



## TL;DR

We train ensemble models with high tail performance using a novel Boosting framework that boosts an unfair learner to a fair learner and demonstrate its efficiency.

## CVaR: Conditional Value at Risk

### Fair Models with High Tail Performance

A fair model should have high **tail performance**, *i.e.* high performance on samples where the performance is the lowest. It can be measured by the  $\alpha$ -CVaR loss:

$$\text{CVaR}_\alpha^\ell(F) = \max_{\mathbf{w} \in \Delta_n, \mathbf{w} \preceq (\alpha n)^{-1}} \sum_{i \in [n]} w_i \ell(F(\mathbf{x}_i), y_i) \quad (1)$$

where  $\{(\mathbf{x}_i, y_i)\}_{i \in [n]}$  is the training set,  $\alpha \in (0, 1)$ ,  $\ell(\hat{y}, y)$  is a loss function and  $\Delta_n$  is the unit simplex in  $\mathbb{R}^n$ .

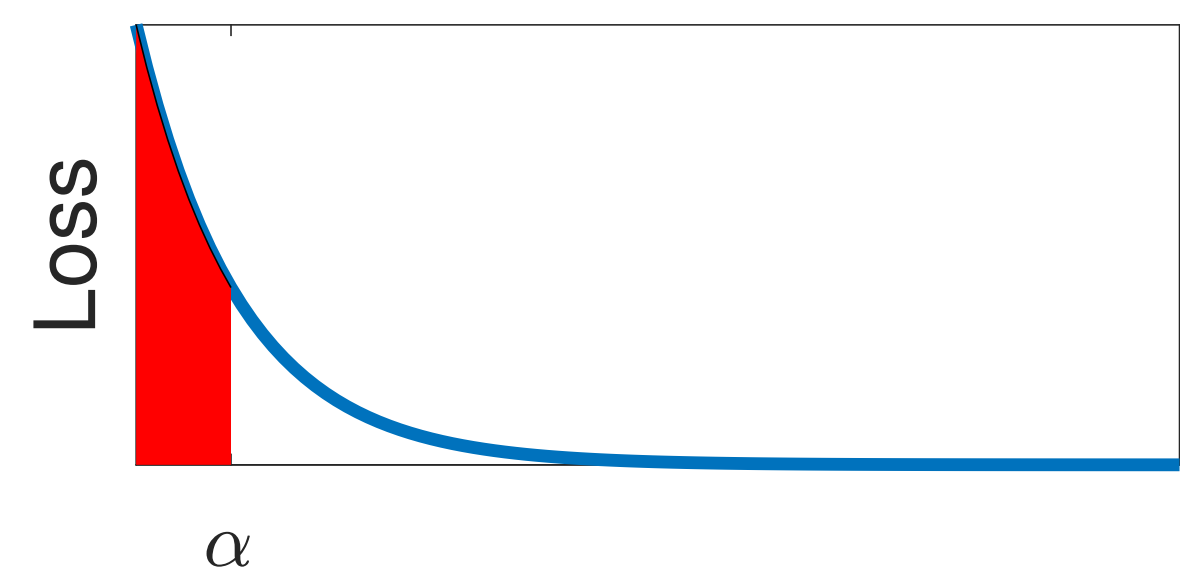


Figure 1:  $\alpha$ -CVaR loss (red region) is the average loss over the worst  $\alpha$  fraction of the data.

Fraction of Data

### ERM Achieves Lowest CVaR in Classification

In classification tasks, we evaluate the models with the zero-one loss  $\ell_{0/1}(\hat{y}, y) = \mathbf{1}_{\{\hat{y} \neq y\}}$ . We can prove that

$$\text{CVaR}_\alpha^{\ell_{0/1}}(F) = \min\{1, \alpha^{-1} \hat{\mathcal{R}}^{\ell_{0/1}}(F)\} \quad (2)$$

where  $\hat{\mathcal{R}}^{\ell_{0/1}}$  is the empirical zero-one loss which ERM minimizes. Thus, ERM also minimizes the CVaR loss, meaning that using CVaR has no gain compared to ERM.

## Training Ensemble Models via Boosting

Eqn. (2) only holds for deterministic models. For randomized models whose outputs are distributions over the output space, using CVaR can improve the tail performance.

