

Summary

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

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

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

Tempered Exponential Measures (TEMs)

$+t$

- Math: **exp. fam. (EF)**, dual of cvx opt. pb
- Algo: training of **linear** & **tree-shaped** models
- Loss: convex surrogates, proper duals
- Comp: fast++ convergence (errors, margins)

exponential families

TEMs 101

Amid, Nock & Warmuth, 2023

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

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

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

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

$$q_t(\mathbf{x}) \propto \exp_t(\boldsymbol{\mu}^\top \Upsilon(\mathbf{x})) \quad \int q_t^{2-t} d\xi = 1 \quad \rightarrow \text{Normalized on co-simplex } \tilde{\Delta}_m$$

\rightarrow co-density $p_t \doteq q_t^{2-t}$

The t -Boosting TEM

Boosting's EF: Kivinen & Warmuth, 1999

- We seek the following TEM (update)

$$q' \doteq \arg \min_{\tilde{q} \in \tilde{\Delta}_m, \tilde{q}^\top \mathbf{u} = 0} D_t(\tilde{q} \| q)$$

$D_t(q \| q') \doteq \sum_{i \in [m]} q_i \cdot (\log_t q_i - \log_t q'_i) - \log_{t-1} q_i + \log_{t-1} q'_i$

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 $t \in \mathbb{R}_{\geq 0} \setminus \{2\}$,

- \rightarrow solutions have the form $q'_i = \frac{q_i \otimes_t \exp_t(-\mu u_i)}{Z_t}$
- unknown μ satisfies $\mu = \arg \min Z_t(\mu)$
- $\rightarrow Z_t(\mu')$ is strictly convex* \rightarrow finding solution to 1 via 2 easy using 3
- (*) +light assumption if $t = 0$

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

Fitting the model part

- Models learned

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

linear model clipped linear model decision tree

- Clipped summation

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

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 $1 = 2 - 1$

(non commutative, "encoding-nice")

- 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 $\mathbf{H}_J, \mathbf{H}_J^{(1/1-t)}$;

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:

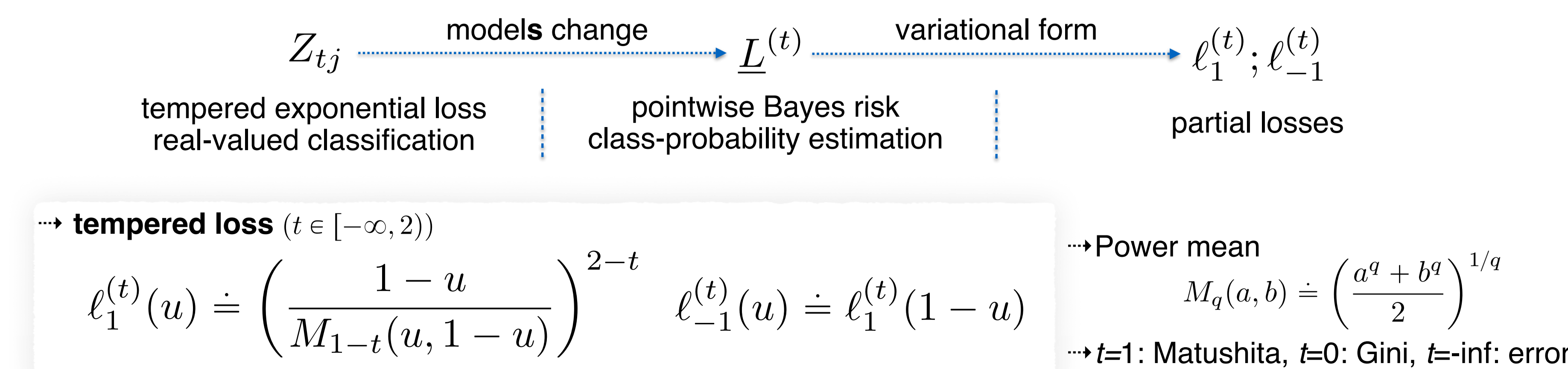
$$q_{(j+1)i} = \frac{q_{ji} \otimes_t \exp_t(-\mu_j u_{ji})}{Z_{tj}}, \quad \text{where } Z_{tj} = \|\mathbf{q}_j \otimes_t \exp_t(-\mu_j \cdot \mathbf{u}_j)\|_{2-t}$$

