

# **H-Consistency Bounds** for Surrogate Loss Minimizers

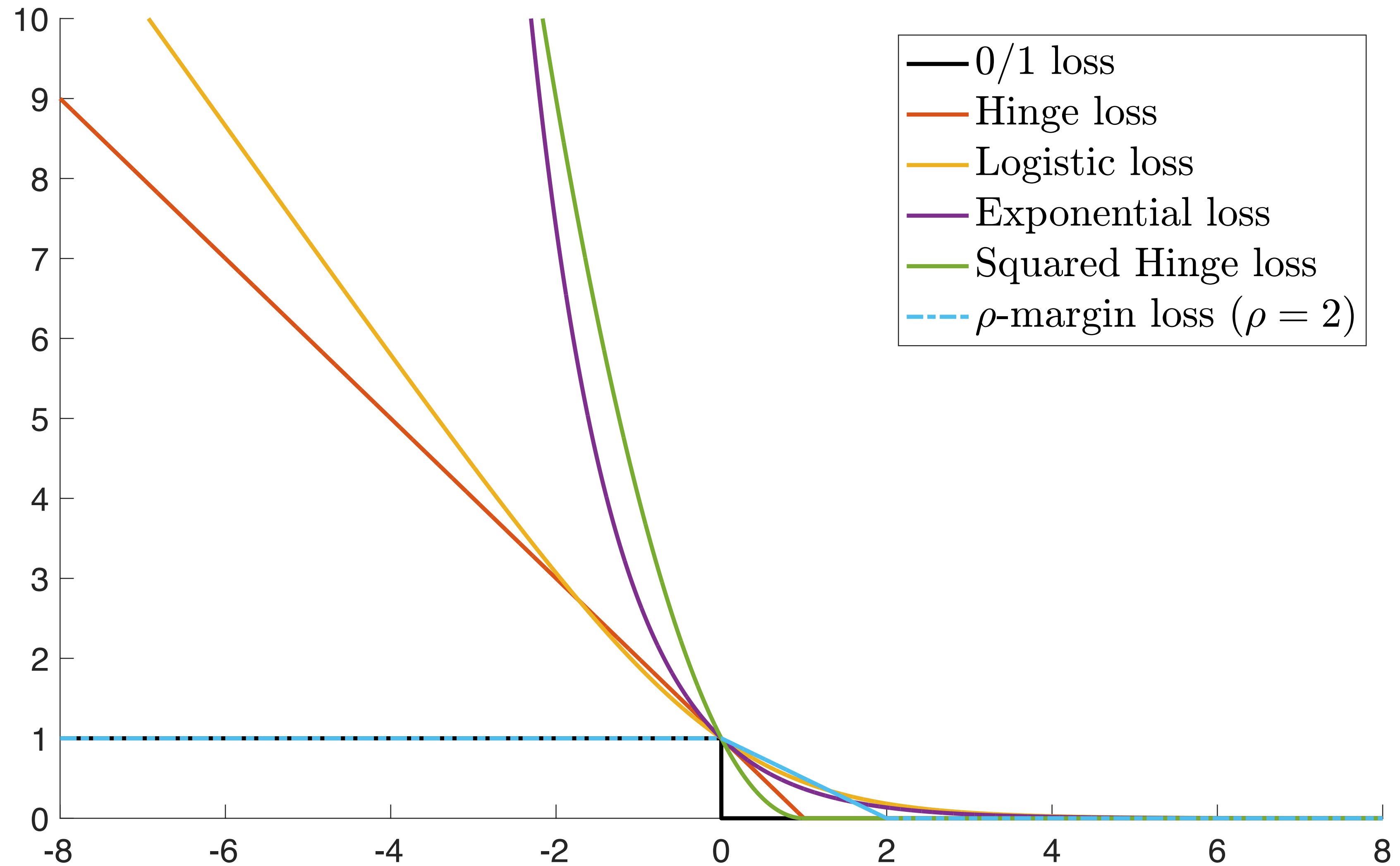
Pranjal Awasthi   Anqi Mao   Mehryar Mohri   Yutao Zhong

