# What loss do you use to train your classifier?*

*Yes, an LLM qualifies as a classifier

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# What is AdaBoost? 🚀

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# Weak learner 💪

$$\mathbf{x}^{(j)} \geq t \; ?$$

True                      False

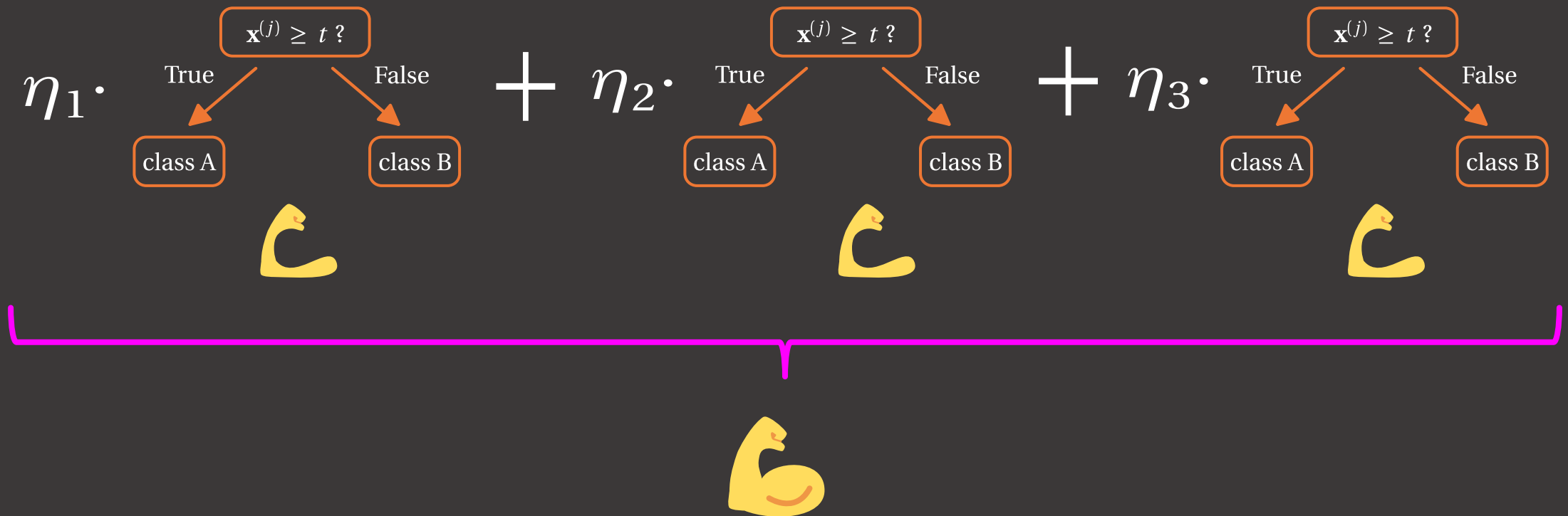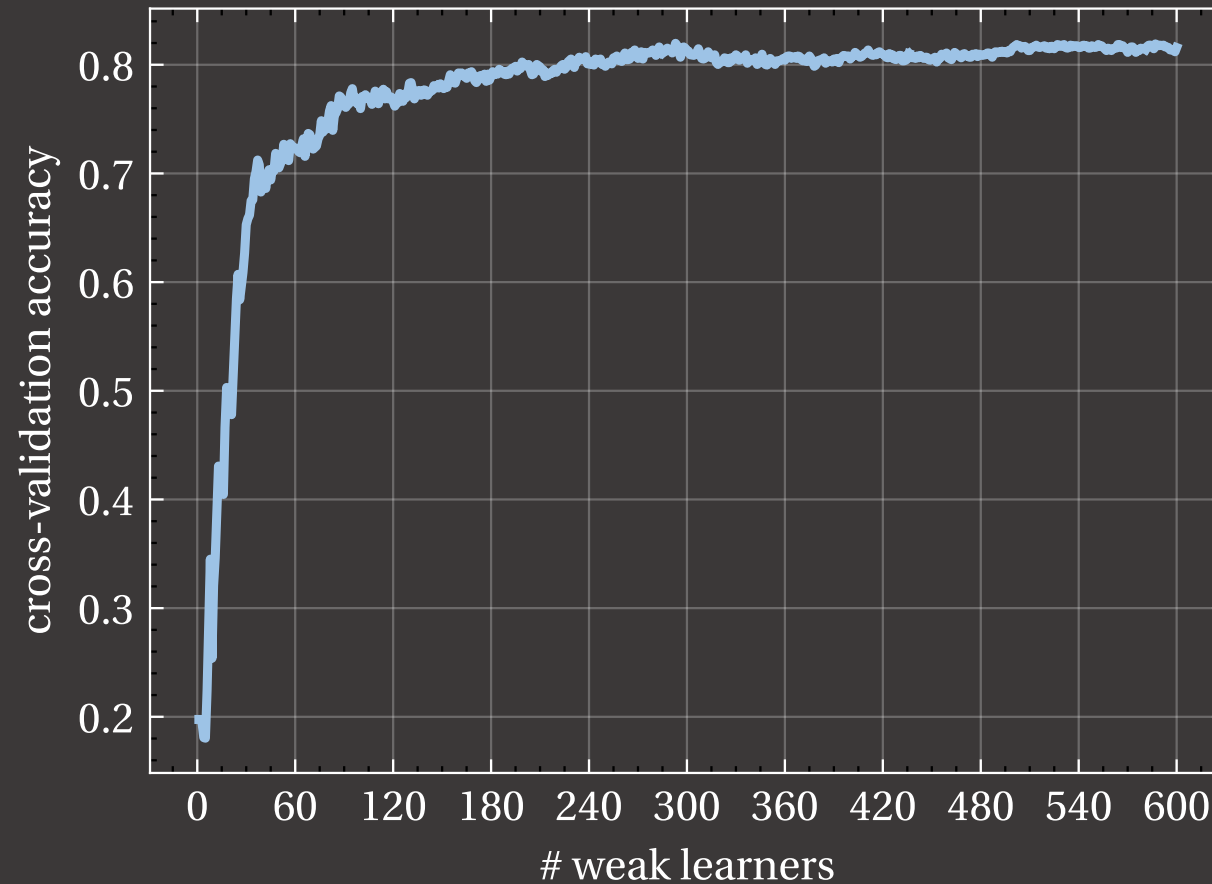class A                class B