### $\alpha$ -CVaR is Equivalent to $\alpha$ -LPBoost

$\alpha$ -LPBoost refers to the following primal/dual LP:

**Dual:**

$$\begin{aligned} \min_{\mathbf{w}, \gamma} \quad & \gamma \\ \text{s.t.} \quad & \langle \mathbf{w}, \ell^s \rangle \geq 1 - \gamma; \forall s \in [t] \\ & \mathbf{w} \in \Delta_n, \mathbf{w} \preceq \frac{1}{\alpha n} \end{aligned}$$

**Primal:**

$$\begin{aligned} \max_{\lambda, \rho} \quad & \rho - \frac{1}{\alpha n} \sum_{i=1}^n (\rho - 1 + \sum_{s=1}^t \lambda_s \ell_i^s)_+ \\ \text{s.t.} \quad & \lambda \in \Delta_t \end{aligned}$$

where  $\ell_i^s$  is the loss of  $f^s$  on sample  $i$ . We can show that:

**Theorem.** The optimal objective value of this LP (same for primal/dual) is equal to  $1 - \min_{\lambda \in \Delta_t} \text{CVaR}_{\alpha}^{\ell_{0/1}}(F)$ , where  $F$  is the ensemble model consisting of  $f^1, \dots, f^t$  and  $\lambda$ .

Thus, minimizing the CVaR loss is equivalent to maximizing the optimal objective of  $\alpha$ -LPBoost, which is equivalent to the problem of **boosting an unfair learner**.

### $\alpha$ -AdaLPBoost: Improving Efficiency with AdaBoost

In  $\alpha$ -LPBoost, for each different  $\alpha$  we need a different set of base models  $\{f^t\}_{t \in [T]}$ . However, in many real tasks we need to tune  $\alpha$  frequently, which would be very inefficient.

To improve efficiency, we pick  $\mathbf{w}^t$  with AdaBoost:

$$w_i^{t+1} \propto \exp\left(\eta \sum_{s=1}^t \ell_i^s\right) \quad (3)$$

Then, for all  $\alpha$  we have the same set of base models, which makes tuning  $\alpha$  much more efficient.

## Boosting an Unfair Learner: Framework

We have an **unfair learner**  $\mathcal{L}$  that outputs models with high average performance but low tail performance.

For  $t = 1, \dots, T$  do

1. Pick a sample weight vector  $\mathbf{w}^t \in \Delta_n$  and feed it to  $\mathcal{L}$ .
2. Receive a base model  $f^t$  from  $\mathcal{L}$  whose weighted average 0/1 loss *w.r.t.*  $\mathbf{w}^t \leq g$  for constant  $g \in (0, 1)$ .

Finally, pick a model weight vector  $\lambda \in \Delta_T$ .

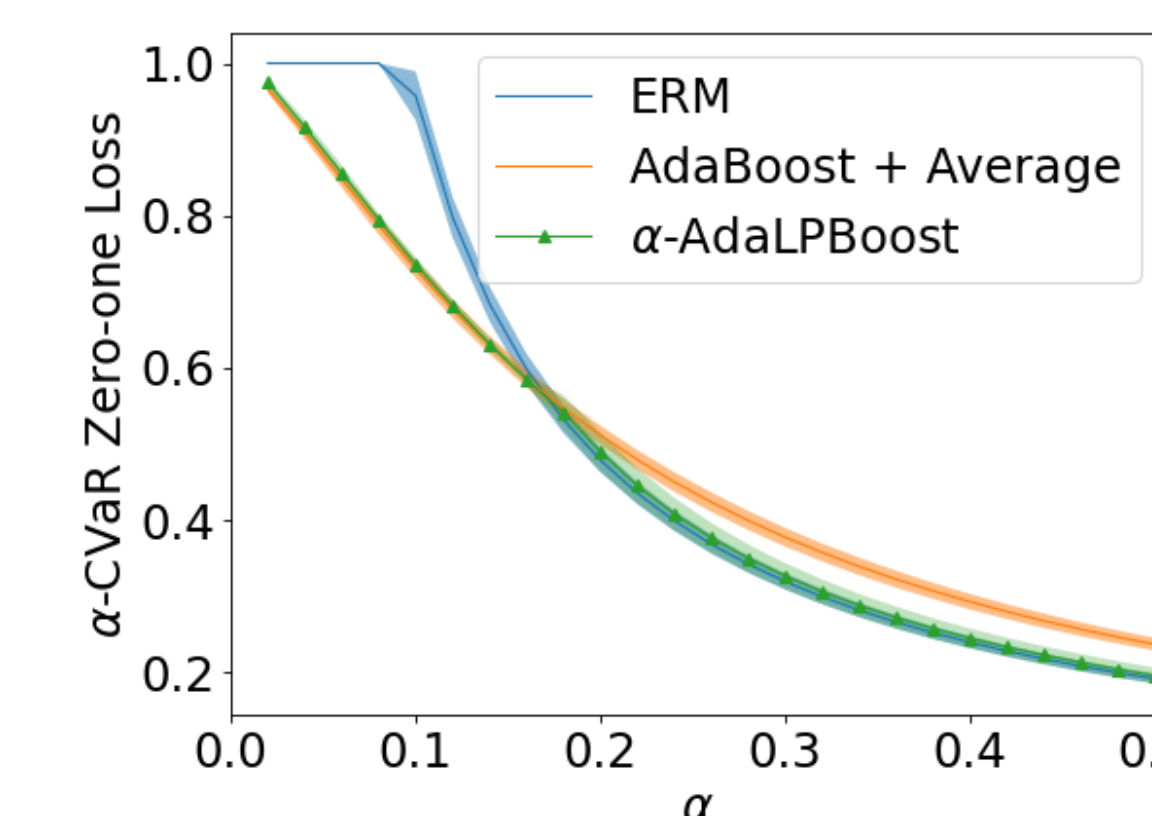
At inference time, first randomly sample a base model  $f^t$  according to  $\lambda$ , and then predict with  $f^t$ .