Google Research and Courant Institute

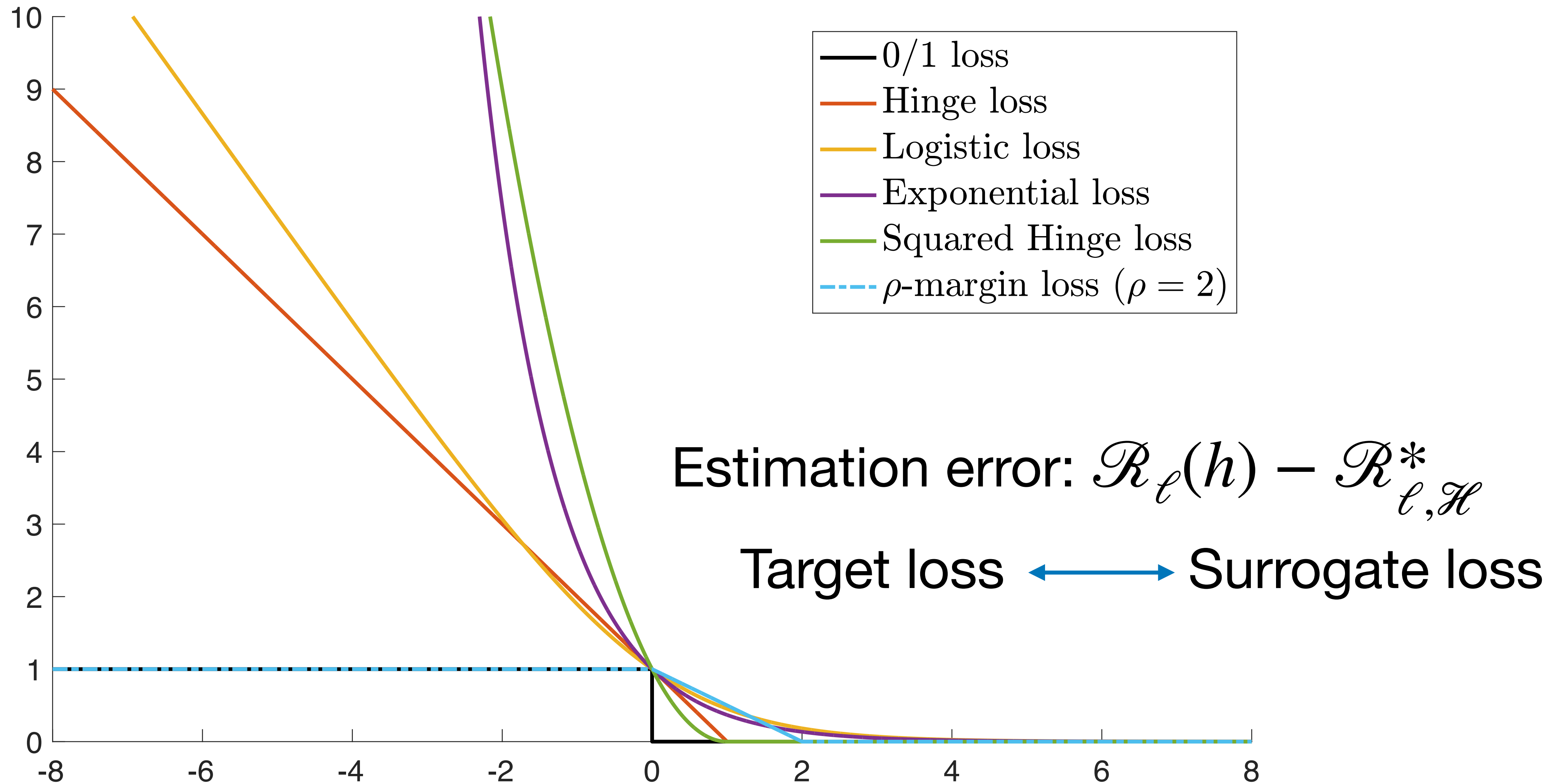


**ICML**  
International Conference  
On Machine Learning

# Surrogate Loss versus Target Loss



# Surrogate Loss versus Target Loss



## Bayes-consistency

$$\lim_{n \rightarrow \infty} \mathcal{R}_{\ell_1}(h_n) - \mathcal{R}_{\ell_1, \mathcal{H}_{\text{all}}}^* = 0 \Rightarrow \lim_{n \rightarrow \infty} \mathcal{R}_{\ell_2}(h_n) - \mathcal{R}_{\ell_2, \mathcal{H}_{\text{all}}}^* = 0$$

## Bayes-consistency

$$\lim_{n \rightarrow \infty} \mathcal{R}_{\ell_1}(h_n) - \mathcal{R}_{\ell_1, \mathcal{H}_{\text{all}}}^* = 0 \Rightarrow \lim_{n \rightarrow \infty} \mathcal{R}_{\ell_2}(h_n) - \mathcal{R}_{\ell_2, \mathcal{H}_{\text{all}}}^* = 0$$

## H-consistency

$$\lim_{n \rightarrow \infty} \mathcal{R}_{\ell_1}(h_n) - \mathcal{R}_{\ell_1, \mathcal{H}}^* = 0 \Rightarrow \lim_{n \rightarrow \infty} \mathcal{R}_{\ell_2}(h_n) - \mathcal{R}_{\ell_2, \mathcal{H}}^* = 0$$

# $\mathcal{H}$ -Consistency Bound

$$\mathcal{R}_{\ell_2}(h) - \mathcal{R}_{\ell_2, \mathcal{H}}^* \leq f\left(\mathcal{R}_{\ell_1}(h) - \mathcal{R}_{\ell_1, \mathcal{H}}^*\right)$$

# $\mathcal{H}$ -Consistency Bound

$$\mathcal{R}_{\ell_2}(h) - \mathcal{R}_{\ell_2, \mathcal{H}}^* \leq f\left(\mathcal{R}_{\ell_1}(h) - \mathcal{R}_{\ell_1, \mathcal{H}}^*\right)$$

$$\forall h \in \mathcal{H}, \mathcal{D} \in \mathcal{P}$$

$\mathcal{P}_{\text{all}}$  — distribution independent

