eta_{\gamma} = \frac{1}{N+1}$$

$$\mathcal{D}_{\mathbf{noisy}}$$

$$\eta_{\mathrm{Bayes}} = \frac{N}{N+1}$$
$$(\forall \boldsymbol{x}!)$$

Loss

Data → Min → Output → Error

**posterior**

Data

train

only positive examples

$\& \ \eta_{\mathrm{Bayes}} > \frac{1}{2}$

$(\gamma, K\gamma)(K=5)$

$(\gamma, -\gamma)$ ×2

$(1,0)$

$$\mathcal{D}_{\mathbf{clean}}$$
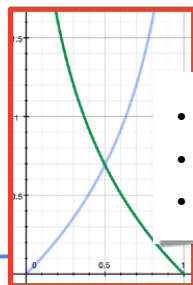
# Savage on L&S (Setting I)



- strictly proper,
- symmetric,
- differentiable

$$\eta_Y = \frac{1}{N+1}$$

$$\mathcal{D}_{\mathbf{noisy}}$$

$$\eta_{\mathrm{Bayes}} = \frac{N}{N+1}$$

$$(\forall \boldsymbol{x}!)$$

Loss

Data → Min → Output → Error

**posterior**

Data

train

$$\mathrm{Err}(\eta_{\mathrm{Bayes}}, \mathcal{D}_{\mathrm{clean}}) = 0$$

$$\eta_{\mathrm{Bayes}} > \frac{1}{2}$$

$(\gamma, K\gamma)(K=5)$

$(\gamma, -\gamma)$ ×2

$(1, 0)$

$$\mathcal{D}_{\mathbf{clean}}$$

search

# Savage on Setting I

# L&S on Setting I

Mansour, Nock & Williamson — ICML'23

# Savage on Setting I



- strictly proper,
- symmetric,
- differentiable

posterior $\eta$

$$\mathrm{Err}(\eta_{\mathrm{Bayes}}, \mathcal{D}_{\mathrm{clean}}) = 0$$

Loss

Data → Min → Output → Error

Data

train

test

# Savage on Setting I

# L&S on Setting I



- strictly proper,
- symmetric,
- differentiable

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

Loss

Min

Data

Output

Error

posterior $\eta$

real-valued $h$

Data

train

test

$$\mathrm{Err}(\eta_{\mathrm{Bayes}}, \mathcal{D}_{\mathrm{clean}}) = 0$$

$$\mathrm{Err}(h^*, \mathcal{D}_{\mathrm{clean}}) = 0.5$$

Google Research

Mansour, Nock & Williamson — ICML'23

# But...



- strictly proper,
- symmetric,
- differentiable

posterior $\eta$

?

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

real-valued $h$

posterior $\eta$

- strictly proper,
- symmetric,
- differentiable

$\subset$

real-valued $h$

- convex, nonincreasing
- differentiable, $\phi'(0) < 0$
- lowerbounded

↳ Minimization of any* strictly proper, symmetric, differentiable CPE loss can be formulated as a convex surrogate minimization for a real valued classifier with a correspondence via the (canonical) link of the loss:

$$\eta \doteq (\ell_{-1} - \ell_1)^{-1}(h)$$

strictly proper,
symmetric,
differentiable

$\subset$

convex, nonincreasing
differentiable, $\phi'(0) < 0$
lowerbounded

posterior $\eta$

real-valued $h$

↳ Minimization of any* strictly proper, symmetric, differentiable CPE loss can be formulated as a convex surrogate minimization for a real valued classifier with a correspondence via the (canonical) link of the loss:

$$\eta \doteq (\ell_{-1} - \ell_1)^{-1}(h)$$

paradox ?

Google Research

# What about properness without symmetry ?

Strict properness **without** symmetry assumption:

↳ asymmetry brings much more freedom to
fine-tune costs



$\ell_{-1}$

$\ell_1$

- strictly proper,
- sym~~metric~~,
- differentiable

**Google** Research

# What about properness without symmetry ?

Strict properness **without** symmetry assumption:

↳ asymmetry brings much more freedom to
  fine-tune costs
↳ no "classical" margin formulation anymore
  -- "escapes" Long & Servedio's setting

$\ell_{-1}$

$\ell_1$

$\not\subset$

- strictly proper,
- sym~~metric~~,
- differentiable

L&S

# What about properness without symmetry ?

Strict properness **without** symmetry assumption:

↳ asymmetry b
   fine-tune co

↳ no "classical"
   -- "escapes"

$\ell_{-1}$

L&S

## Setting II

$\text{Err}(h_{\text{boost}}, \mathcal{D}_{\text{clean}}) = 0.5$

already after #iterations
as small as **2**

## Setting I

$\text{Err}(h^*, \mathcal{D}_{\text{clean}}) = 0.5$

properness as a whole "useless" !

Google Research

# Let us cut to the chase...

In-context, hardness has nothing to do with
↪ the convexity of the loss
↪ nor the fact that algorithm = boosting



Mansour, Nock & Williamson — ICML'23

# Let us cut to the chase...

In-context, hardness has nothing to do with
↳ the convexity of the loss
↳ nor the fact that algorithm = boosting

Culprit = model class

🔶 **Linear Separators**
"break" the guarantee of properness
under the "simplest" noise model

... and ...



$$(\gamma, K\gamma) \quad (K = 5)$$

$$\boldsymbol{\theta}_{\text{noisy}}$$

$$\boldsymbol{\theta}_{\text{clean}}$$

$$(1, 0)$$

$$(\gamma, -\gamma) \times 2$$

Loss | Data | Min | Algo (min) | Output | Error | Data | $\boldsymbol{\theta}_{\text{noisy}}$ | $\boldsymbol{\theta}_{\text{clean}}$ | train | test

Google Research

# Let us cut to the chase...

In-context, hardness has nothing to do with
↳ the convexity of the loss
↳ nor the fact that algorithm = boosting

Culprit = model class

**Linear Separators**
"break" the guarantee of properness
under the "simplest" noise model

... and we are also going to show it **constructively**

requires a new convex booster...

$(\gamma, K\gamma)$  $(K = 5)$

$\boldsymbol{\theta}_{\text{noisy}}$

$\boldsymbol{\theta}_{\text{clean}}$

$(1, 0)$

$\times 2$
$(\gamma, -\gamma)$

Loss

$\boldsymbol{\theta}_{\text{noisy}}$

$\boldsymbol{\theta}_{\text{clean}}$

Data      Min      Output      Error

Algo
(min)

Data

train                                    test

# Convex boosting, model-adaptive

ModaBoost (Model-Adaptive Boosting)

↳ **Start**: Adaboost-style boosting for
***strictly proper, symmetric, differentiable losses***


proper losses with "margin form" (function-of(y * H))

# Convex boosting, model-adaptive

ModaBoost (Model-Adaptive Boosting)

↪ **Start**: Adaboost-style boosting for
***strictly proper, symmetric, differentiable losses***

- Weights $\boldsymbol{w}$ = record of past performances

- ...

- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random

$$|\mathbb{E}_{\boldsymbol{w}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$

- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↪ Returns a linear model $H \doteq \sum_t \alpha_t \cdot h_t$



proper losses with
"margin form"
(function-of(y * H))

Weak Learning
Assumption (WLA)

Google Research

# Convex boosting, model-adaptive

ModaBoost (Model-Adaptive Boosting)

↪ **Step 1**: lift the applicable losses to all
***strictly proper, sym❌etric, differentiable loss***

- Weights $\boldsymbol{w}$ = record of past performances

- ...

- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random
  $$|\mathbb{E}_{\boldsymbol{w}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$

- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↪ Returns a linear model $H \doteq \sum_t \alpha_t \cdot h_t$

no more "margin form"

two convex surrogates instead of 1

Google Research

# Convex boosting, model-adaptive

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

  • Weights $\boldsymbol{w}$ = record of past performances

  • <u>Architecture Emulation Oracle</u> : outputs $\underline{\mathcal{S} \subseteq \mathcal{X}}$

# Convex boosting, model-adaptive

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$

$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$

# Convex boosting, model-adaptive

ModaBoost (Model-Adaptive Boosting)

↪ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances

- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$

- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random on $\mathcal{S}$
  $$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$

- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↪ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

**Theorem 1.** *Suppose the following assumptions are satisfied on the loss and weak learner:*

**LOSS** *the loss is strictly proper differentiable; its partial losses are such that* $\exists \kappa > 0, C \in \mathbb{R}$,

$$\ell_{-1}(0), \ell_1(1) \geqslant C, \tag{18}$$
$$\inf\{\ell'_{-1} - \ell'_1\} \geqslant \kappa. \tag{19}$$

**WLA** *There exists a constant* $\gamma_{\mathrm{WL}} > 0$ *such that at each iteration* $t \in [T]$, *the weak hypothesis* $h_t$ *returned by* WL *satisfies*

$$\left| \sum_{i \in [m]_t} \frac{w_{t,i}}{\sum_{j \in [m]_t} w_{t,j}} \cdot y_i \cdot \frac{h_t(\boldsymbol{x}_i)}{\max_{j \in [m]_t} |h_t(\boldsymbol{x}_j)|} \right| \geqslant \gamma_{\mathrm{WL}}. \tag{20}$$

**AEOC** *there exists a sequence* $\{u_t\}_{t \in \mathbb{N}_{>0}}$ *of strictly positive reals such that the choice of* $\mathcal{X}_t$ *in Step 2.1 is* $u_t$ *compliant.*

*Then for any* $\theta \geqslant 0, \varepsilon > 0$, *letting* $\underline{w}(\theta) \doteq \min\{1 - (-\underline{L}')^{-1}(\theta), (-\underline{L}')^{-1}(-\theta)\}$, *if* MODABOOST *is run for at least*

$$T \geqslant U^{-1}\left( \frac{2\left(\Phi(H_0, \mathcal{S}) - C\right)}{\kappa \cdot \varepsilon^2 \underline{w}(\theta)^2 \gamma_{\mathrm{WL}}^2} \right) \tag{21}$$

*iterations, then we are guaranteed*

$$\mathbb{P}_{i \sim [m]}[y_i \, H_T(\boldsymbol{x}_i) \leqslant \theta] < \varepsilon. \tag{22}$$

*Here,* $U$ *is crafted as in* (17).

Mansour, Nock & Williamson − ICML'23

**Theorem 1.** *Suppose the following assumptions are satisfied on the loss and weak learner:*

**LOSS** *the loss is strictly proper differentiable; its partial losses are such that $\exists \kappa > 0, C \in \mathbb{R}$,*

$$\ell_{-1}(0), \ell_1(1) \geq C, \tag{18}$$
$$\inf\{\ell'_{-1} - \ell'_1\} \geq \kappa. \tag{19}$$

Assumptions on loss, "necessary"

**WLA** *There exists a constant $\gamma_{\mathrm{WL}} > 0$ such that at each iteration $t \in [T]$, the weak hypothesis $h_t$ returned by $\mathrm{WL}$ satisfies*

$$\left| \sum_{i \in [m]_t} \frac{w_{t,i}}{\sum_{j \in [m]_t} w_{t,j}} \cdot y_i \cdot \frac{h_t(\boldsymbol{x}_i)}{\max_{j \in [m]_t} |h_t(\boldsymbol{x}_j)|} \right| \geq \gamma_{\mathrm{WL}}. \tag{20}$$

**AEOC** *there exists a sequence $\{u_t\}_{t \in \mathbb{N}_{>0}}$ of strictly positive reals such that the choice of $\mathcal{X}_t$ in Step 2.1 is $u_t$ compliant.*

*Then for any $\theta \geq 0, \varepsilon > 0$, letting $\underline{w}(\theta) \doteq \min\{1 - (-\underline{L}')^{-1}(\theta), (-\underline{L}')^{-1}(-\theta)\}$, if $\mathrm{MODABOOST}$ is run for at least*

$$T \geq U^{-1}\left( \frac{2\left(\Phi(H_0, \mathcal{S}) - C\right)}{\kappa \cdot \varepsilon^2 \underline{w}(\theta)^2 \gamma_{\mathrm{WL}}^2} \right) \tag{21}$$

*iterations, then we are guaranteed*

$$\mathbb{P}_{i \sim [m]}[y_i \, H_T(\boldsymbol{x}_i) \leq \theta] < \varepsilon. \tag{22}$$

*Here, $U$ is crafted as in* (17).

**Theorem 1.** *Suppose the following assumptions are satisfied on the loss and weak learner:*

**LOSS** *the loss is strictly proper differentiable; its partial losses are such that $\exists \kappa > 0, C \in \mathbb{R}$,*

$$\ell_{-1}(0), \ell_1(1) \geqslant C, \tag{18}$$
$$\inf\{\ell'_{-1} - \ell'_1\} \geqslant \kappa. \tag{19}$$

**WLA** *There exists a constant $\gamma_{\mathrm{WL}} > 0$ such that at each iteration $t \in [T]$, the weak hypothesis $h_t$ returned by $\mathrm{WL}$ satisfies*

$$\left| \sum_{i \in [m]_t} \frac{w_{t,i}}{\sum_{j \in [m]_t} w_{t,j}} \cdot y_i \cdot \frac{h_t(\boldsymbol{x}_i)}{\max_{j \in [m]_t} |h_t(\boldsymbol{x}_j)|} \right| \geqslant \gamma_{\mathrm{WL}}. \tag{20}$$

**AEOC** *there exists a sequence $\{u_t\}_{t \in \mathbb{N}_{>0}}$ of strictly positive reals such that the choice of $\mathfrak{X}_t$ in Step 2.1 is $u_t$ compliant. Then for any $\theta \geqslant 0, \varepsilon > 0$, letting $\underline{w}(\theta) \doteq \min\{1 - (-\underline{L}')^{-1}(\theta), (-\underline{L}')^{-1}(-\theta)\}$, if $\mathrm{MODABOOST}$ is run for at least*

$$T \geqslant U^{-1}\left( \frac{2\left(\Phi(H_0, \mathcal{S}) - C\right)}{\kappa \cdot \varepsilon^2 \underline{w}(\theta)^2 \gamma_{\mathrm{WL}}^2} \right) \tag{21}$$

*iterations, then we are guaranteed*

$$\mathbb{P}_{i \sim [m]}[y_i^* H_T(\boldsymbol{x}_i) \leqslant \theta] < \varepsilon. \tag{22}$$

*Here, $U$ is crafted as in* (17).

**Theorem 1.** *Suppose the following assumptions are satisfied on the loss and weak learner:*

**LOSS** *the loss is strictly proper differentiable; its partial losses are such that $\exists \kappa > 0, C \in \mathbb{R},$*

$$\ell_{-1}(0), \ell_1(1) \geq C, \tag{18}$$
$$\inf\{\ell'_{-1} - \ell'_1\} \geq \kappa. \tag{19}$$

**WLA** *There exists a constant $\gamma_{\text{WL}} > 0$ such that at each iteration $t \in [T]$, the weak hypothesis $h_t$ returned by $\text{WL}$ satisfies*

$$\left| \sum_{i \in [m]_t} \frac{w_{t,i}}{\sum_{j \in [m]_t} w_{t,j}} \cdot y_i \cdot \frac{h_t(\boldsymbol{x}_i)}{\max_{j \in [m]_t} |h_t(\boldsymbol{x}_j)|} \right| \geq \gamma_{\text{WL}}. \tag{20}$$

**AEOC** *there exists a sequence $\{u_t\}_{t \in \mathbb{N}_{>0}}$ of strictly positive reals such that the choice of $\mathcal{X}_t$ in Step 2.1 is $u_t$ compliant.*

*Then for any $\theta \geq 0, \varepsilon > 0$, letting $\underline{w}(\theta) \doteq \min\{1 - (-\underline{L}')^{-1}(\theta), (-\underline{L}')^{-1}(-\theta)\}$, if $\text{MODABOOST}$ is run for at least*

$$T \geq U^{-1}\left( \frac{2\left( \Phi(H_0, \mathcal{S}) - C \right)}{\kappa \cdot \varepsilon^2 \underline{w}(\theta)^2 \gamma_{\text{WL}}^2} \right) \tag{21}$$

*iterations, then we are guaranteed*

$$\mathbb{P}_{i \sim [m]}[y_i \, H_T(\boldsymbol{x}_i) \leq \theta] < \varepsilon. \quad \cdots \quad \boxed{\text{guarantee on edges / margins}} \tag{22}$$

*Here, $U$ is crafted as in* (17).

# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final
    emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random on $\mathcal{S}$
  $$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t \mathbb{1}_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

AEO → Models != Linear Models

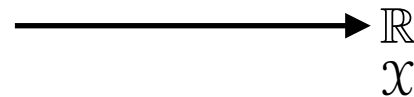# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t \mathbf{1}_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

|  | AEO | Weak Learner | Model |

$$\mathbb{R}$$
$$\mathcal{X}$$

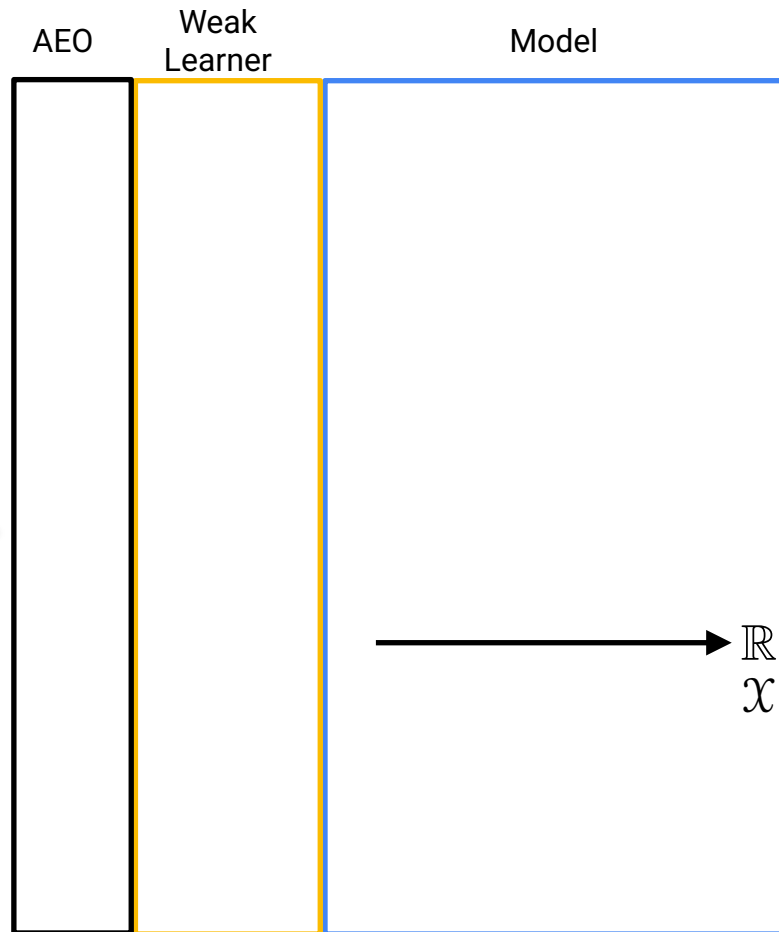# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
  $$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t \mathbb{1}_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ | | |

$\mathbb{R}$
$\mathcal{X}$

# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is ⟺ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\underline{\mathcal{S} \subseteq \mathcal{X}}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\underline{\mathcal{S}}$
$$|\mathbb{E}_{\underline{\boldsymbol{w}_{|\mathcal{S}}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t \underline{1_{\boldsymbol{x} \in \mathcal{S}_t}} \alpha_t \cdot h_t(\boldsymbol{x})$

| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ | constant ┈┈┈➤ | $z_0$ |

$$\xrightarrow{\hspace{3cm}} \begin{array}{c} \mathbb{R} \\ \mathcal{X} \end{array}$$

