

Boosting with Tempered Exponential Measures

Richard Nock

Ehsan Amid

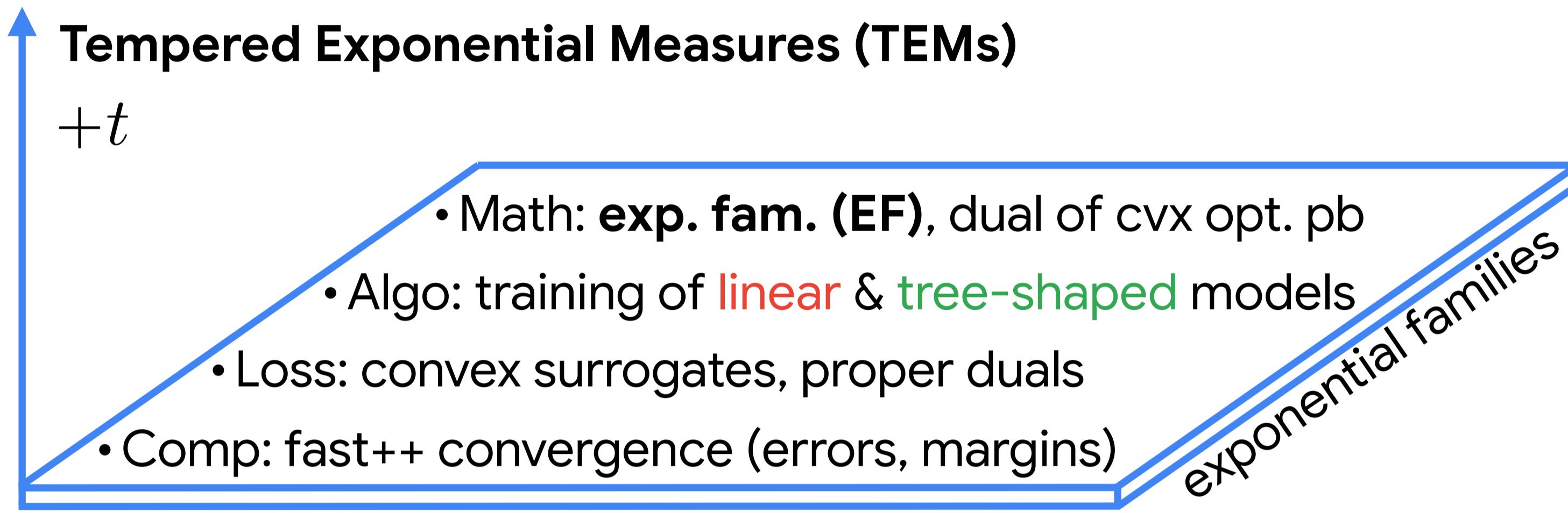
Manfred K. Warmuth

Summary

- Boosting fits parameters of an **exponential family (EF)**

$$\text{pdf}(\mathbf{x}) \propto \exp(\boldsymbol{\mu}^\top \boldsymbol{\Upsilon}(\mathbf{x}))$$

- We generalize it to a superset recently introduced, with ML benefits:



TEMs 101

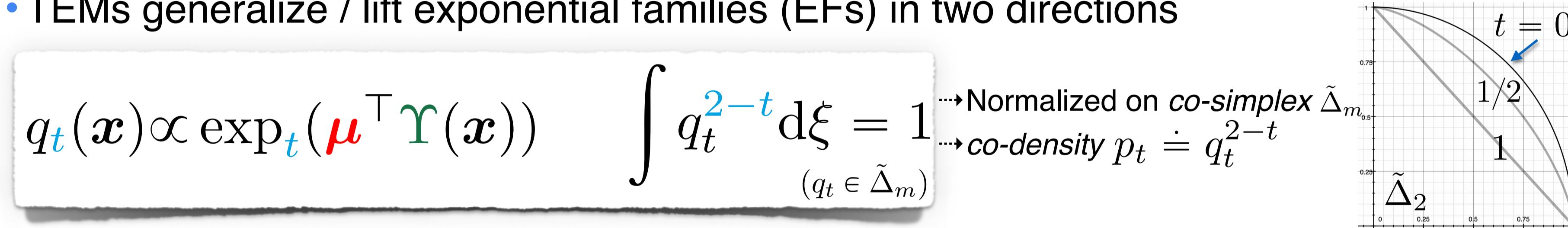
Amid, Nock & Warmuth, 2023

- t -logarithm, t -exponential, t -arithmetic (here, $t \in [0, 1]$)

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1), \quad \exp_t(z) \doteq [1 + (1-t)z]_+^{1/(1-t)}, \quad a \otimes_t b \doteq [a^{1-t} + b^{1-t} - 1]_+^{\frac{1}{1-t}}$$

$[z]_+ \doteq \max\{0, z\}$ ($t \rightarrow 1$: become "log", "exp" and ".")