# Standard Binary Classification

Family of all measurable functions:  $\mathcal{H}_{\text{all}}$

All distributions:  $\mathcal{P}_{\text{all}}$

(Zhang, 04a; Bartlett et al., 06; Mohri et al., 18)

$\ell_1$ : margin-based loss  $\phi$   $\longleftrightarrow$  ?  $\longrightarrow$   $\ell_2$ : standard 0/1 loss  $\ell_{0-1}$

Excess error bound ( $\mathcal{H}_{\text{all}}$ -consistency bound)



# Standard Binary Classification

General hypothesis sets:  $\mathcal{H}$

Distribution-independent and distribution-dependent  
(Our contribution)

$\ell_1$ : margin-based loss  $\phi$   $\longleftrightarrow$  ?  $\longrightarrow$   $\ell_2$ : standard 0/1 loss  $\ell_{0-1}$

$\mathcal{H}$ -consistency bound



# Adversarial Attacks (Szegedy et al., 13)



Correct

Attack

Ostrich

Correct

Attack

Ostrich



# Adversarially Robust Classification

General hypothesis sets:  $\mathcal{H}$

Distribution-independent and distribution-dependent  
(Our contribution)

$\ell_1$ : supremum-based margin-based loss  $\tilde{\phi}$   $\longleftrightarrow$  ?  $\longrightarrow$   $\ell_2$ : adversarial 0/1 loss  $\ell_\gamma$

$$\tilde{\phi} = \sup_{x': \|x-x'\| \leq \gamma} \phi(yf(x'))$$

$$\ell_\gamma = \sup_{x': \|x-x'\| \leq \gamma} \ell_{0-1}(yf(x'))$$

Adversarial  $\mathcal{H}$ -consistency bound

# $\mathcal{H}$ -Consistency Bounds Analysis

$$\mathcal{R}_{\ell_2}(h) - \mathcal{R}_{\ell_2, \mathcal{H}}^* \leq f\left(\mathcal{R}_{\ell_1}(h) - \mathcal{R}_{\ell_1, \mathcal{H}}^*\right)$$



Inverse of  $\mathcal{H}$ -estimation error transformation + Minimizability gap

# Minimizability Gap

$$\mathcal{M}_{\ell, \mathcal{H}} = \mathcal{R}_{\ell, \mathcal{H}}^* - \mathbb{E}_X \left[ \mathcal{C}_{\ell, \mathcal{H}}^*(x) \right]$$

difference of the best-in class error and  
the expectation of the minimal conditional-risk

$$\mathcal{C}_{\ell, \mathcal{H}}^*(x) = \inf_{h \in \mathcal{H}} \left[ \mathcal{D}(Y = 1 | X = x) \ell(h, x, +1) + (1 - \mathcal{D}(Y = 1 | X = x)) \ell(h, x, -1) \right]$$

# Minimizability Gap

$$\mathcal{M}_{\ell, \mathcal{H}} = \mathcal{R}_{\ell, \mathcal{H}}^* - \mathbb{E}_X \left[ \mathcal{C}_{\ell, \mathcal{H}}^*(x) \right]$$

difference of the best-in class error and  
the expectation of the minimal conditional-risk

$$\mathcal{C}_{\ell, \mathcal{H}}^*(x) = \inf_{h \in \mathcal{H}} \left[ \mathcal{D}(Y = 1 | X = x) \ell(h, x, +1) + (1 - \mathcal{D}(Y = 1 | X = x)) \ell(h, x, -1) \right]$$

we cannot hope to estimate or minimize.

# $\mathcal{H}$ -Estimation Error Transformation

$$\forall t \in [0, 1], \mathcal{T}_{\Phi}(t) = \mathcal{T}(t) \mathbb{1}_{t \in [\epsilon, 1]} + (\mathcal{T}(\epsilon)/\epsilon) t \mathbb{1}_{t \in [0, \epsilon)}$$

$$\text{Where } \mathcal{T}(t) := \inf_{x \in \mathcal{X}, h \in \mathcal{H} : h(x) < 0} \Delta \mathcal{C}_{\Phi, \mathcal{H}} \left( h, x, \frac{t+1}{2} \right)$$

# $\mathcal{H}$ -Estimation Error Transformation

$$\forall t \in [0, 1], \mathcal{T}_{\Phi}(t) = \mathcal{T}(t) \mathbb{1}_{t \in [\epsilon, 1]} + (\mathcal{T}(\epsilon)/\epsilon) t \mathbb{1}_{t \in [0, \epsilon)}$$

