

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

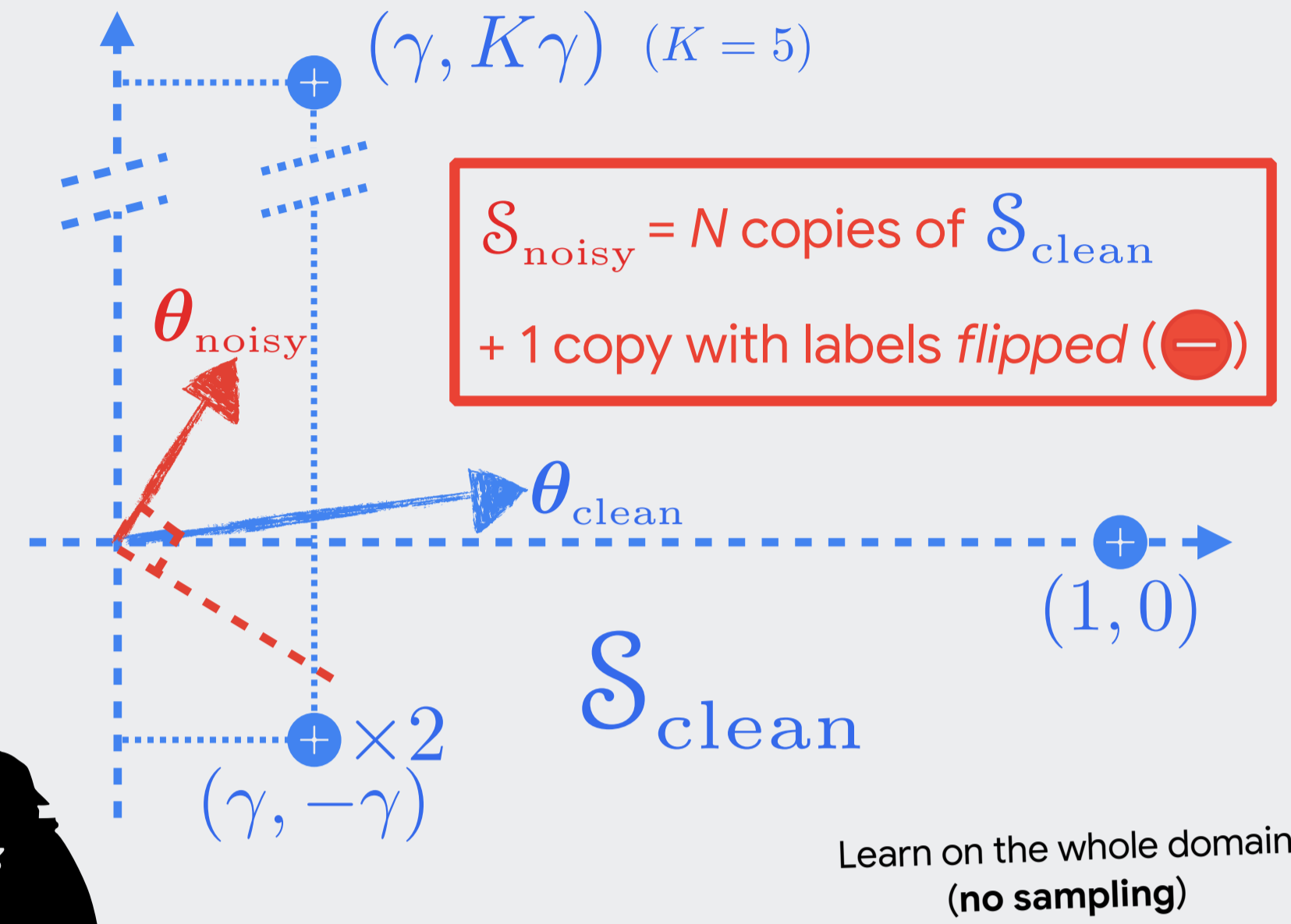
Yishay Mansour
Tel Aviv U. & Google Research