# Constructing a strong learner

# Optimization objective: exponential loss

$$f(\mathbf{x}) = \sum_{i=1}^{N} \eta_i \, g_i(\mathbf{x})$$

$g_i$  greedily minimize

$$\hat{\mathbb{E}}\big[\exp\{-y\,f(\mathbf{x})\}\big], \qquad y \in \{-1, 1\}$$

Breiman, Leo. "Prediction Games and Arcing Algorithms." *Neural Computation* 11.7 (1999): 1493-1517.

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# AdaBoost is resilient to "overfitting"

**MAX PLANCK INSTITUTE**
FOR INTELLIGENT SYSTEMS

*"[AdaBoost with trees is] the best off-the-shelf classifier in the world."*

\- Leo Breiman

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. "Additive logistic regression: a statistical view of boosting. *The Annals of Statistics* 28.2 (2000): 337-407.

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# Can we translate the "AdaBoost magic" to neural networks?

# Multiclass exponential loss

$$\underbrace{\hat{\mathbb{E}}\left[\exp\left\{-(K-1)^{-1}f^{(y)}(\mathbf{x})\right\}\right]}_{\text{exponential loss}} \quad \text{subject to} \quad \underbrace{\sum_{j=1}^{K}f^{(j)}(\mathbf{x})=0,\ \forall\mathbf{x}.}_{\substack{\text{prevents logits from} \\ \text{diverging}}}$$

Constraints 😧

Zhu, Ji, et al. "Multi-class adaboost." *Statistics and its Interface* 2.3 (2009): 349-360.