## Theoretical and Empirical Results

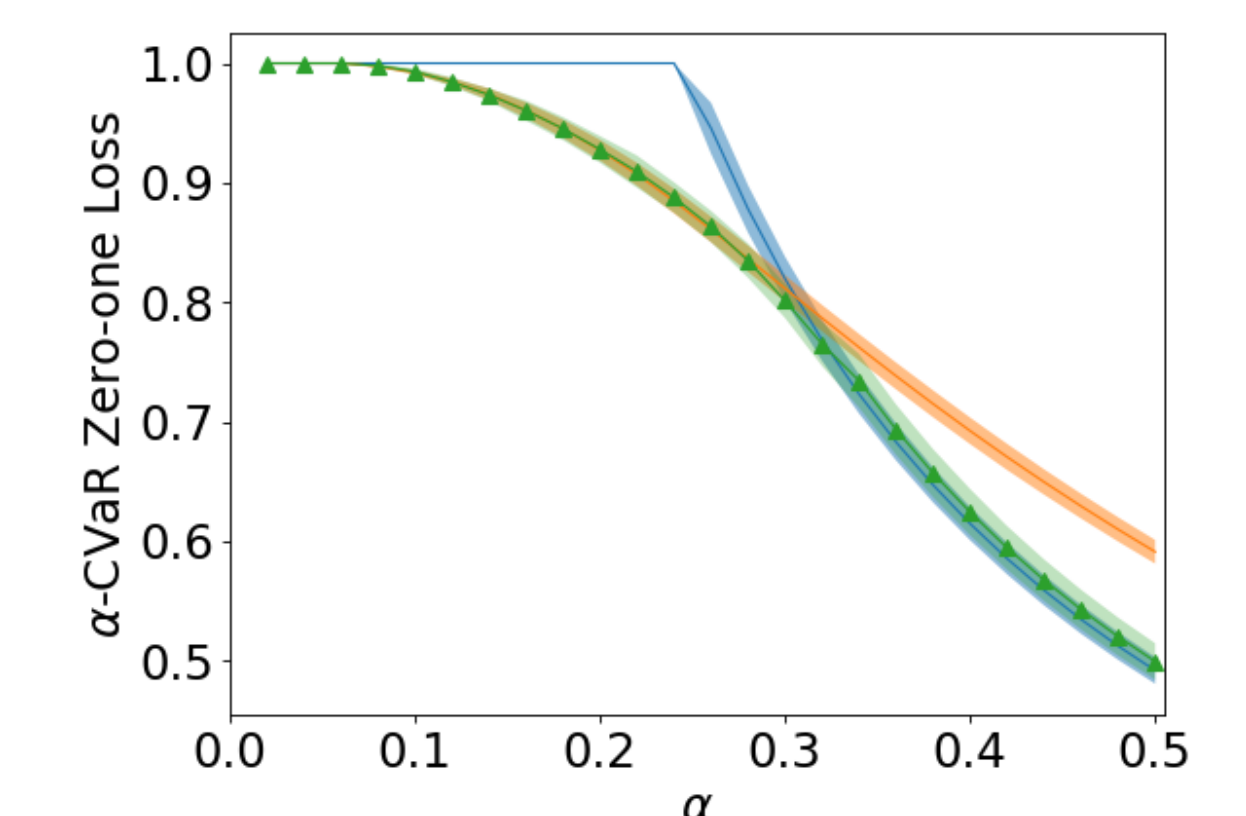
**Theorem.** For any  $\delta > 0$ , and for  $T = O(\frac{\log n}{\delta^2})$ , the training  $\alpha$ -CVaR zero-one loss of the ensemble model given by  $\alpha$ -AdaLPBoost is at most  $g + \delta$  if we set  $\eta = \sqrt{8 \log n / T}$ .

In the worst case, the zero-one loss of the ensemble model is at least  $g$ . This theorem shows that  $\alpha$ -AdaLPBoost can get as close to this lower bound as possible.

We also empirically show the efficiency of the framework:



(a) CIFAR-10



(b) CIFAR-100