Richard Nock
Google Research

Robert C. Williamson
U. of Tübingen & Tübingen AI center

Why this work

- "Random Classification Noise Defeats All Convex Potential Boosters" -- Long and Servedio (ICML'08 & MLJ'10)
- 1-class trivial data, **noisify** it with symmetric label noise
- Compute the respective minimizers of any* **convex loss** or of any* convex **boosting** algorithm $(\theta_{\text{clean}}, \theta_{\text{noisy}})$
- On $\mathcal{S}_{\text{clean}}$, θ_{clean} is (predictably) perfect but $\theta_{\text{noisy}} = \text{!}$
- Sizeable impact on design of boosting algorithms,
- Recurrent stress of key culprit: convex boosting. Or is it?



The Setting: *properness*

Shufford et al., 1966; Savage, 1971

- Loss for Class Probability Estimation (CPE)

$$\ell(y^*, u) \doteq \begin{cases} [y^* = 1] \cdot \ell_1(u) \\ + \\ [y^* = -1] \cdot \ell_{-1}(u) \end{cases}$$

- Functions ℓ_1, ℓ_{-1} = partial losses
- Conditional risk of u wrt ground truth $v \in [0, 1]$:

$$L(u, v) \doteq v \cdot \ell_1(u) + (1 - v) \cdot \ell_{-1}(u)$$

- Bayes risk: $\underline{L}(v) \doteq \inf_u L(u, v)$

- A CPE loss is

- *symmetric* iff $\ell_1(u) = \ell_{-1}(1 - u), \forall u \in [0, 1]$
No class-dependent misclassification cost

- *differentiable* iff ℓ_1, ℓ_{-1} differentiable

- *lowerbounded* iff ℓ_1, ℓ_{-1} lowerbounded

- *proper* iff:

$$\forall v \in [0, 1], L(v, v) = \inf_u L(u, v)$$

- *strictly proper* iff v = sole minimizer

\forall proper loss, $\underline{L}(v) = v \cdot \ell_1(v) + (1 - v) \cdot \ell_{-1}(v)$

Ex. of **Strictly Proper Differentiable** - SPD - (and *symmetric*) losses: square, log / crossentropy, Matusita

ML: given sample $\mathcal{S} \doteq \{(\mathbf{x}_i, y_i^*)\}_{i=1}^m$, learn *posterior* $\tilde{\eta} : \mathcal{X} \rightarrow [0, 1]$ min. $\Phi(\tilde{\eta}, \mathcal{S}) \doteq \mathbb{E}_{i \sim [m]} [\ell(y_i^*, \tilde{\eta}(\mathbf{x}_i))]$

A "proper paradox"

- In Long and Servedio's setting, does not learn a posterior directly but real-valued *classifier* $h : \mathcal{X} \rightarrow \mathbb{R}$

Theorem (Nock & Menon, 2020): any SPD loss admits a dual form for real-valued classification,

$$\ell(y, h(\mathbf{x})) = -\underline{L}(y) + \phi_\ell(-h(\mathbf{x})) - yh(\mathbf{x})$$

$$\phi_\ell(z) \doteq (-\underline{L})^*(z)$$

& if *+symmetric*,

$$\ell(y^*, h(\mathbf{x})) = -\underline{L}\left(\frac{1 + y^*}{2}\right) + \phi_\ell(y^* h(\mathbf{x}))$$

- Lemma:** for any SPD+symmetric loss, ϕ_ℓ meets the blueprint to Long and Servedio's convex losses

On Long and Servedio's data, [boosted] optimum ends up eliciting **not Bayes rule**, but a **fair coin**!

... and there is "worse":

Theorem: Long and Servedio's results survive to dropping the *symmetry* constraint on SPD loss...

Properness virtually *useless* to learn, while it should lead to maximal accuracy on **noise-free data**...