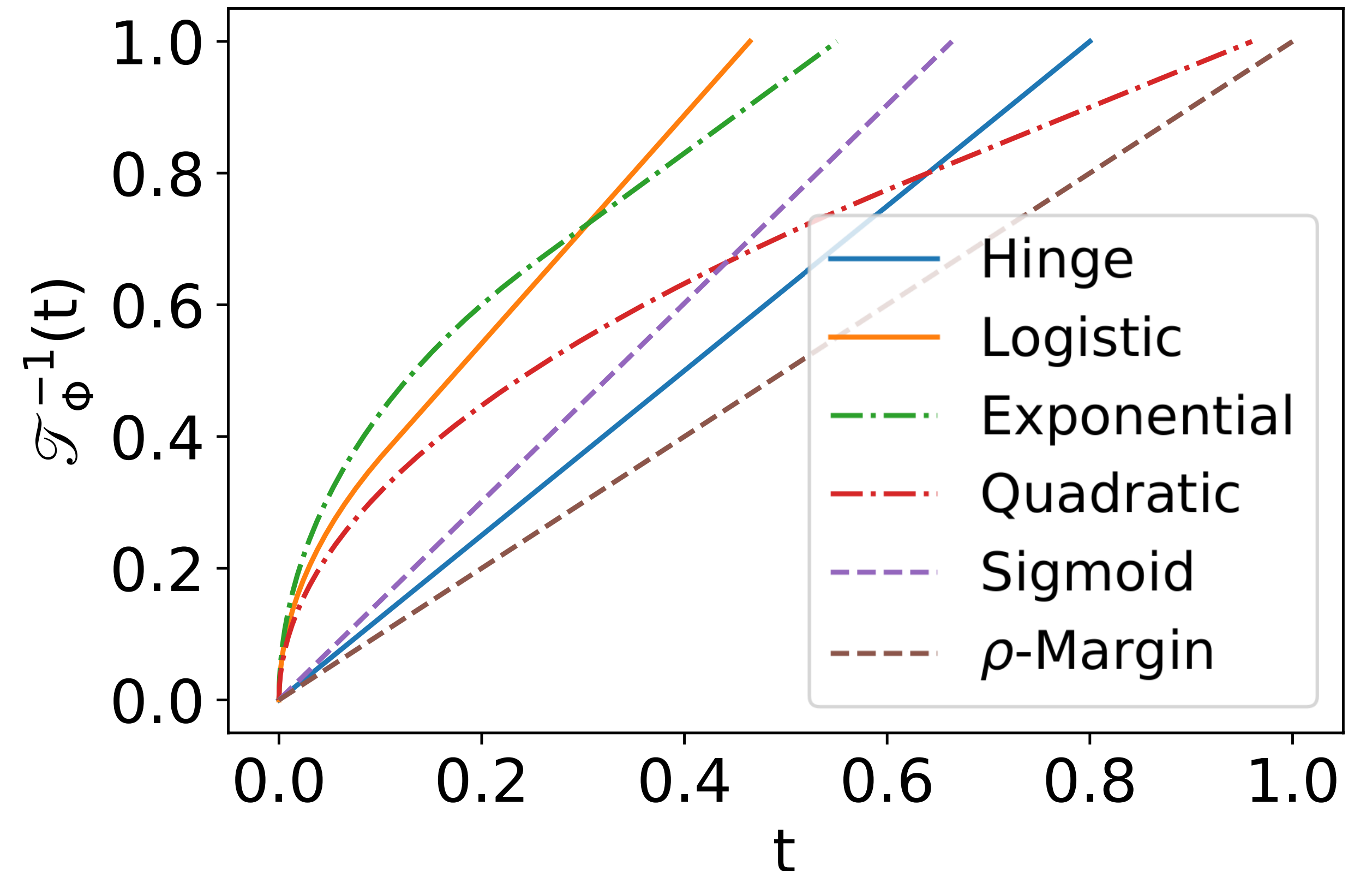
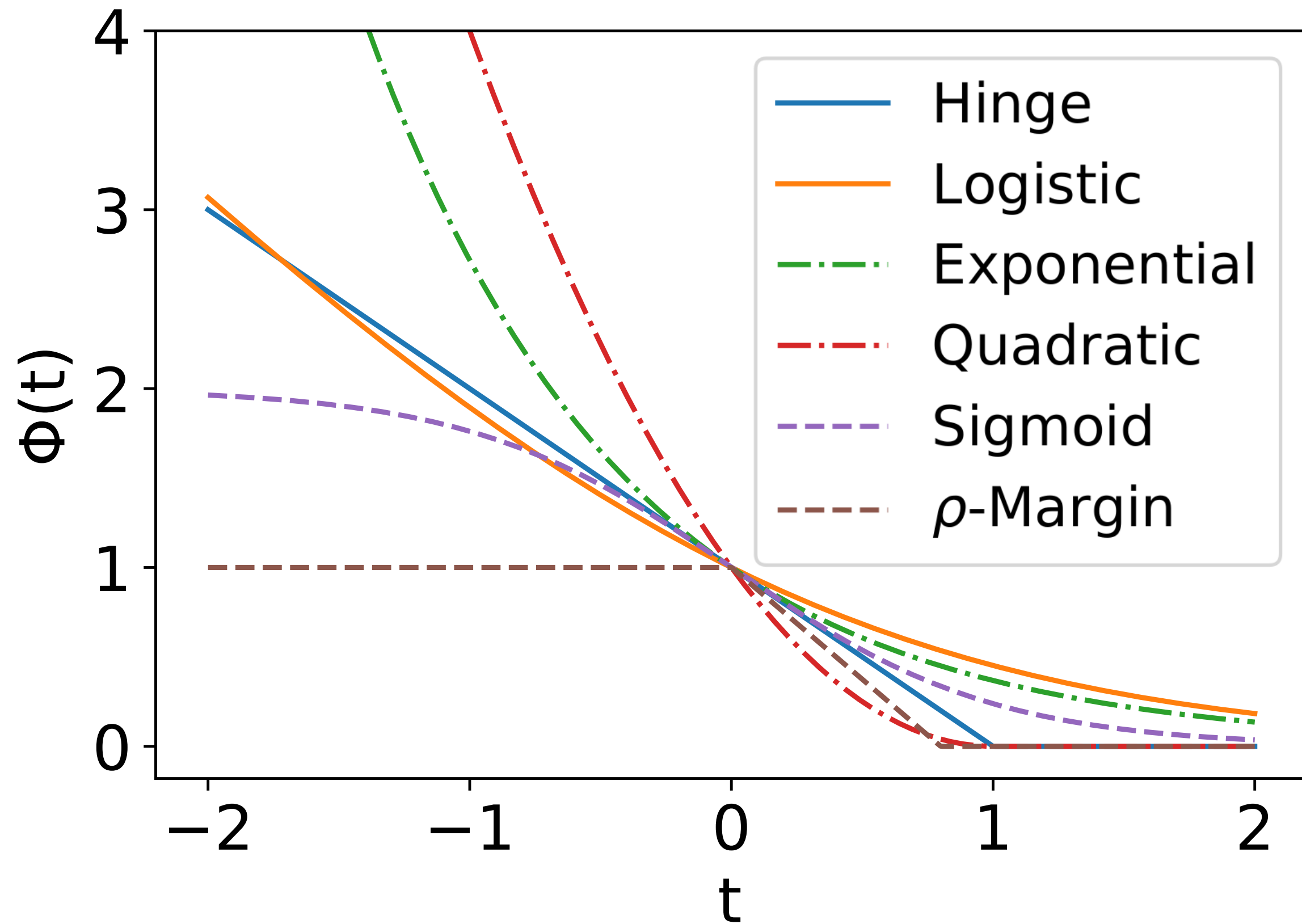
Where  $\mathcal{T}(t) := \inf_{x \in \mathcal{X}, h \in \mathcal{H} : h(x) < 0} \Delta \mathcal{C}_{\Phi, \mathcal{H}} \left( h, x, \frac{t+1}{2} \right)$



Tightness!



# $\mathcal{H}$ -Estimation Error Transformation



# Conclusion

$\mathcal{H}$ -consistency bounds for both standard and adversarial binary classification

- New estimation error guarantees for both the non-adversarial 0/1 loss function and the adversarial 0/1 loss function
- Compare different surrogate loss functions of the 0/1 loss or adversarial loss, given the specific hypothesis set used
- Theoretical and conceptual tools helpful for the analysis of other loss functions and other hypothesis sets

# Conclusion

$\mathcal{H}$ -consistency bounds for both standard and adversarial binary classification

- New estimation error guarantees for both the non-adversarial 0/1 loss function and the adversarial 0/1 loss function
- Compare different surrogate loss functions of the 0/1 loss or adversarial loss, given the specific hypothesis set used
- Theoretical and conceptual tools helpful for the analysis of other loss functions and other hypothesis sets

# Future Directions

- $\mathcal{H}$ -consistency bounds for other loss functions and other hypothesis sets
- Incorporating the trade-off of the optimization and  $\mathcal{H}$ -consistency bounds

Poster #1206 in Hall E