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):



Properties

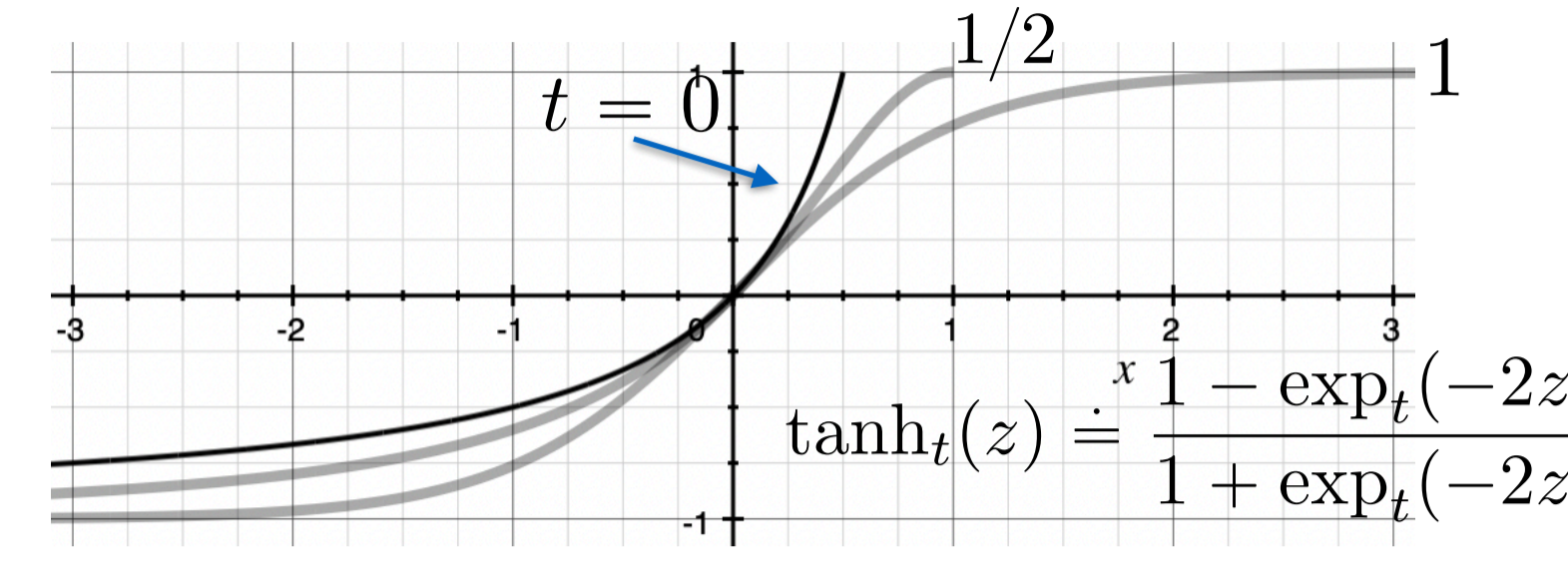
- (t)-margins

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

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



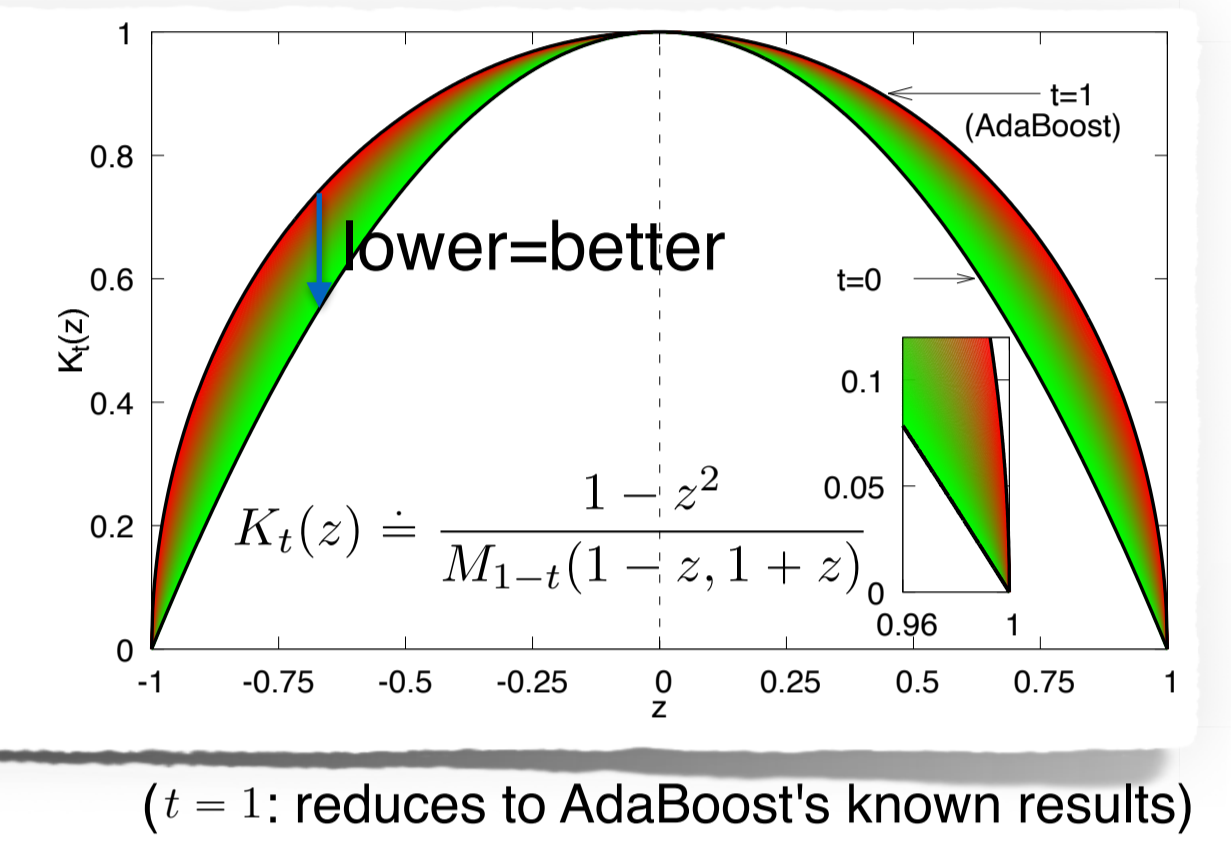
- Theorem (fast convergence for margins of t -AdaBoost, *simplified*), for $t \in [0, 1]$

\rightarrow In t -AdaBoost, fix

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

\rightarrow If there is no "forgetting weights" (e.g. $\mathbf{H}_J, \mathbf{H}_J^{(1/1-t)}$ 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

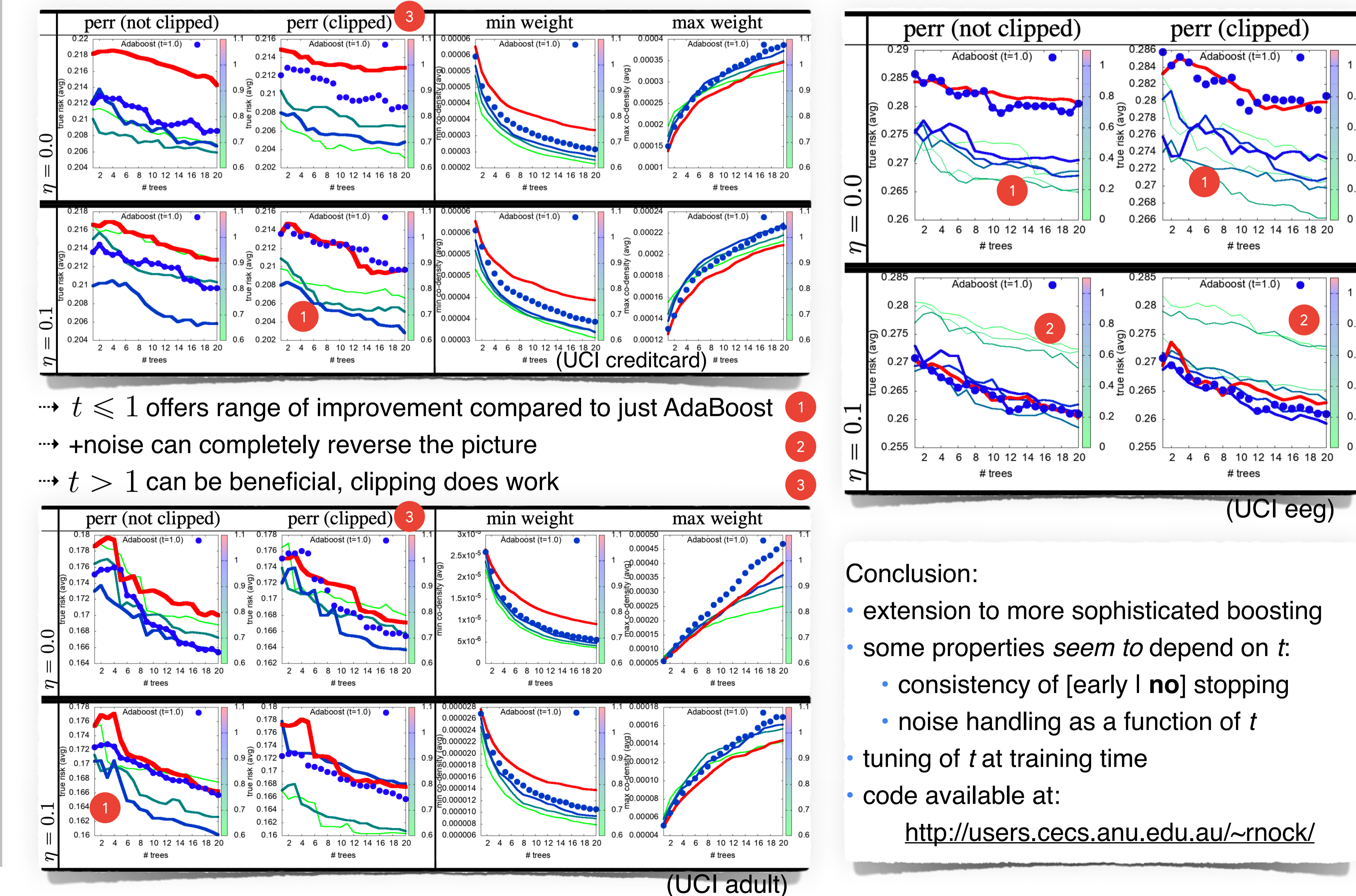
\rightarrow For any $t \in (-\infty, 2)$, the tempered loss is symmetric, differentiable and **strictly proper** (Bayes rule optimal)

\rightarrow For $t = 2$, "just" symmetric and proper

\rightarrow Spans the full spectrum of know boosting rates for $t \in [-\infty, 1]$, near-optimal for $t = 1 \dots$ what about $t \in (1, 2)$?

Experiments

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



Conclusion:

- extension to more sophisticated boosting
- some properties *seem* to depend on t :
 - consistency of [early | 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/>