→ class-dependent partial losses cannot solve problem → negative result holds for losses "without margin form" as well

The culprit & how to address the negative result

- Long and Servedio's setting relates to [boosted] optimum of a (margin) **convex** loss... but supervised ML involves training a **model** in the pipeline (Long and Servedio: linear)... so what if we *just replace the model*?
- +make it more general: any* **Strictly Proper Differentiable** loss (not necessarily *symmetric*, no margin form)

Algorithm 1 MODABOOST($\mathcal{S}, \ell, \text{WL}, \text{AEO}, T$)

Input: Dataset $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, SPD loss ℓ , weak learner WL, architecture emulation oracle AEO, iteration number $T \geq 1$;

Output: PLM H_T ;

Step 1 : $\forall i \in [m], w_{i,1} \doteq w((\mathbf{x}_i, y_i), H_0)$

Step 2 : **for** $t = 1, 2, \dots, T$

Step 2.1 : $\mathcal{X}_t \leftarrow \text{AEO}(\mathcal{X}, H_{t-1})$;

Step 2.2 : $h_t \leftarrow \text{WL}(w_t^*, \mathcal{S} \cap \mathcal{X}_t)$;

Step 2.3 : compute α_t as the solution to:

$$\sum_{i \in [m]_t} w((\mathbf{x}_i, y_i), H_t) \cdot y_i^* h_t(\mathbf{x}_i) = 0;$$

Step 2.4 : $\forall i \in [m]_t, w_{t+1,i} \doteq w((\mathbf{x}_i, y_i), H_t)$

return $H_T(\mathbf{x}) \doteq \sum_{t=1}^T [\mathbf{x} \in \mathcal{X}_t] \cdot \alpha_t h_t(\mathbf{x})$;

- PLM = Partition Linear Model:

$$H_t(\mathbf{x}) \doteq \sum_{t=1}^T [\mathbf{x} \in \mathcal{X}_t] \cdot \alpha_t h_t(\mathbf{x})$$

- h_t returned by classical Weak Learner
- \mathcal{X}_t by *Architecture Emulation Oracle*

- AEO [+WL] chosen so that H_T emulates a specific model architecture (if AEO returns \mathcal{X} , learns linear model)

- Weight function uses both class notations:

$$w((\mathbf{x}, y), H) \doteq y - y^* \cdot (-\underline{L}')^{-1}(H(\mathbf{x}))$$

Weak learning assumption (WLA): any h_t satisfies

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

AEO Compliance (AEOC): any \mathcal{X}_t satisfies

$$J([\mathbf{m}]_t, t) \geq u_t \cdot J([\mathbf{m}], t) \quad \text{for some } u_t > 0$$

$J(\mathcal{W}, t) \doteq \text{Card}(\mathcal{W}) \cdot (\mathbb{E}_{i \sim \mathcal{W}} [w_{t,i}])^2$
"sufficiently many significant weights chosen"

Theorem : Suppose WLA + AEOC hold & the SPD loss has $\ell_{-1}(0) \geq C, \ell_1(1) \geq C, \inf\{\ell'_{-1} - \ell'_1\} \geq \kappa$. Then for any $\theta \geq 0, \varepsilon > 0$, the following holds for ModaBoost's output H_T :

$$\left(T \geq U^{-1} \left(\frac{2(\Phi(H_0, \mathcal{S}) - C)}{\kappa \cdot \varepsilon^2 \underline{w}(\theta)^2 \gamma_{\text{WL}}^2} \right) \right) \Rightarrow (\mathbb{P}_{i \sim [m]} [y_i^* H_T(\mathbf{x}_i) \leq \theta] < \varepsilon)$$

$$U(T) \doteq \sum_{t=1}^T u_t \quad \underline{w}(\theta) = \min\{1 - (-\underline{L}')^{-1}(\theta), (-\underline{L}')^{-1}(-\theta)\}$$

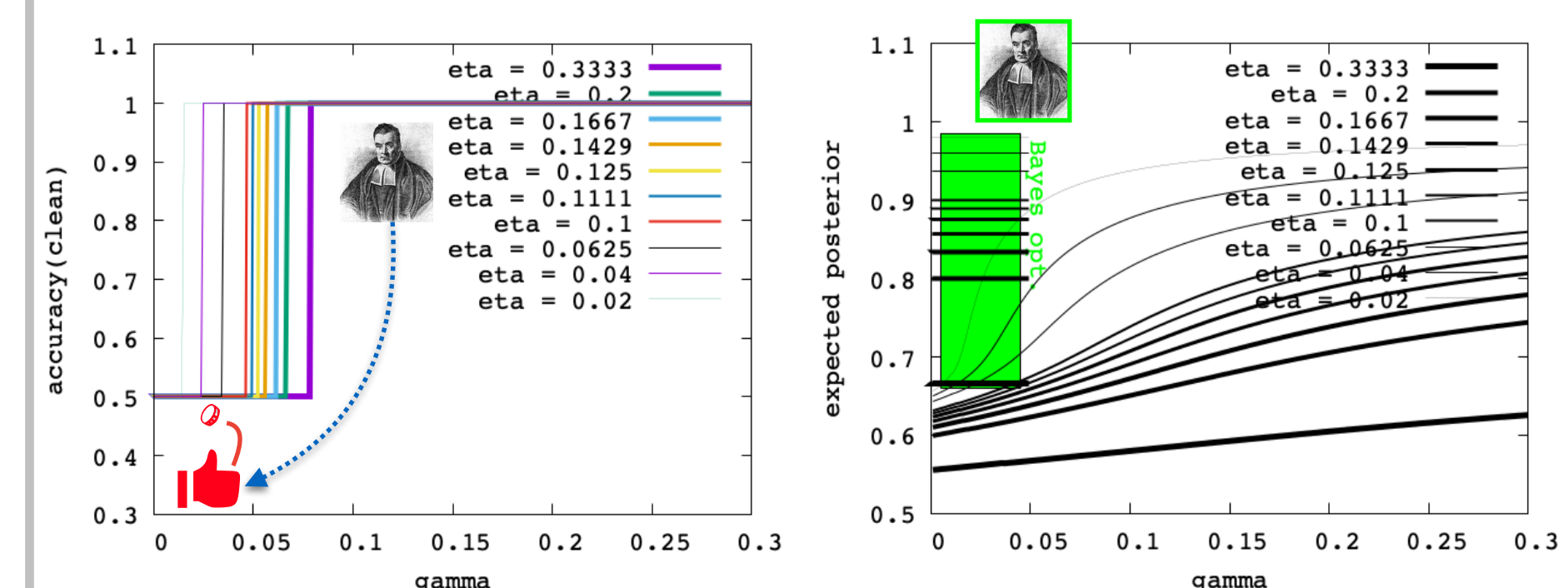
ML models "boostable" by emulation from ModaBoost (see the paper for translated boosting rates) include: Linear separators, Decision trees, Alternating decision trees, Nearest neighbors, Labeled branching programs

- ModaBoost stands an independent contribution to our work, so what does it bring in the context of our paper?

Lemmata: for *any* values of the triple $(N, K, \gamma) \in \mathbb{N}_{>0}^2 \times \mathbb{R}_{>0}$ in Long and Servedio's data, **ModaBoost, trained on $\mathcal{S}_{\text{noisy}}$, is Bayes optimal in 1 iteration** if it emulates any of [alternating] decision trees, nearest neighbors or labeled branching programs... but it can so blatantly fail with linear separators that it can hit fair coin prediction on $\mathcal{S}_{\text{clean}}$ in just 2 iterations (depending on loss)

Toy Experiment

- ModaBoost + linear separators on Long and Servedio's data, loss = Matusita



Conclusion: mind the parameterization !