- TEMs generalize / lift exponential families (EFs) in two directions

The t -Boosting TEM

Boosting's EF: Kivinen & Warmuth, 1999

- We seek the following TEM (update)

$$q' \doteq \arg \min_{\tilde{\mathbf{q}} \in \tilde{\Delta}_m} D_t(\tilde{\mathbf{q}} \parallel \mathbf{q}) \quad D_t(\mathbf{q} \parallel \mathbf{q}') \doteq \sum_{i \in [m]} q_i \cdot (\log_t q_i - \log_t q'_i) - \log_{t-1} q_i + \log_{t-1} q'_i$$

tempered relative entropy (generalizes KL divergence)
edge vector $u_i \doteq y_i h(\mathbf{x}_i)$ (m examples, supervised learning +1/-1)

- Theorem: for $t \in \mathbb{R}_{\geq 0} \setminus \{2\}$,

• solutions have the form $q'_i = \frac{q_i \otimes_t \exp_t(-\boldsymbol{\mu}^\top u_i)}{Z_t}$ $\boldsymbol{\mu} = \arg \min Z_t(\boldsymbol{\mu}')$ $\boldsymbol{q}'(\boldsymbol{\mu})^\top \mathbf{u} = 0$ (*) +light assumption if $t = 0$

• $Z_t(\boldsymbol{\mu}') \doteq \|\mathbf{q} \otimes_t \exp_t(-\boldsymbol{\mu}' \cdot \mathbf{u})\|_{2-t}$

Fitting the model part

- Models learned

$$\mathbf{H}_J(\mathbf{x}) \doteq \sum_{j \in [J]} \alpha_j h_j(\mathbf{x}) \quad \mathbf{H}_J^{(\delta)}(\mathbf{x}) \doteq \sum_{j \in [J]} \alpha_j h_j^{(\delta)}(\mathbf{x})$$

linear model clipped linear model decision tree

- Clipped summation

$$\sum_{j \in [J]} v_j \doteq \min \left\{ \delta, \max \left\{ -\delta, v_J + \sum_{j \in [J-1]} v_j \right\} \right\} \quad (\in [-\delta, \delta])$$

(non commutative, "encoding-nice")

Example: $a = -1, b = 3, \delta = 2$
 $\rightarrow v_1 = a, v_2 = b$
 Clipped sum is $2 = -1 + 3$
 $\rightarrow v_1 = b, v_2 = a$
 Clipped sum is $1 = 2 - 1$

- Training the linear part: t -AdaBoost

Algorithm t -ADABoost(t, \mathcal{S}, J)

Input: $t \in [0, 1]$, training sample \mathcal{S} , #iterations J ;

Output: classifiers $\mathbf{H}_J, \mathbf{H}_J^{(1-t)}$;

Step 1 : initialize tempered weights: $\mathbf{q}_1 = (1/m^{1/(2-t)}) \cdot \mathbf{1}$ ($\in \tilde{\Delta}_m$);

Step 2 : for $j = 1, 2, \dots, J$

Step 2.1 : get weak classifier $h_j \leftarrow \text{weak_learner}(\mathbf{q}_j, \mathcal{S})$;

Step 2.2 : choose weight update coefficient $\mu_j \in \mathbb{R}$;

Step 2.3 : $\forall i \in [m]$, for $u_{ji} \doteq y_i h_j(\mathbf{x}_i)$, update tempered weights:

$$q_{(j+1)i} = \frac{q_{ji} \otimes_t \exp_t(-\mu_j u_{ji})}{Z_{tj}}, \quad \text{where } Z_{tj} = \|\mathbf{q}_j \otimes_t \exp_t(-\boldsymbol{\mu}_j \cdot \mathbf{u}_j)\|_{2-t}.$$

Step 2.4 : choose leveraging coefficient $\alpha_j \in \mathbb{R}$;

(remark that we allow $\alpha_j \neq \mu_j$)

Training the decision tree (DT) part = top-down splitting s with a twist-on-the-loss

- $(t = 1)$ AdaBoost \rightarrow the most efficient splitting criterion for DT induction (Matushita's loss)
- We generalize to $t \neq 1$ and elicit a wide family of new losses for posterior estimation
- Process summarized (see paper for details):

models change $\rightarrow L^{(t)}$ variational form $\rightarrow \ell_1^{(t)}, \ell_{-1}^{(t)}$

Z_{tj}	tempered exponential loss real-valued classification	pointwise Bayes risk class-probability estimation	partial losses
----------	---	--	----------------

\rightarrow tempered loss ($t \in [-\infty, 2)$)

$$\ell_1^{(t)}(u) \doteq \left(\frac{1-u}{M_{1-t}(u, 1-u)} \right)^{2-t} \quad \ell_{-1}^{(t)}(u) \doteq \ell_1^{(t)}(1-u)$$

\rightarrow Power mean

$$M_q(a, b) \doteq \left(\frac{a^q + b^q}{2} \right)^{1/q}$$

\rightarrow $t=1$: Matushita, $t=0$: Gini, $t=\infty$: error

Properties

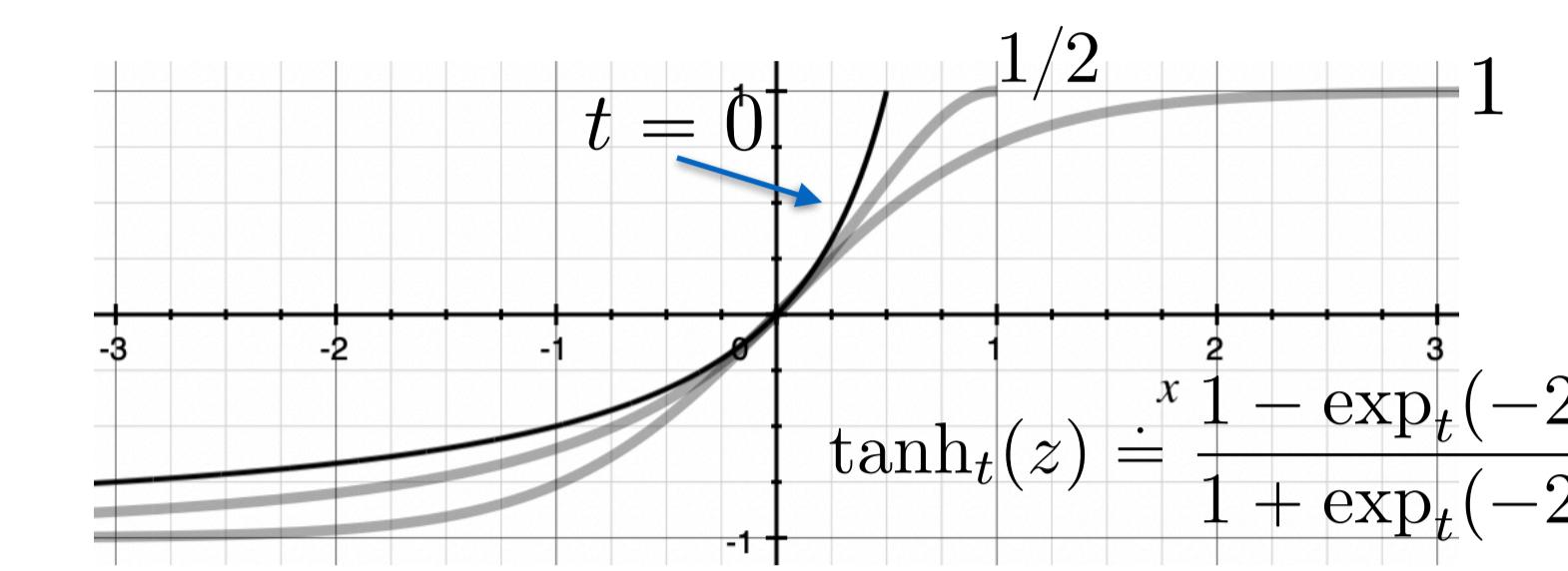
- (t)-margins

empirical margin risk

$$F_{t,\theta}(H, \mathcal{S}) \doteq \mathbb{E}_i[\nu_t((\mathbf{x}_i, y_i), H) \leq \theta]$$

training sample

$$\nu_t((\mathbf{x}, y), H) \doteq \tanh_t(yH(\mathbf{x})/2)$$



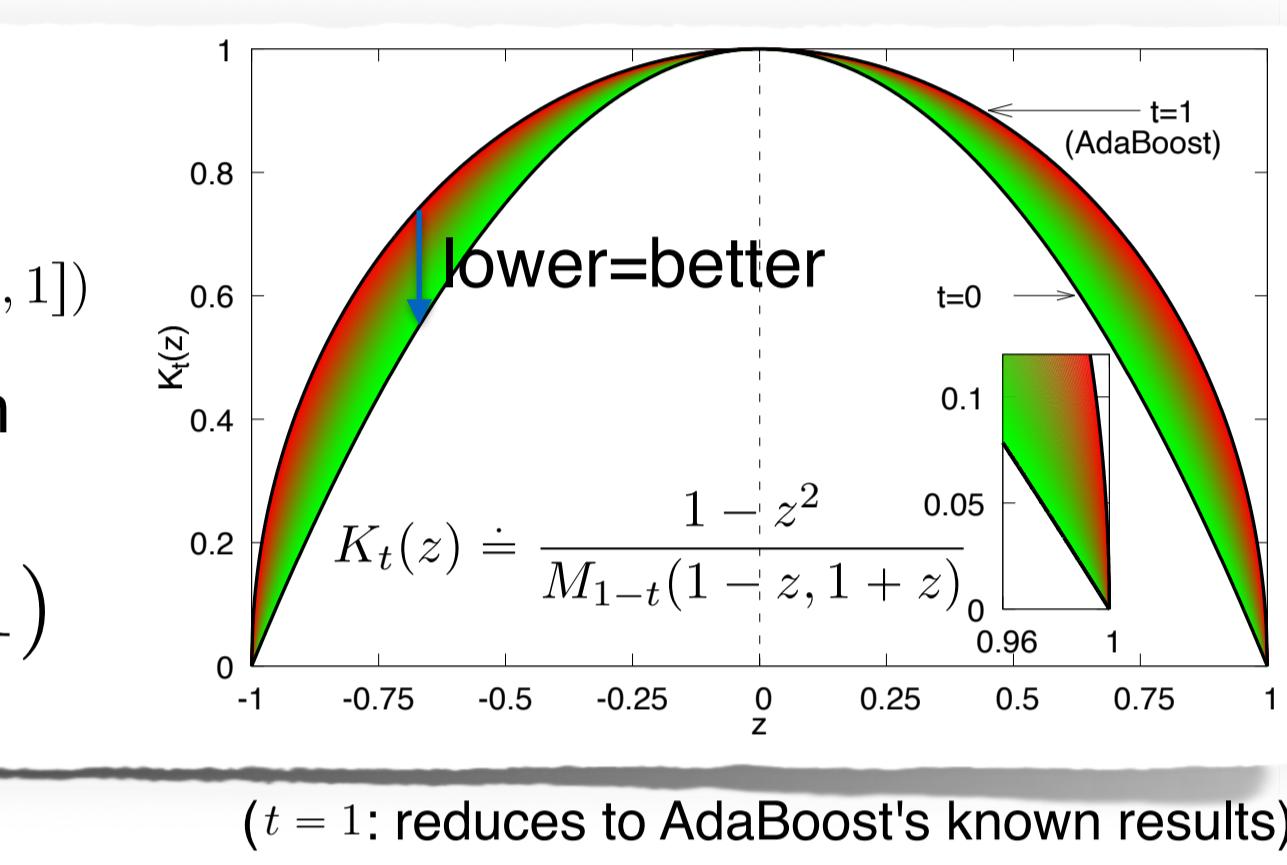
- Theorem (fast convergence for margins of t -AdaBoost, simplified), for $t \in [0, 1]$

In t -AdaBoost, fix

$$\alpha_j \propto \mu_j \quad \mu_j \propto -\log_t \left(\frac{1-\rho_j}{M_{1-t}(1-\rho_j, 1+\rho_j)} \right) \quad \rho_j \propto \mathbb{E}_{\mathbf{q}_j} [y_i h_j(\mathbf{x}_i)]_{i \in [-1, 1]}$$