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# Penalized exponential loss (PENEX)

$$\hat{\mathbb{E}}\left[ \exp\left\{ -\alpha f^{(y)}(\mathbf{x}) \right\} + \rho \sum_{j=1}^{K} \exp\left\{ f^{(j)}(\mathbf{x}) \right\} \right]$$

exponential loss

prevents logits from diverging

No constraints! 🥳

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# ... and it works 🔥



training loss:

— PENEX

— label smoothing

— cross-entropy

(CIFAR-100)

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# PENEX often works better than other methods

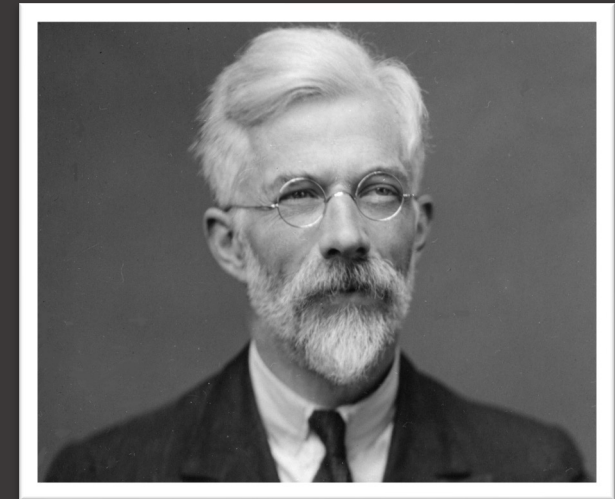| Method | Metric | CIFAR-10 | Noisy CIFAR-10 | CIFAR-100 | PathMNIST | BBC News |
|---|---|---|---|---|---|---|
| CE | ACC | $0.785 \pm 0.004$ | $0.724 \pm 0.004$ | $0.443 \pm 0.004$ | $0.826 \pm 0.004$ | $0.967 \pm 0.007$ |
| | -ECE | $-0.162 \pm 0.003$ | $-0.179 \pm 0.003$ | $-0.287 \pm 0.003$ | $-0.151 \pm 0.004$ | $\mathbf{-0.032} \pm 0.006$ |
| | -CE | $-1.004 \pm 0.024$ | $-1.125 \pm 0.019$ | $-3.072 \pm 0.034$ | $-2.018 \pm 0.130$ | $-0.109 \pm 0.024$ |
| | -BRIER | $-0.346 \pm 0.006$ | $-0.424 \pm 0.006$ | $-0.794 \pm 0.006$ | $-0.300 \pm 0.007$ | $-0.051 \pm 0.011$ |
| label smoothing | ACC | $0.789 \pm 0.004$ | $0.747 \pm 0.004$ | $0.451 \pm 0.005$ | $0.829 \pm 0.004$ | $0.970 \pm 0.006$ |
| | -ECE | $-0.112 \pm 0.002$ | $-0.183 \pm 0.003$ | $\mathbf{-0.147} \pm 0.002$ | $-0.109 \pm 0.002$ | $-0.033 \pm 0.006$ |
| | -CE | $-0.657 \pm 0.011$ | $-0.889 \pm 0.008$ | $-2.292 \pm 0.019$ | $\mathbf{-0.589} \pm 0.012$ | $-0.115 \pm 0.022$ |
| | -BRIER | $-0.300 \pm 0.005$ | $-0.384 \pm 0.004$ | $-0.692 \pm 0.004$ | $-0.255 \pm 0.005$ | $-0.049 \pm 0.010$ |
| confidence penalty | ACC | $0.786 \pm 0.004$ | $0.733 \pm 0.004$ | $0.449 \pm 0.006$ | $0.828 \pm 0.004$ | $\mathbf{0.974} \pm 0.006$ |
| | -ECE | $-0.130 \pm 0.002$ | $-0.149 \pm 0.003$ | $-0.152 \pm 0.002$ | $-0.110 \pm 0.003$ | $-0.050 \pm 0.005$ |
| | -CE | $-0.731 \pm 0.015$ | $-0.866 \pm 0.009$ | $-2.254 \pm 0.018$ | $-0.917 \pm 0.047$ | $-0.094 \pm 0.015$ |
| | -BRIER | $-0.317 \pm 0.005$ | $-0.385 \pm 0.004$ | $-0.695 \pm 0.005$ | $-0.262 \pm 0.005$ | $\mathbf{-0.042} \pm 0.008$ |
| focal loss | ACC | $0.778 \pm 0.004$ | $0.708 \pm 0.004$ | $0.428 \pm 0.005$ | $0.803 \pm 0.004$ | $0.970 \pm 0.006$ |
| | -ECE | $-0.117 \pm 0.002$ | $-0.165 \pm 0.003$ | $-0.161 \pm 0.003$ | $-0.112 \pm 0.003$ | $-0.051 \pm 0.005$ |
| | -CE | $-0.661 \pm 0.010$ | $-0.905 \pm 0.008$ | $-2.341 \pm 0.022$ | $-0.939 \pm 0.050$ | $\mathbf{-0.092} \pm 0.014$ |
| | -BRIER | $-0.313 \pm 0.005$ | $-0.423 \pm 0.004$ | $-0.723 \pm 0.005$ | $-0.291 \pm 0.006$ | $\mathbf{-0.042} \pm 0.008$ |
| PENEX | ACC | $\mathbf{0.793} \pm 0.004$ | $\mathbf{0.766} \pm 0.004$ | $\mathbf{0.460} \pm 0.005$ | $\mathbf{0.833} \pm 0.004$ | $0.968 \pm 0.006$ |
| | -ECE | $\mathbf{-0.109} \pm 0.002$ | $\mathbf{-0.131} \pm 0.002$ | $\mathbf{-0.147} \pm 0.003$ | $\mathbf{-0.100} \pm 0.003$ | $-0.034 \pm 0.006$ |
| | -CE | $\mathbf{-0.646} \pm 0.012$ | $\mathbf{-0.716} \pm 0.009$ | $\mathbf{-2.140} \pm 0.018$ | $-1.200 \pm 0.089$ | $-0.124 \pm 0.025$ |
| | -BRIER | $\mathbf{-0.299} \pm 0.005$ | $\mathbf{-0.332} \pm 0.004$ | $\mathbf{-0.685} \pm 0.004$ | $\mathbf{-0.251} \pm 0.006$ | $-0.055 \pm 0.011$ |