$z_0$

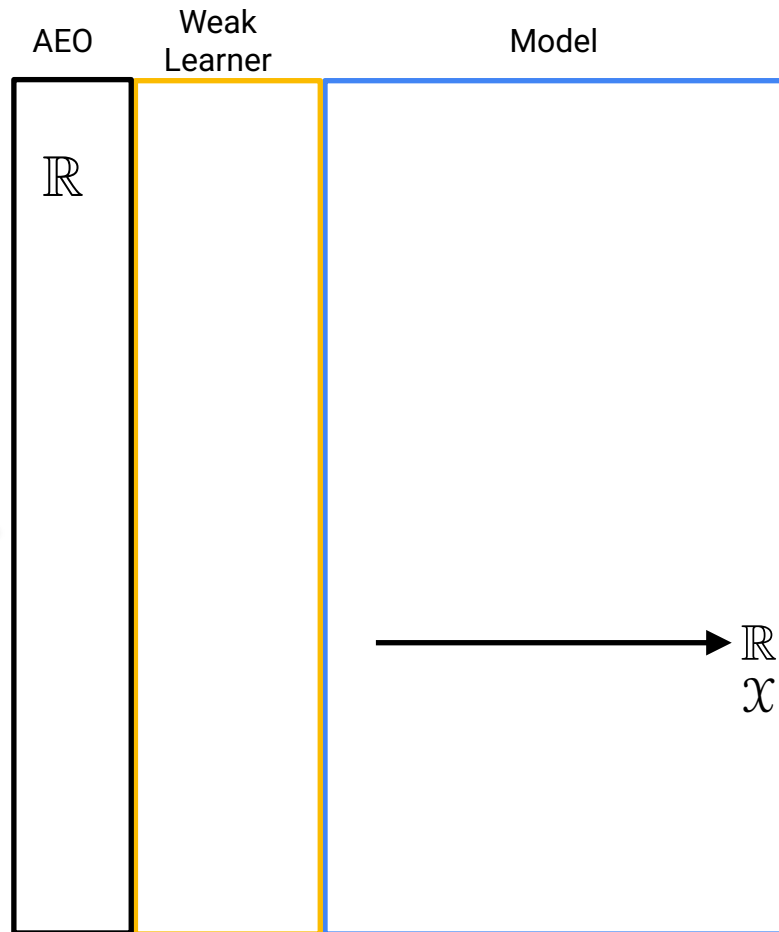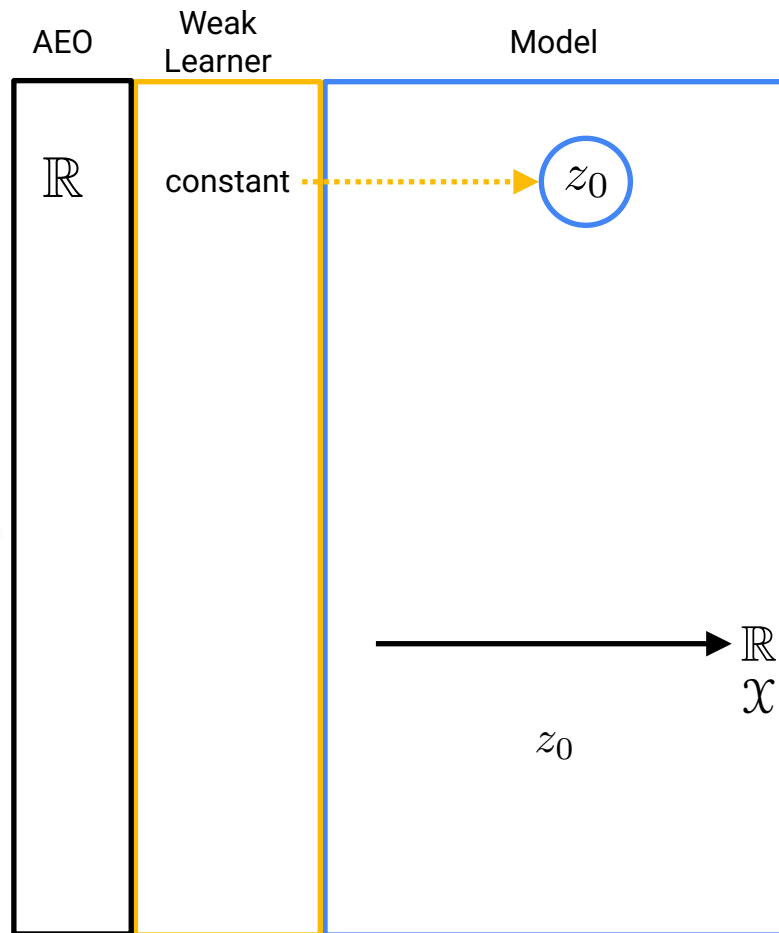# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t \mathbb{1}_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ $\mathbb{R}$ | constant | $z_0$ |

$$\longrightarrow \begin{array}{c} \mathbb{R} \\ \mathcal{X} \end{array}$$

$z_0$

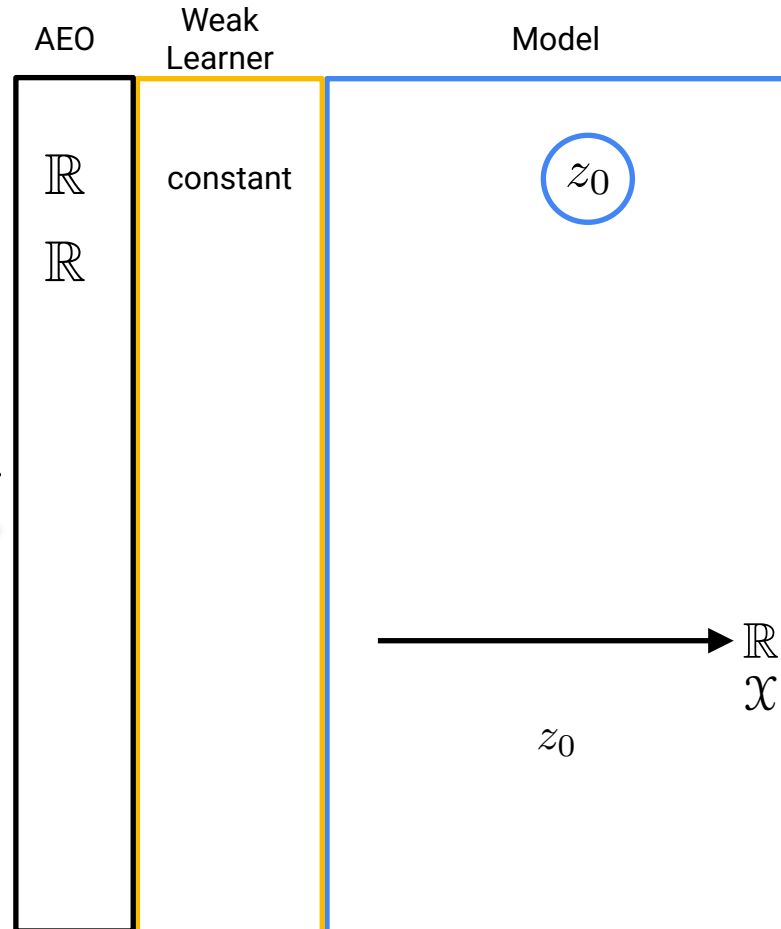# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random on $\mathcal{S}$
  $$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ | constant | $z_0$ |
| $\mathbb{R}$ | $1_{x \geq a_1} \cdot z_1$ | |

$\mathbb{R}$
$\mathcal{X}$

$z_0$

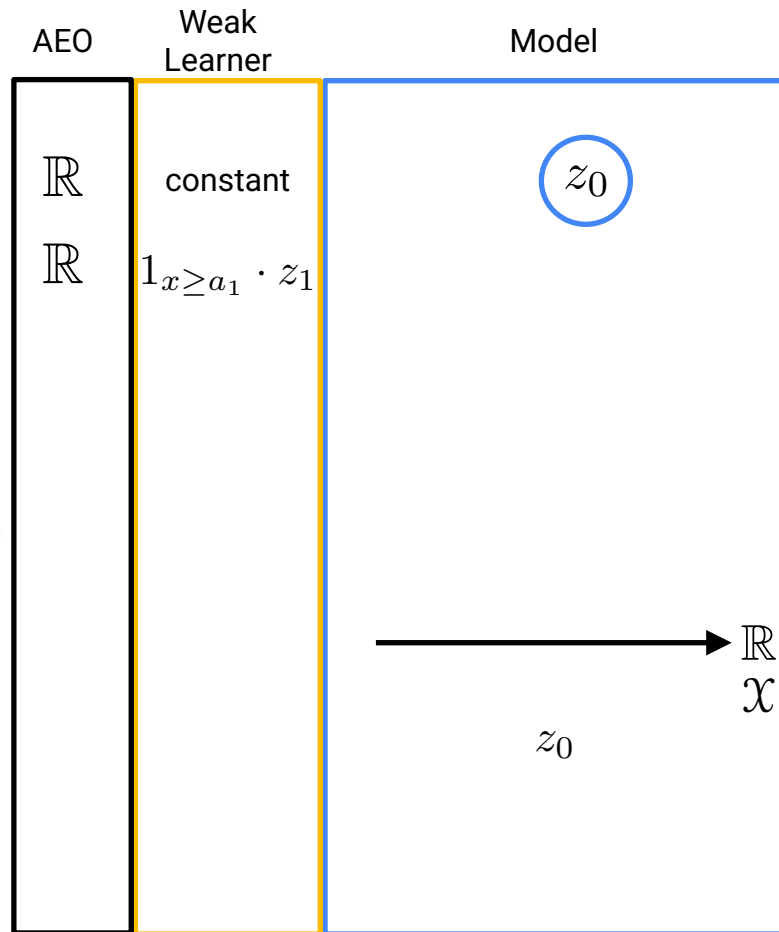# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↪ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances

- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$

- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random on $\mathcal{S}$
  $$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$

- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↪ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$



| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ | constant | |
| $\mathbb{R}$ | $1_{x \geq a_1} \cdot z_1$ | |

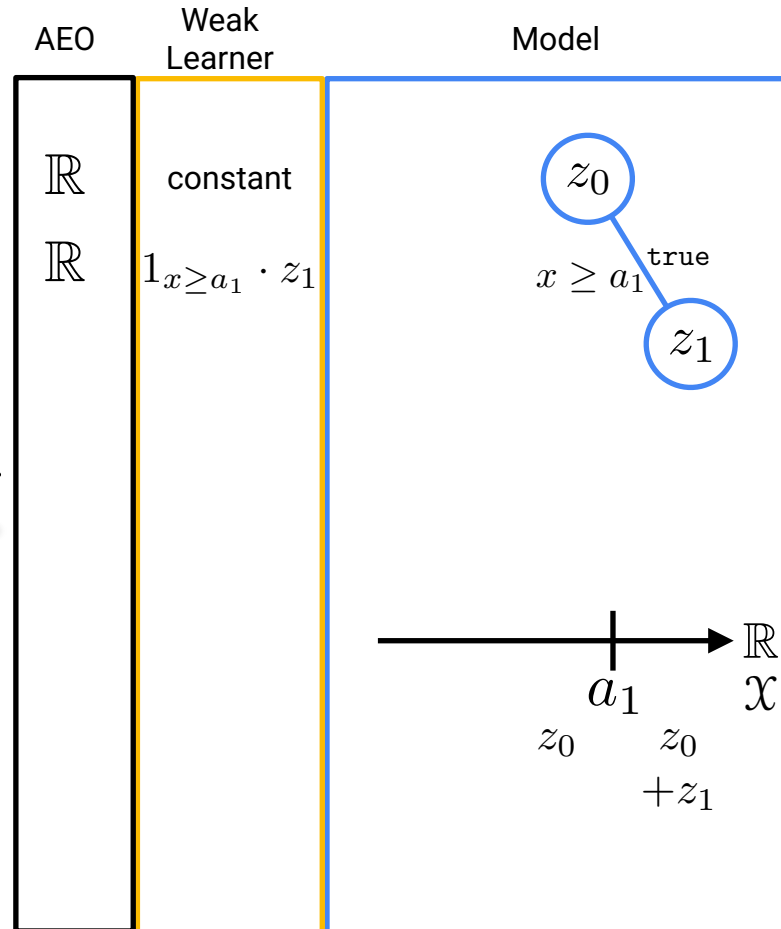# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

AEO

Weak Learner

Model

$\mathbb{R}$

$\mathbb{R}$

constant

$1_{x \geq a_1} \cdot z_1$

$1_{x < a_1} \cdot -z_1$

$z_0$

$x \geq a_1$ true

$z_1$

Meets the WLA

$\mathbb{R}$
$\mathcal{X}$

$a_1$

$z_0 \qquad z_0$

$+z_1$

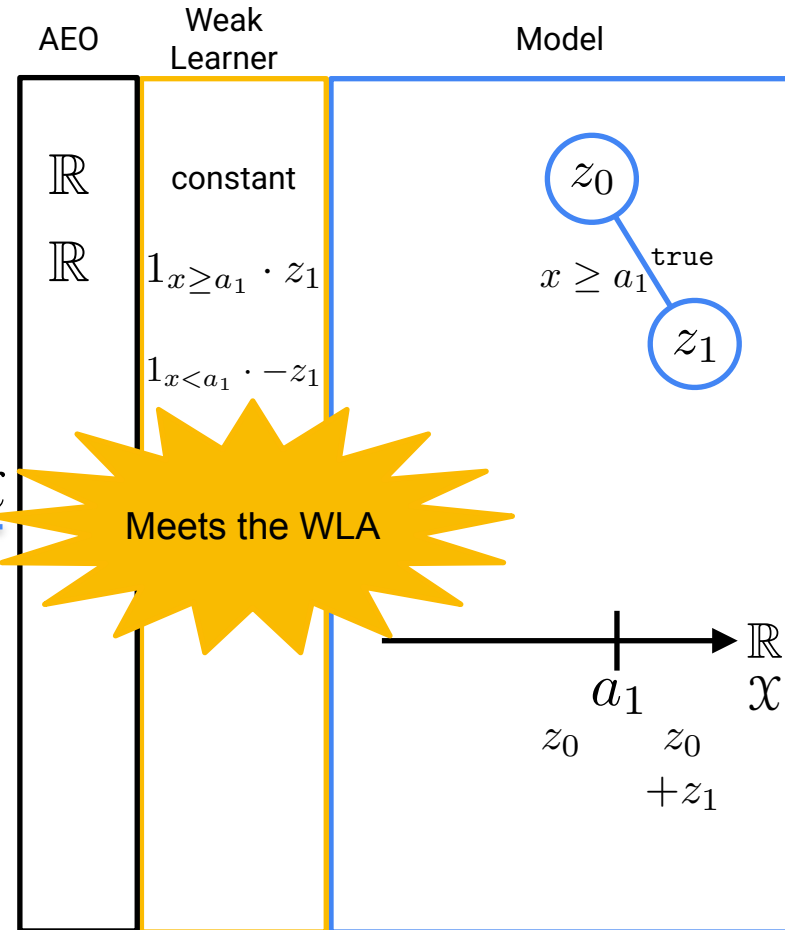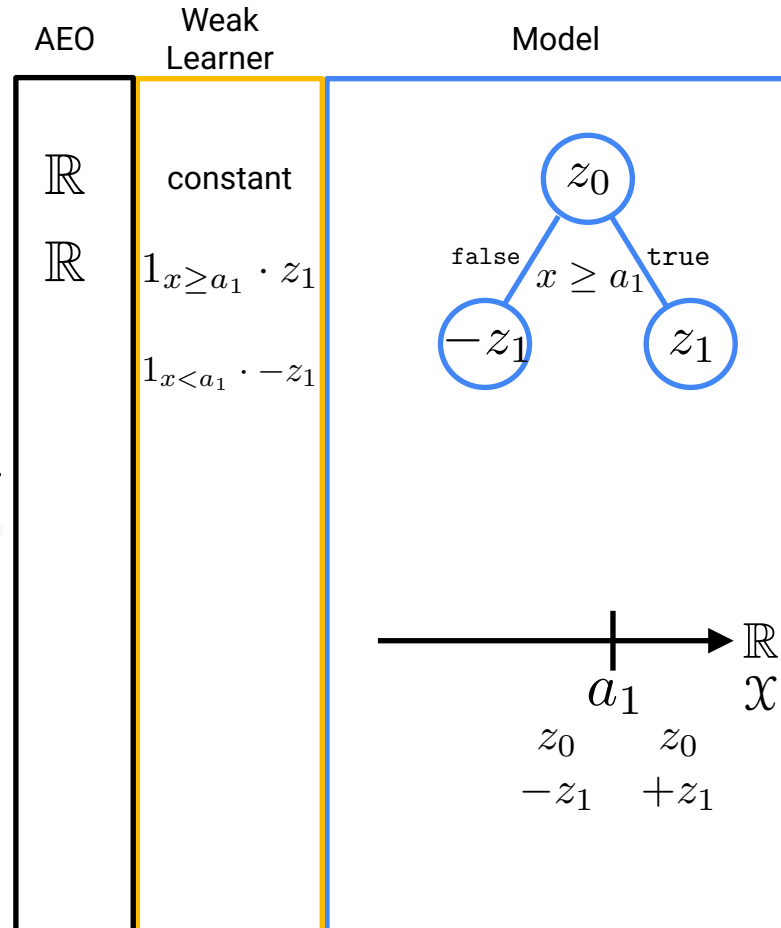# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$



AEO | Weak Learner | Model

$\mathbb{R}$

$\mathbb{R}$    constant

$1_{x \geq a_1} \cdot z_1$

$1_{x < a_1} \cdot -z_1$

$z_0$

false   $x \geq a_1$   true

$-z_1$    $z_1$

$\mathbb{R}$
$\mathcal{X}$

$a_1$

$z_0 \quad z_0$
$-z_1 \quad +z_1$

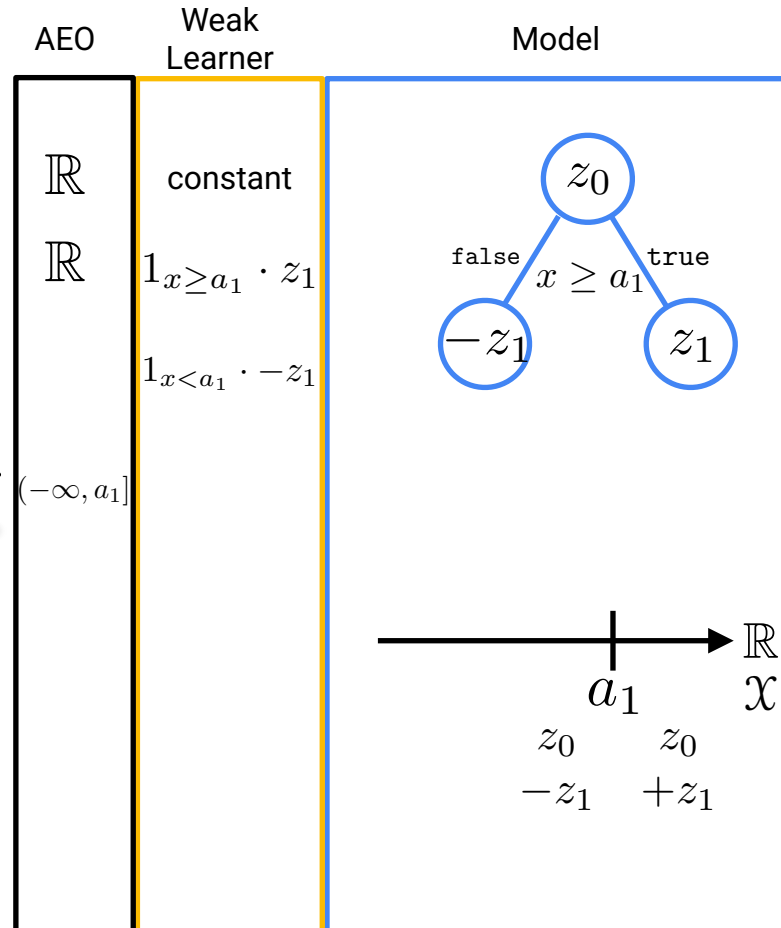# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↪ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↪ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ | constant | |
| $\mathbb{R}$ | $1_{x \geq a_1} \cdot z_1$ | |
| | $1_{x < a_1} \cdot -z_1$ | |
| $(-\infty, a_1]$ | | |

$z_0$

false $\quad x \geq a_1 \quad$ true

$-z_1 \qquad z_1$

$a_1$

$z_0 \qquad z_0$

$-z_1 \qquad +z_1$

$\mathbb{R}$
$\mathcal{X}$

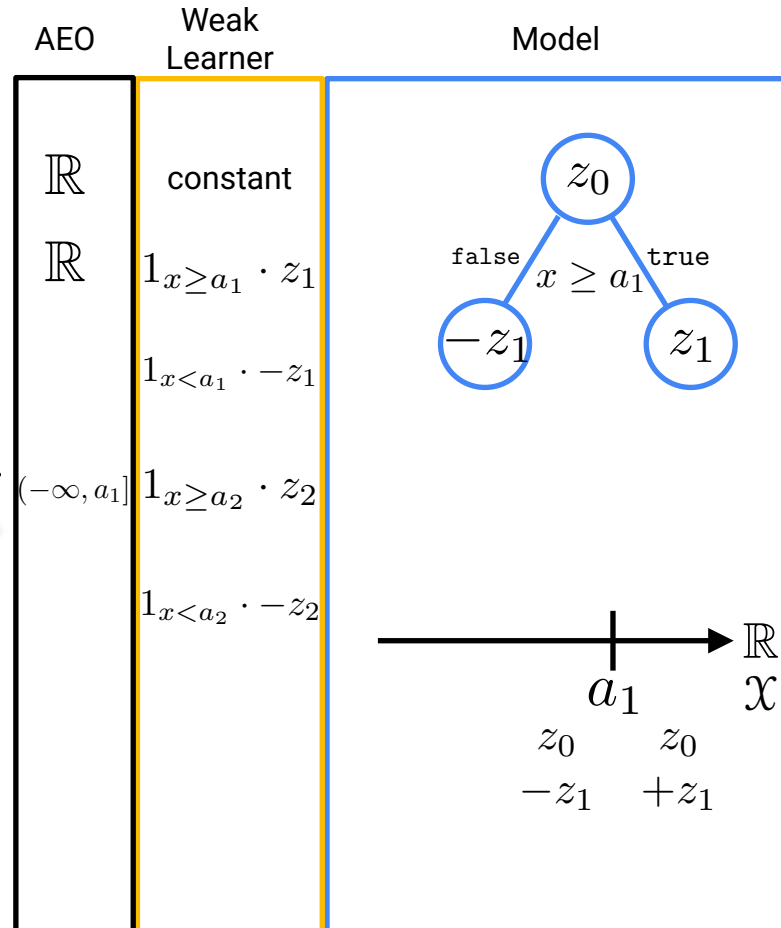# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

| AEO | Weak Learner | Model |
|---|---|---|
| $\mathbb{R}$ | constant | |
| $\mathbb{R}$ | $1_{x \geq a_1} \cdot z_1$ | |
| | $1_{x < a_1} \cdot -z_1$ | |
| $(-\infty, a_1]$ | $1_{x \geq a_2} \cdot z_2$ | |
| | $1_{x < a_2} \cdot -z_2$ | |

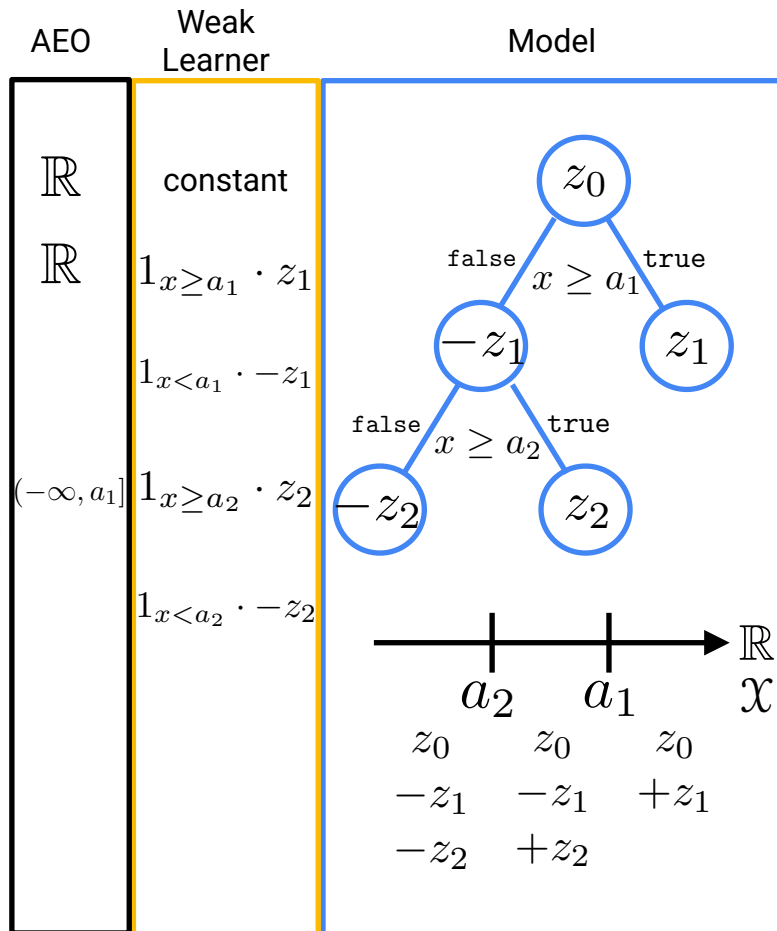# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\Longleftrightarrow$ to) a specific *model architecture*
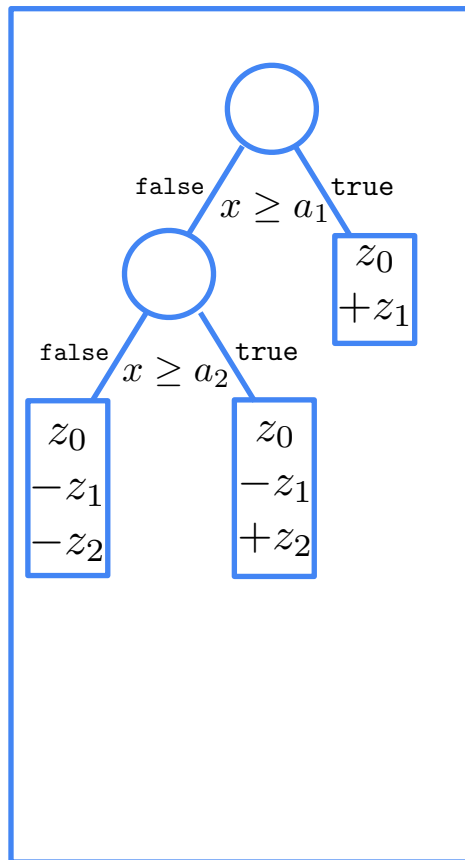
- Weights $\boldsymbol{w}$ = record of past performances
- <u>Architecture Emulation Oracle</u> : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random on $\mathcal{S}$
  $$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t \mathbb{1}_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↪ **Step 2**: introduce a new oracle ensuring the final emulates (is $\iff$ to) a specific *model architecture*

- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$
  at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↪ Returns model $H(\boldsymbol{x}) \doteq \sum_t \mathbb{1}_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

Equivalent representation

# Convex boosting: which models ?

ModaBoost (Model-Adaptive Boosting)

↳ **Step 2**: introduce a new oracle ensuring the final emulates (is $\iff$ to) a specific *model architecture*
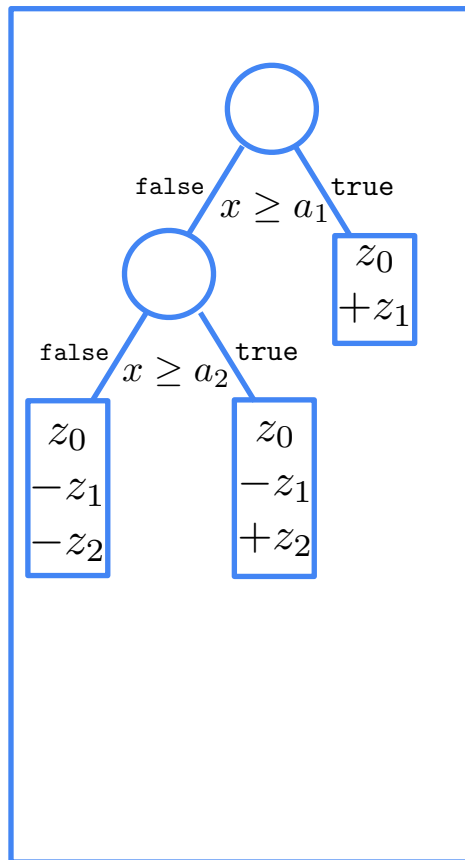
- Weights $\boldsymbol{w}$ = record of past performances
- Architecture Emulation Oracle : outputs $\mathcal{S} \subseteq \mathcal{X}$
- Weak learner : outputs hypotheses $h \in \mathbb{R}^{\mathcal{X}}$ at least $(\gamma > 0)$ different from random on $\mathcal{S}$
$$|\mathbb{E}_{\boldsymbol{w}_{|\mathcal{S}}}[y \cdot h(\boldsymbol{x})]| \geq \gamma$$
- Fits leveraging coefficients $\alpha \in \mathbb{R}$

↳ Returns model $H(\boldsymbol{x}) \doteq \sum_t 1_{\boldsymbol{x} \in \mathcal{S}_t} \alpha_t \cdot h_t(\boldsymbol{x})$

**Decision Tree**



$x \geq a_1$ — false / true

$z_0 + z_1$

$x \geq a_2$ — false / true

$z_0 - z_1 - z_2$

$z_0 - z_1 + z_2$

# Convex boosting with ModaBoost: which models ?

**Decision Trees**

# Convex boosting with ModaBoost: which models ?

**Decision Trees**

**Linear Separators**

# Convex boosting with ModaBoost: which models ?

**Decision Trees**

**Linear Separators**

**Alternating Decision Trees**

# Convex boosting with ModaBoost: which models ?

Decision Trees

Linear Separators

Alternating Decision Trees

Nearest Neighbors

# Convex boosting with ModaBoost: which models ?

**Decision Trees**

**Linear Separators**

**Alternating Decision Trees**

**Nearest Neighbors**

**Labeled Branching Programs**

...

# Convex boosting with ModaBoost vs Long & Servedio

**ModaBoost's output on Long & Servedio's setting**

Decision Trees

Linear Separators

Alternating Decision Trees

Nearest Neighbors

Labeled Branching Programs

# Convex boosting with ModaBoost vs Long & Servedio

**ModaBoost's output on Long & Servedio's setting**



Decision Trees
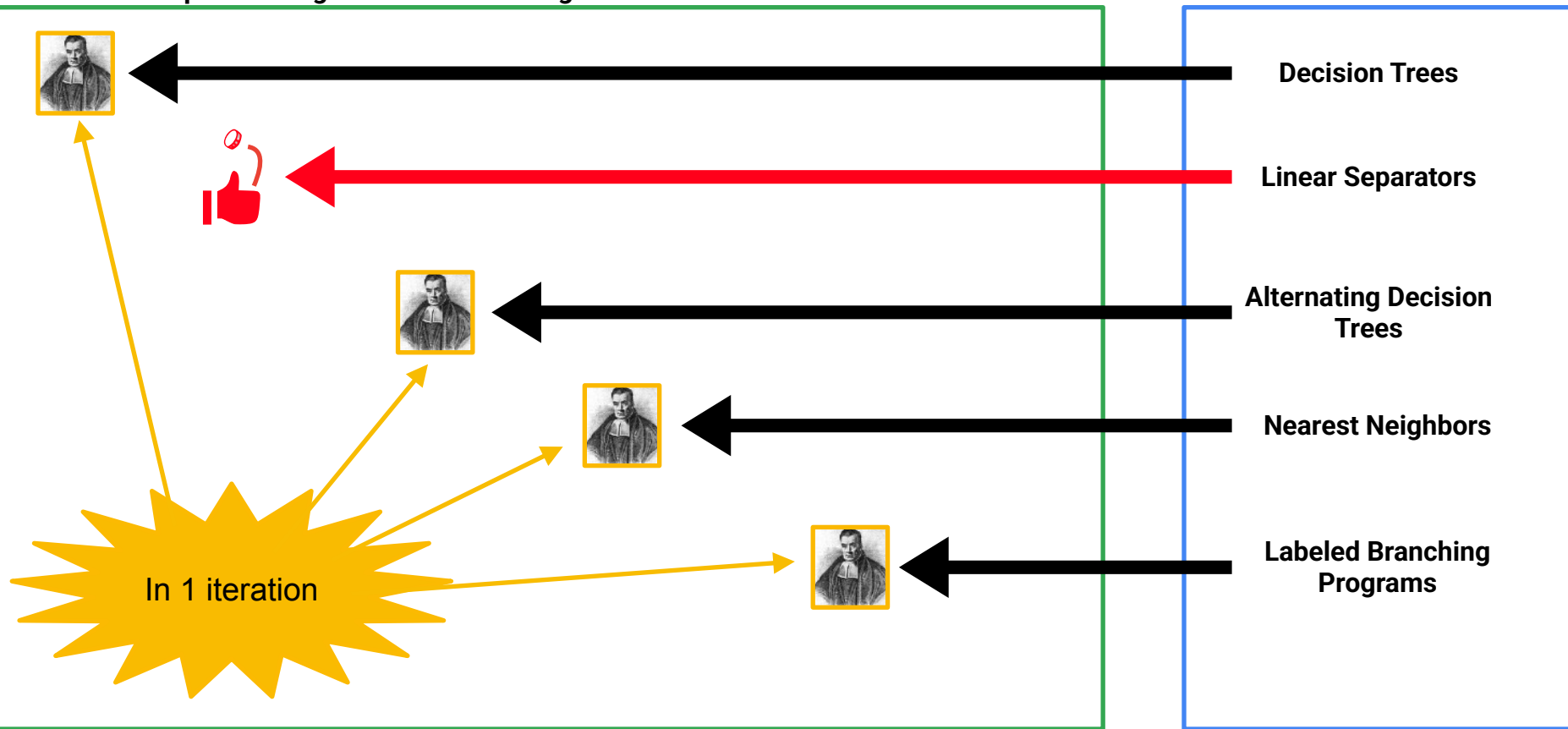
Linear Separators

Alternating Decision Trees

Nearest Neighbors

Labeled Branching Programs

# Convex boosting with ModaBoost vs Long & Servedio

**ModaBoost's output on Long & Servedio's setting**

# Convex boosting with ModaBoost vs Long & Servedio

**ModaBoost's output on Long & Servedio's setting**

# Convex boosting with ModaBoost vs Long & Servedio

**ModaBoost's output on Long & Servedio's setting**

# Convex boosting with ModaBoost vs Long & Servedio

**ModaBoost's output on Long & Servedio's setting**

# Conclusion

↳ Long and Servedio's paper has has a lasting impact on boosting / optimization

↳ Its impact should broaden on / shift to **models**, because it shows that

Linear Models can derail a whole ML pipeline otherwise optimal as soon as the "simplest" form of noise affects training data

↳ Suggests a broader question: given a class of models (more complex ?), what is its simplest "nemesis" noise model ?

Google Research

# Thank you !

Google Research