If there is no "forgetting weights" (e.g. $H_J, H_J^{(1-t)}$ not good enough) then

$$F_{t,\theta}(H, \mathcal{S}) \leq \left(\frac{1+\theta}{1-\theta} \right)^{2-t} \cdot \prod_{j=1}^J K_t(\rho_j) \quad \text{for any } \theta \in (-1, 1)$$

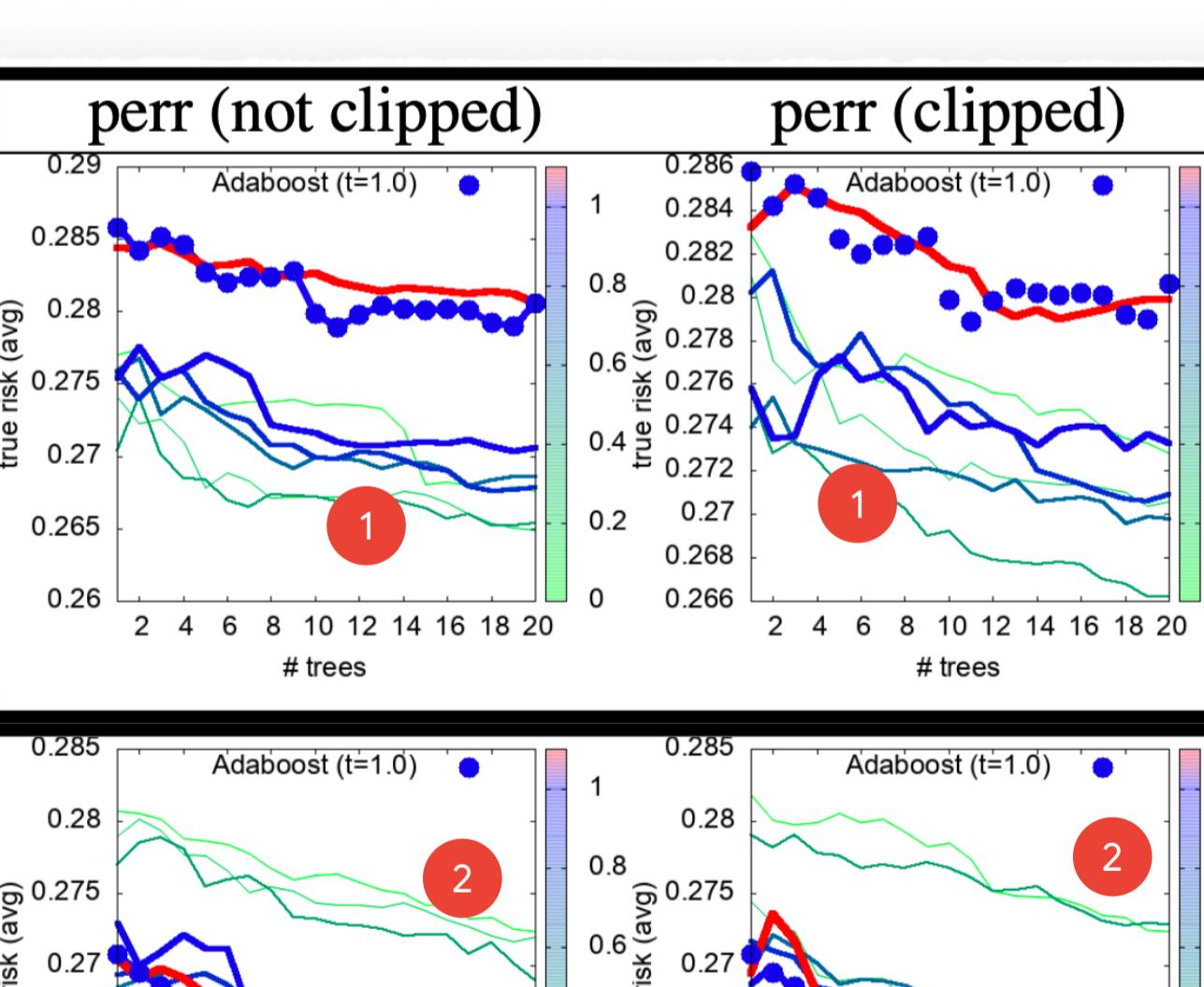
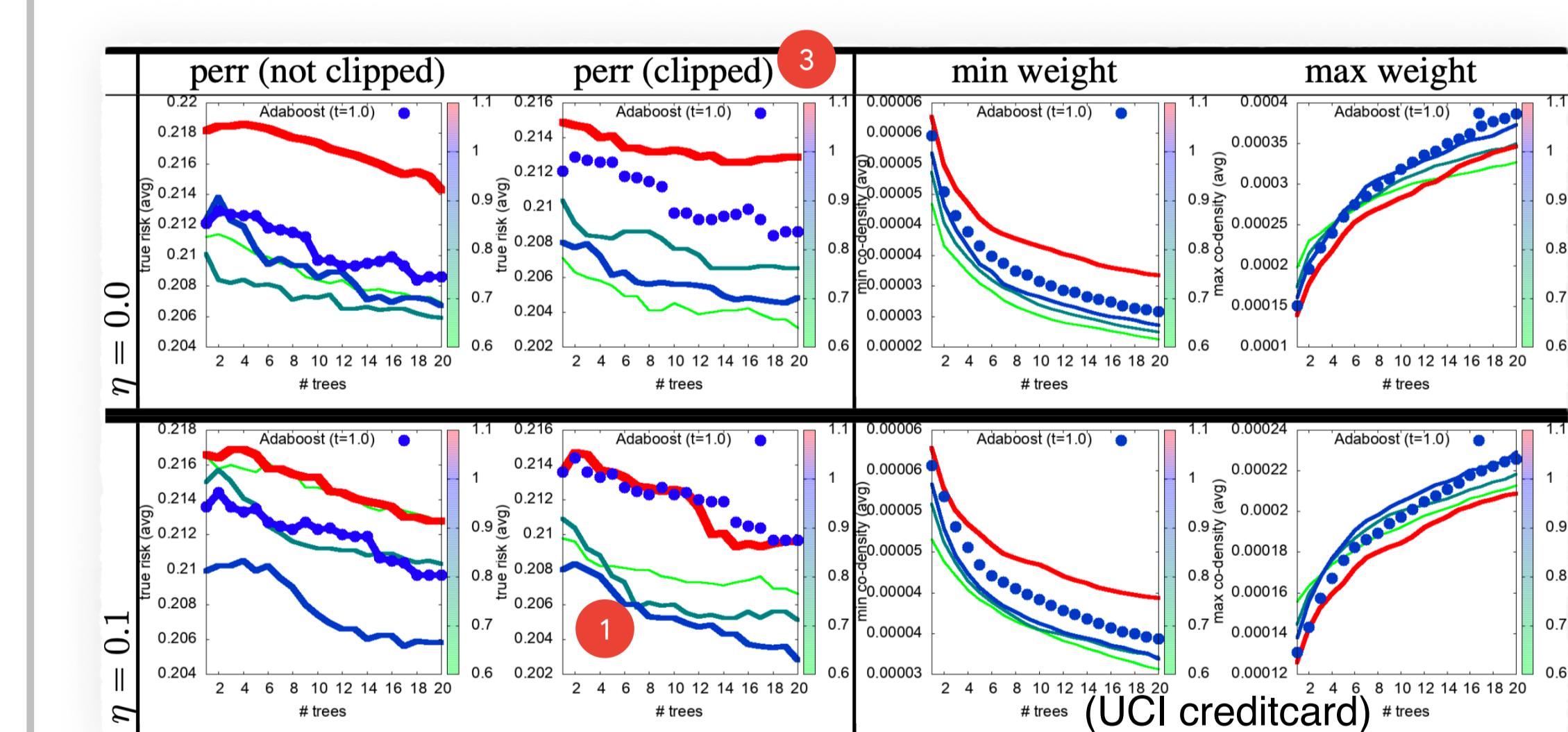


- Properties of tempered loss (DT induction), summarized

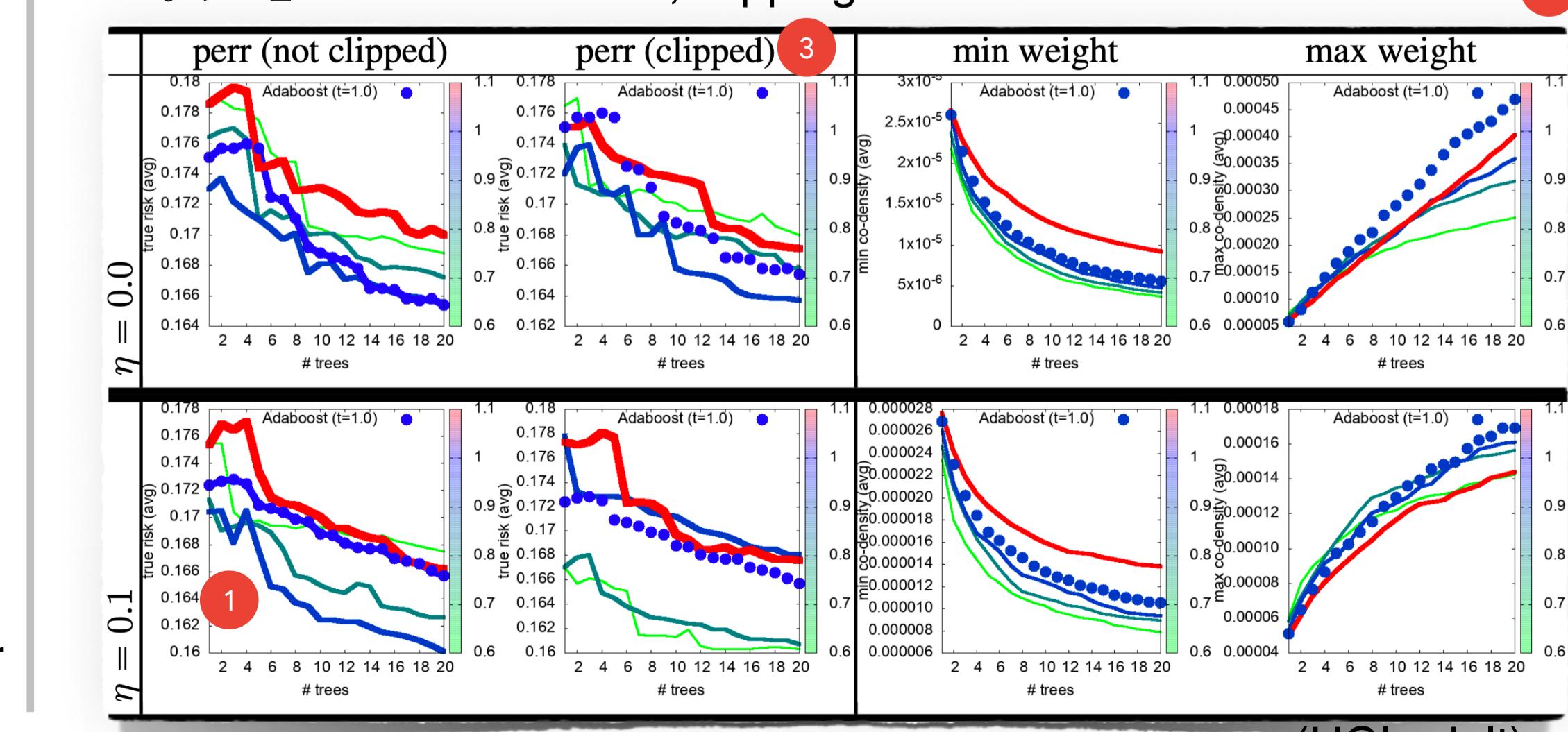
- For any $t \in (-\infty, 2)$, the tempered loss is symmetric, differentiable and **strictly proper** (Bayes rule optimal)
- For $t = 2$, "just" symmetric and proper
- Spans the full spectrum of known boosting rates for $t \in [-\infty, 1]$, near-optimal for $t = 1$... what about $t \in (1, 2)$?

Experiments

- 10-folds stratified CV, different t s, t not in $[0, 1]$ and symmetric label noise $\eta \in \{0, 0.1\}$



- $t \leq 1$ offers range of improvement compared to just AdaBoost
- +noise can completely reverse the picture
- $t > 1$ can be beneficial, clipping does work



- Conclusion:
- extension to more sophisticated boosting
 - some properties seem to depend on t :
 - consistency of [early l no] stopping
 - noise handling as a function of t
 - tuning of t at training time
 - code available at: <http://users.cecs.anu.edu.au/~rnock/>