# Theoretical properties of PENEX

# Fisher consistency

*"When applied to the whole population the derived statistic should be equal to the parameter."*

- Ronald A. Fisher

**MAX PLANCK INSTITUTE**
FOR INTELLIGENT SYSTEMS

# PENEX is Fisher consistent

$$\hat{\mathbb{E}}\left[\exp\left\{-\alpha f^{(y)}(\mathbf{x})\right\} + \rho \sum_{j=1}^{K} \exp\left\{f^{(j)}(\mathbf{x})\right\}\right]$$

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# PENEX is Fisher consistent

$$\mathbb{E}\left[\exp\left\{-\alpha f^{(y)}(\mathbf{x})\right\} + \rho \sum_{j=1}^{K} \exp\left\{f^{(j)}(\mathbf{x})\right\}\right]$$

Minimize w.r.t. $f \to f_*$

$$P(y \mid \mathbf{x}) \propto \exp\left\{(1+\alpha)f_*^{(y)}(\mathbf{x})\right\}, \quad \forall \mathbf{x}$$

# Common regularizers fail Fisher consistency

$$\mathscr{L}_{\mathrm{CE}}(f) + \lambda \Omega(f)$$

Encompasses label smoothing, L2 regularization, confidence penalty, ...

Intuition: Regularization term $\Omega(f)$ pushes the solution off the Bayes-optimal predictor

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

$$\mathscr{L}_{\text{PENEX}}(f) \qquad\qquad \mathscr{L}_{\text{CE}}(f) + \lambda\Omega(f)$$


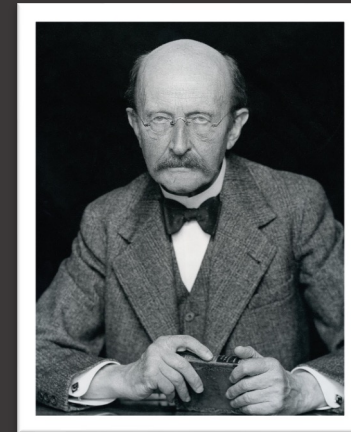
(nice car with five seats)



(nice car with two seats
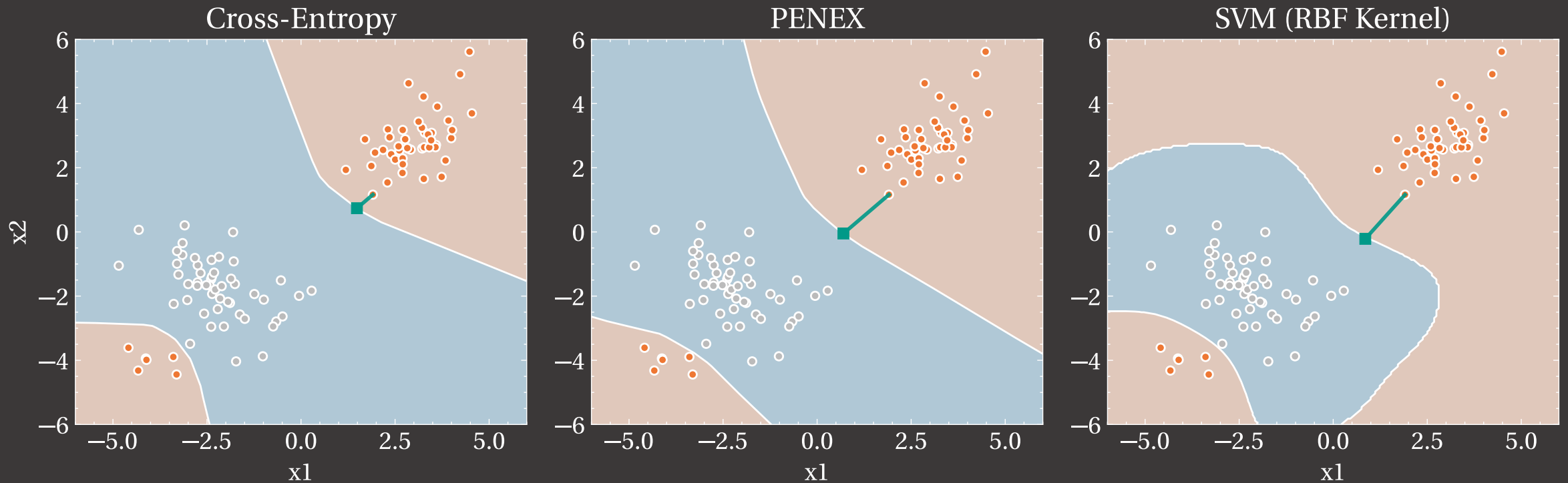and two extra seats
mounted on top)

# So why does it work?

*"Insight must precede application."*

\- Max Planck

# Implicit margin maximization

# PENEX provably maximizes margins

Defining the margin for example $(\mathbf{x}, y)$ as

$$m_f(\mathbf{x}, y) := f^{(y)}(\mathbf{x}) - \max_{j \neq y} f^{(j)}(\mathbf{x}),$$

we show that

$$\mathbb{P}(m_f(\mathbf{x}, y) \leq \gamma) \leq e^{\gamma \frac{\alpha}{\alpha+1}} \rho^{-\frac{\alpha}{\alpha+1}} \mathbb{E}[\mathscr{L}_{\text{PENEX}}(f; \alpha, \rho)].$$

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

# Key take-aways

- The "AdaBoost magic" can be translated to NNs

- Regularization is not at odds with Fisher consistency

- PENEX implicitly maximizes margins

- Let's question the very foundations!

Preprint:

Code:

E-Mail: kkladny@tuebingen.mpg.de

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS