

Calibrated Surrogate Losses and Robust Learning

Dec. 14, 2020 (in US) & Dec. 15, 2020 (in Japan)

Han Bao



東京大学
THE UNIVERSITY OF TOKYO



The content is based on our joint work with Clayton Scott (UMich) and Masashi Sugiyama (RIKEN / UTokyo).

Outline

- **Part 1: Calibrated surrogate losses**
 - ❖ Q. What are minimum requirements for loss functions?
- **Part 2: Loss functions in robust learning**
 - ❖ Q. Is it possible to design robust loss functions?

Setting: Binary Classification

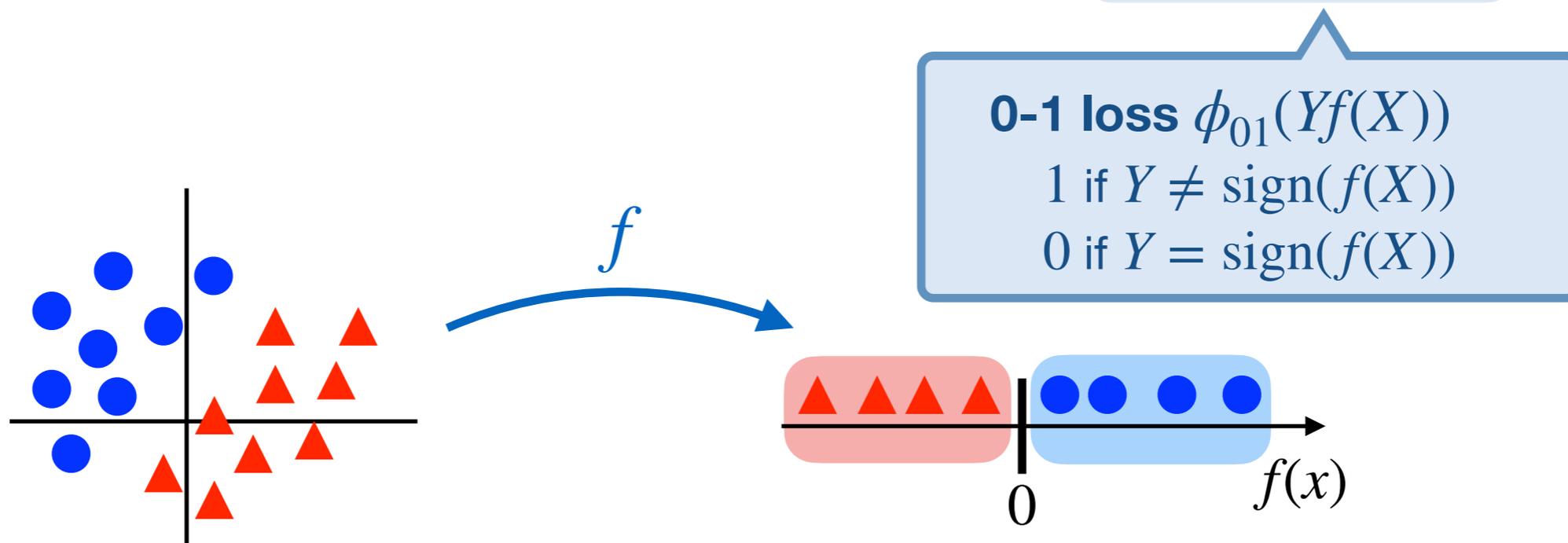
- Input

- ❖ sample $\{(x_i, y_i)\}_{i=1}^n$: feature $x_i \in \mathcal{X}$ and label $y_i \in \{\pm 1\}$

- Output: classifier $f: \mathcal{X} \rightarrow \mathbb{R}$

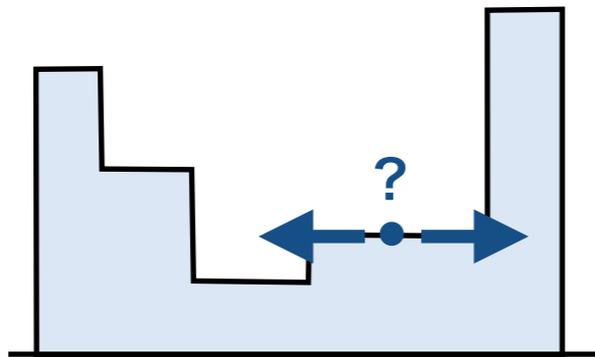
- ❖ use $\text{sign}(f(\cdot))$ to predict labels

- ❖ criterion: misclassification rate $R_{01}(f) = \mathbb{E} [\mathbf{1}[Y \neq \text{sign}(f(X))]]$

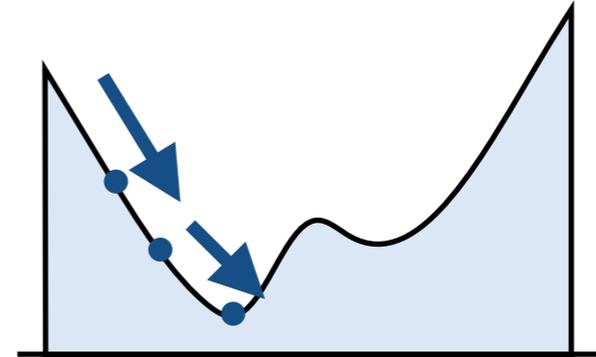


Surrogate Losses

- Motivation: minimizing 0-1 loss is NP-hard

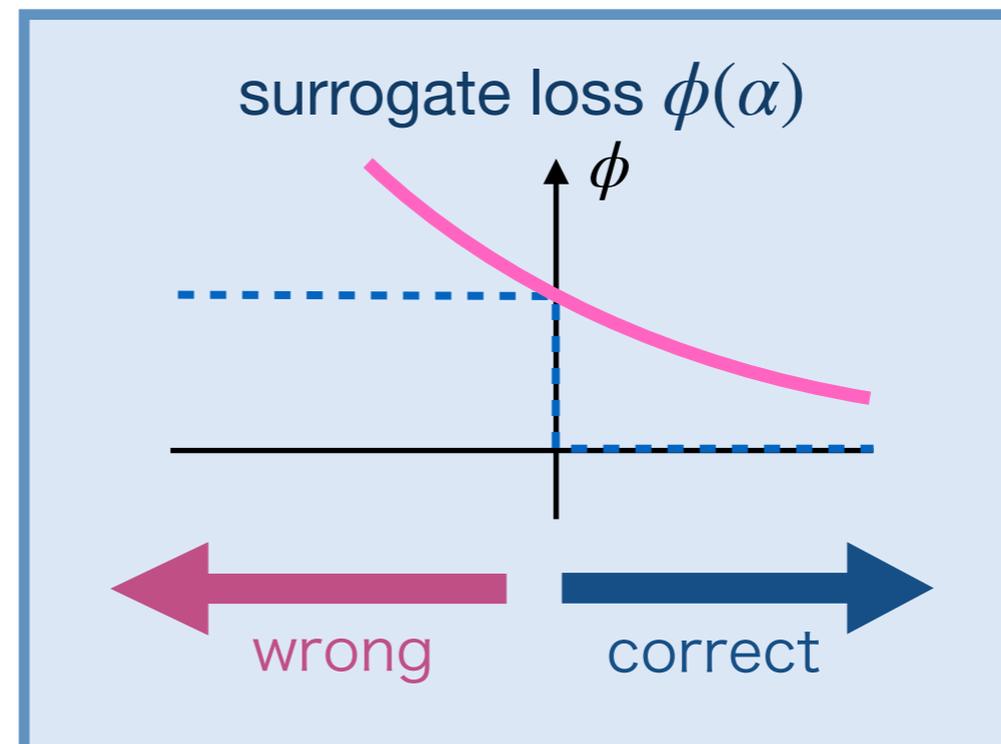
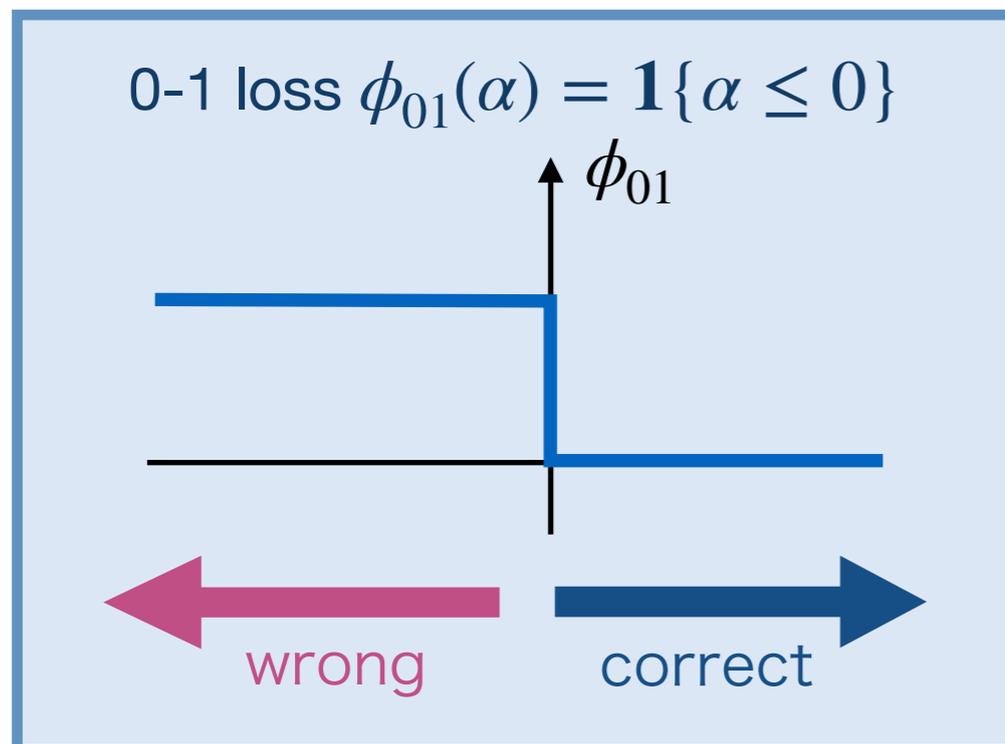


no gradient for discrete function



easily optimizable if convex and smooth

- Replace 0-1 loss with surrogate loss



hinge loss,
logistic loss, etc.

Elements of Learning Theory

(empirical) surrogate risk

$$\hat{R}_\phi(f) = \frac{1}{n} \sum_{i=1}^n \phi(y_i f(x_i))$$

(population) surrogate risk

$$R_\phi(f) = \mathbb{E}[\phi(Yf(X))]$$

target risk

$$R_{\phi_{01}}(f) = \mathbb{E}[\phi_{01}(Yf(X))]$$

Generalization theory:

If model is not too complicated,
then justified (roughly speaking)

Our interests: **Calibration theory**

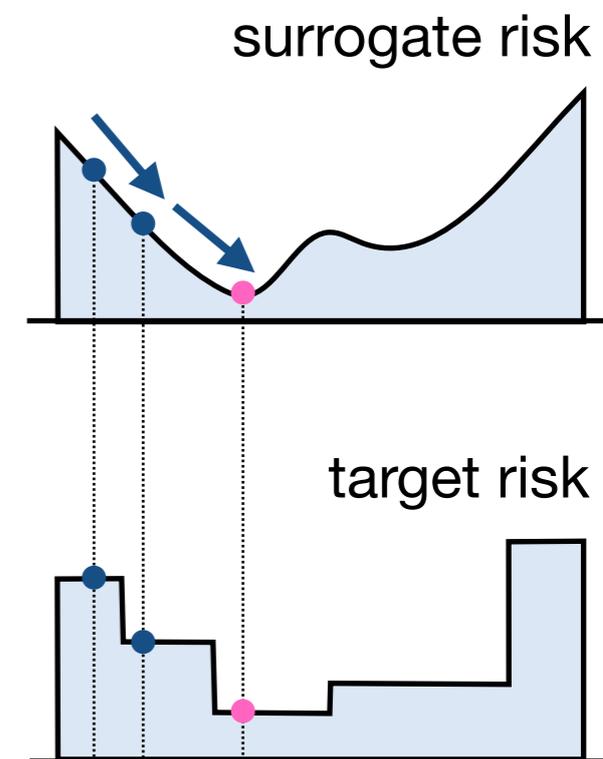
Q. What is a desirable surrogate?

- A. surrogate risk minimizer should be target risk minimizer

For two losses ψ (target) and ϕ (surrogate),

Definition. Surrogate ϕ is **calibrated** to target ψ if for any $\varepsilon > 0$, there exists $\delta > 0$ such that for all f ,

$$R_{\phi}(f) - R_{\phi}^* < \delta \implies R_{\psi}(f) - R_{\psi}^* < \varepsilon.$$



How to check calibration?

[Steinwart 2007]

Definition. Surrogate ϕ is calibrated to target ψ if for any $\varepsilon > 0$, there exists $\delta > 0$ such that for all f ,

$$R_{\phi}(f) - R_{\phi}^* < \delta \implies R_{\psi}(f) - R_{\psi}^* < \varepsilon.$$



$$R_{\psi}(f) - R_{\psi}^* \geq \varepsilon \implies R_{\phi}(f) - R_{\phi}^* \geq \delta$$

ε - δ definition of limit

contraposition

Definition. (calibration function)

$$\delta(\varepsilon) = \inf_f \left(R_{\phi}(f) - R_{\phi}^* \right) \text{ s.t. } \left(R_{\psi}(f) - R_{\psi}^* \geq \varepsilon \right)$$

smallest possible δ given lower bound of target

easy to ask existence of $\delta > 0$ given ε

Disclaimer: calibration function is defined over class-conditional risk to be precise

Main Tool: Calibration Function

[Steinwart 2007]

Definition. (calibration function)

$$\delta(\varepsilon) = \inf_f R_\phi(f) - R_\phi^* \quad \text{s.t.} \quad R_\psi(f) - R_\psi^* \geq \varepsilon$$

smallest possible surrogate given lower bound of target

- Provides iff condition

❖ calibrated to $\psi \iff \delta(\varepsilon) > 0$ for all $\varepsilon > 0$

- Provides excess risk bound

❖ calibrated to $\psi \iff R_\psi(f) - R_\psi^* \leq (\delta^{**})^{-1} (R_\phi(f) - R_\phi^*)$

target risk monotone surrogate risk

minimizing surrogate risk = minimizing target risk
(we know convergence rate in addition)

Case: Binary Classification

[Bartlett *et al.* 2006]

Calibration function for 0-1 loss

$$\delta(\varepsilon) = \inf_f R_\phi(f) - R_\phi^* \quad \text{s.t.} \quad R_{\phi_{01}}(f) - R_{\phi_{01}}^* \geq \varepsilon$$

smallest possible surrogate
given lower bound of 0-1

iff condition: $\delta(\varepsilon) > 0 \quad \forall \varepsilon > 0$

calibrated iff

$$\inf_f \left\{ R_\phi(f) \mid R_{\phi_{01}}(f) > R_{\phi_{01}}^* \right\} > \inf_f R_\phi(f)$$

f is non-optimal wrt 0-1 loss

minimum risk of non-optimal classifiers
 \vee
 minimum risk of all classifiers

- Check the latter condition to see if calibrated
- More simple equivalent conditions available (next slide)

Disclaimer: several literature defines calibration by the latter condition

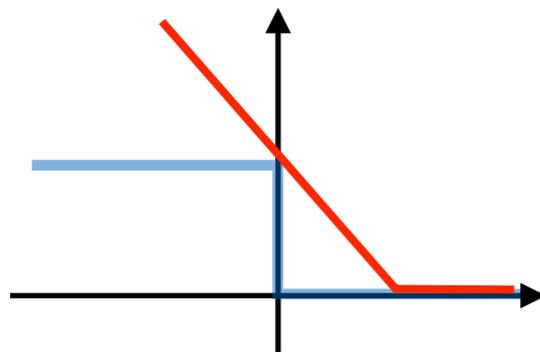
Case: Binary Classification

[Bartlett *et al.* 2006]

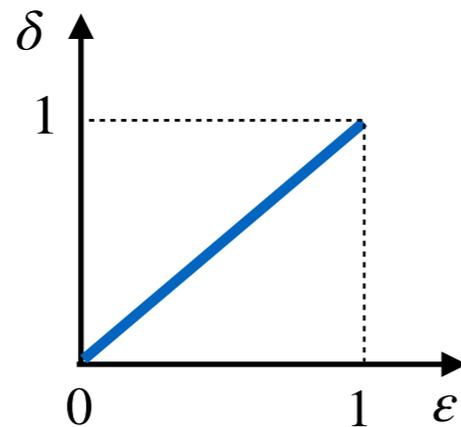
Theorem. If surrogate ϕ is convex, it is calibrated to ϕ_{01} iff

- differentiable at 0,
- $\phi'(0) < 0$.

hinge loss

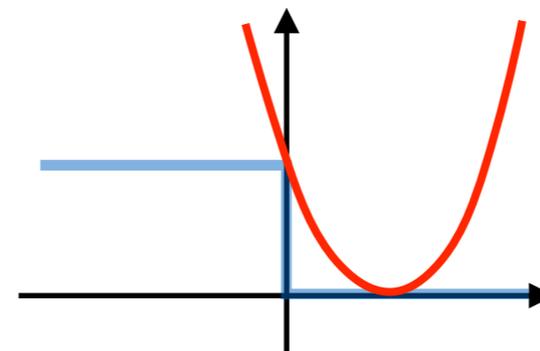


$$\phi(\alpha) = [1 - \alpha]_+$$

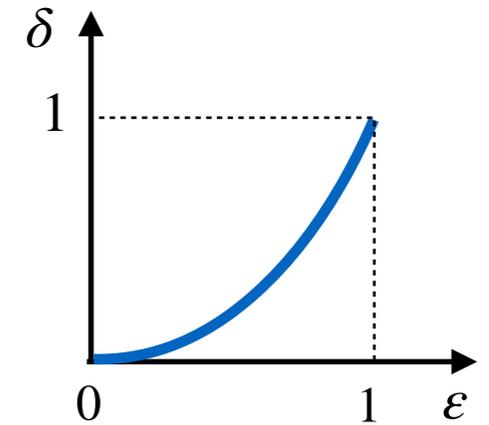


$$\delta(\epsilon) = \epsilon$$

squared loss



$$\phi(\alpha) = (1 - \alpha)^2$$



$$\delta(\epsilon) = \epsilon^2$$

- Most of well-known losses are calibrated

❖ perceptron loss $\phi(\alpha) = [-\alpha]_+$ is not

Outline

- Part 1: Calibrated surrogate losses

- ❖ Q. What are minimum requirements for loss functions?

A. calibration: surrogate minimizer = target minimizer

- confirmed via calibration function
- simple iff conditions for binary classification

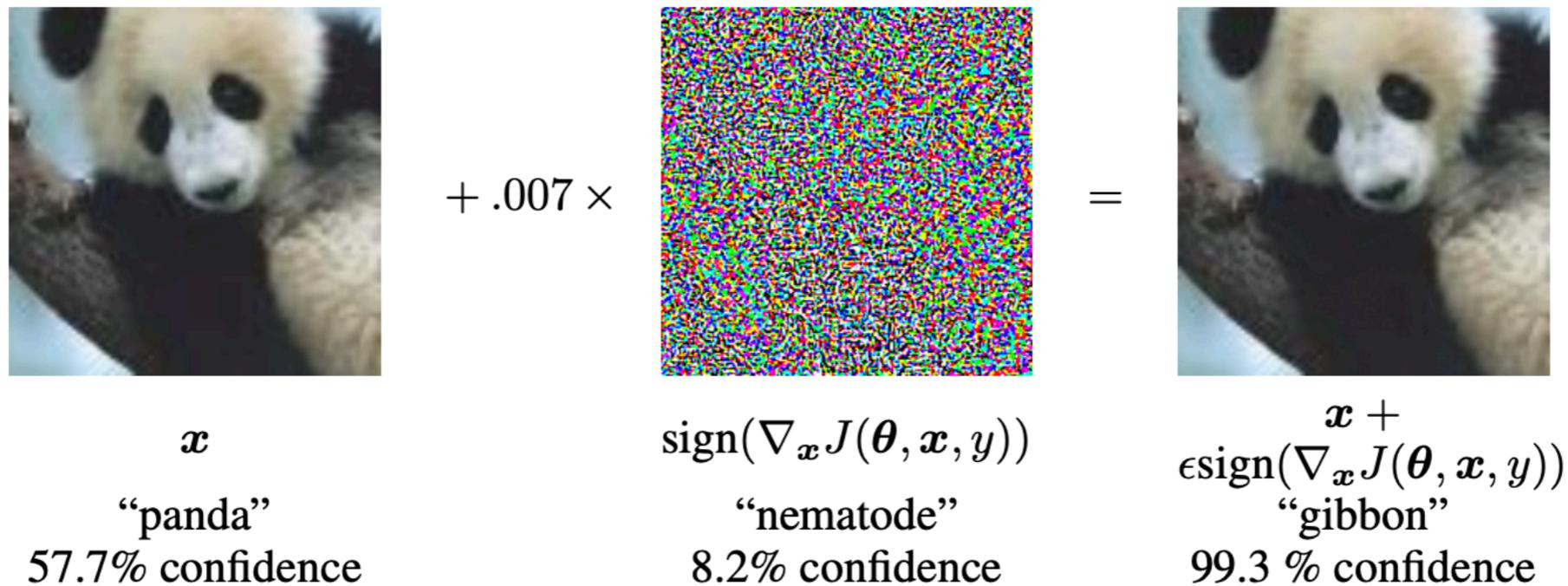
- Part 2: Loss functions in robust learning

- ❖ Q. Is it possible to design robust loss functions?

Classifier is vulnerable to “attacks” ¹²

[Goodfellow *et al.* 2015]

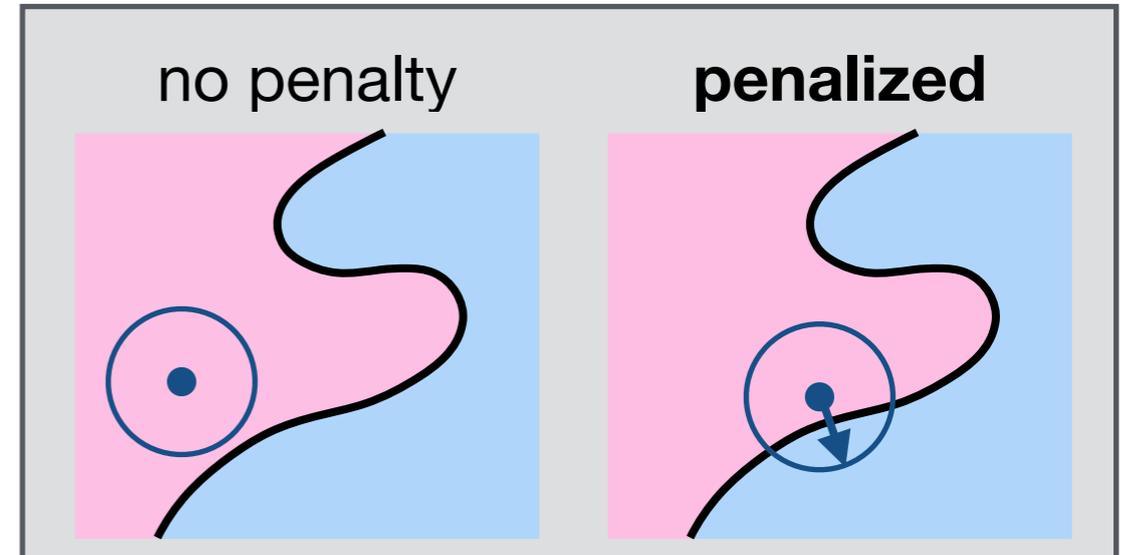
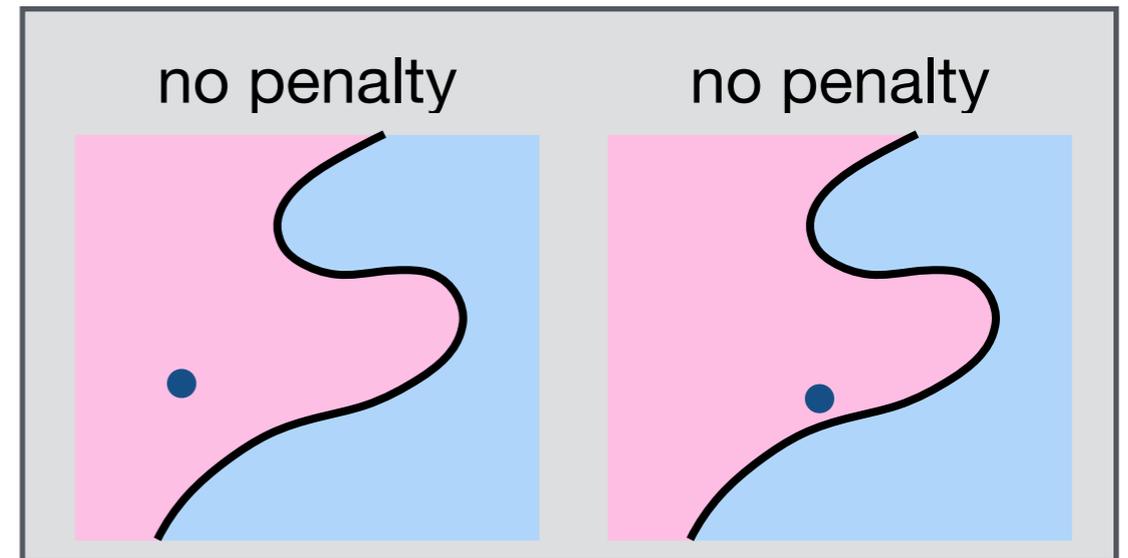
- Adversarial attacks:
manipulate predictions by adding imperceptible small noise



- More interests in whether our learning method are robust
 - ❖ important in applications such as autonomous driving, medical diagnosis

Formulation of Adversary

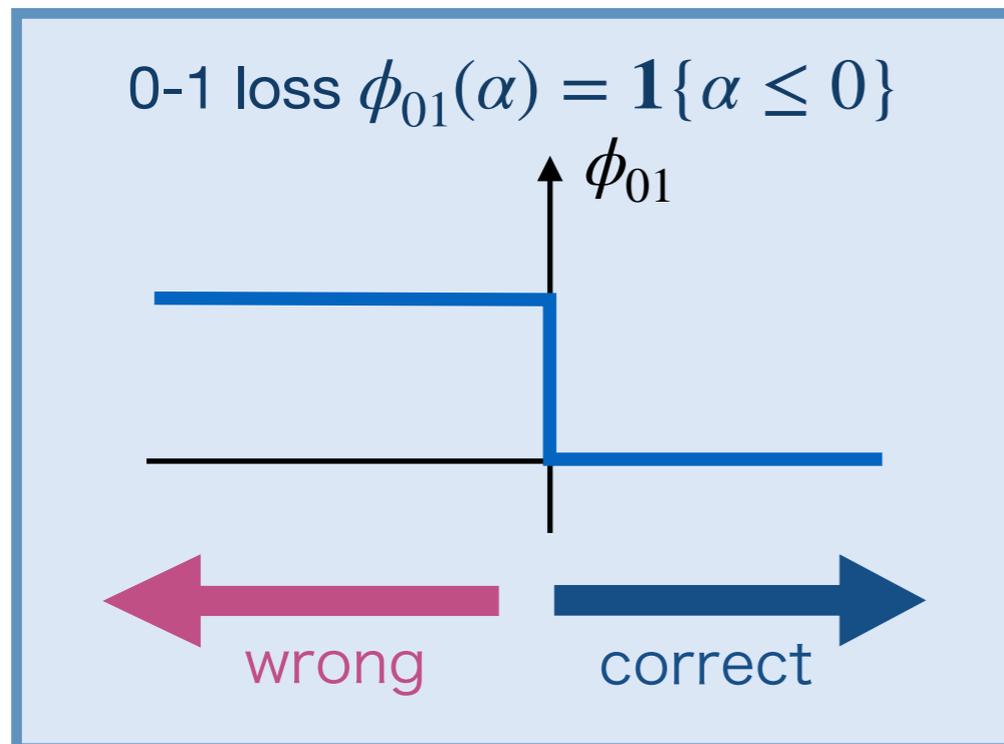
- Standard learning:
no penalty if classified to the correct side of the boundary
- Robust learning:
prediction close to the boundary will be penalized even if correctly classified
 - ❖ the boundary will be crossed over by attacks
 - ❖ assume L_2 -ball attack



Standard vs. Robust Learning

- Standard learning:
minimize 0-1 loss

$$R_{\phi_{01}}(f) = \mathbb{E} [\phi_{01}(Yf(X))]$$



- Robust learning:
minimize robust 0-1 loss

$$R_{\phi_{\gamma}}(f) = \mathbb{E} \left[\max_{\Delta \in B_2(\gamma)} \phi_{01}(Yf(X + \Delta)) \right]$$

worst L_2 -attack

learn best (min) classifier
under worst-case (max) attack
= **robust optimization**

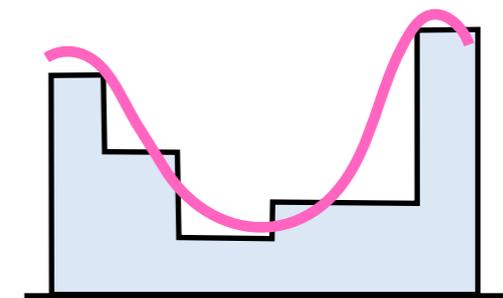
Relaxation of Robust Optimization

- Direct optimization of robust 0-1 loss is hard
- Existing relaxation

Not necessarily calibrated to robust 0-1 loss!

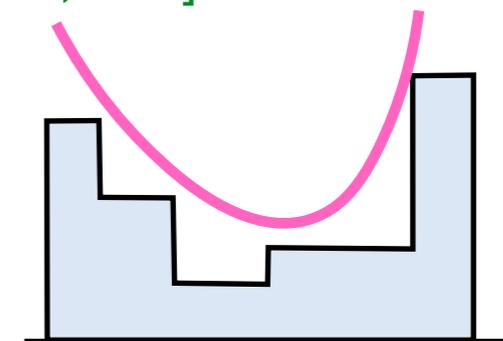
❖ Taylor approximation [Shaham *et al.* 2018; etc.]

local approximation of original objective
does not necessarily lead to global minimum



❖ Minimize convex upper bound [Wong & Kolter 2018; etc.]

global minimum of upper bound
does not necessarily equal to minima
of original objective

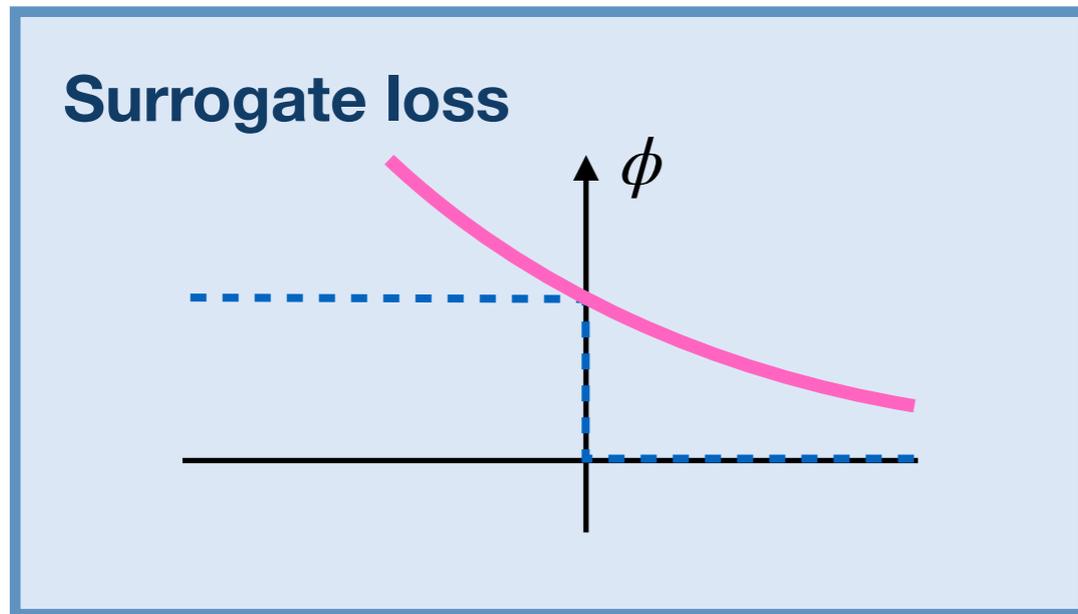


Shaham, U., Yamada, Y., & Negahban, S. (2018). Understanding adversarial training: Increasing local stability of supervised models through robust optimization. *Neurocomputing*, 195-204.

Wong, E., & Kolter, Z. (2018,). Provable Defenses against Adversarial Examples via the Convex Outer Adversarial Polytope. In *International Conference on Machine Learning* (pp. 5286-5295).

What is calibrated surrogates?

- Standard learning



calibrated
[Bartlett *et al.* 2006]

Target loss = 0-1 loss

$$R_{\phi_{01}}(f) = \mathbb{E} [\phi_{01}(Yf(X))]$$

- Robust learning

Q. What kind of losses are calibrated?

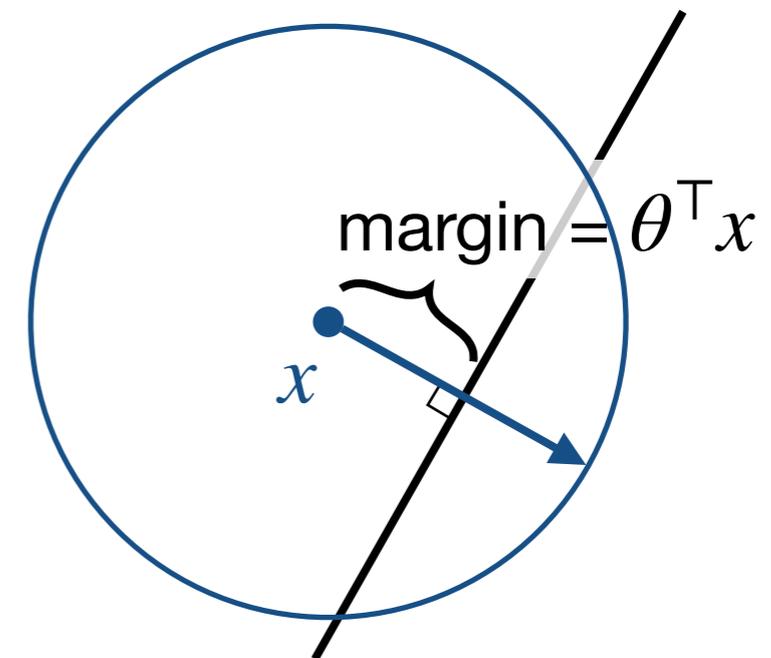
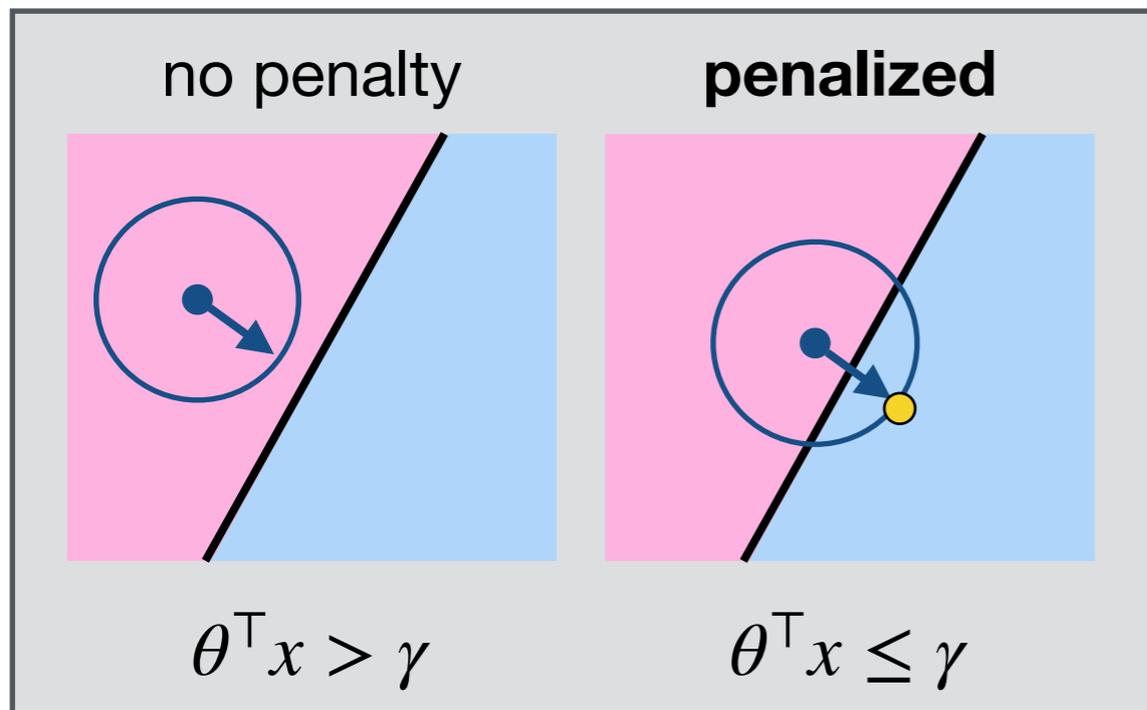
calibrated

Target loss = robust 0-1 loss

$$R_{\phi_{\gamma}}(f) = \mathbb{E} \left[\max_{\Delta \in B_2(\gamma)} \phi_{01}(Yf(X + \Delta)) \right]$$

Special case: linear model + L₂-attack 17

- Linear model $f_\theta(x) = \theta^\top x$ where $\|\theta\|_2 = 1$



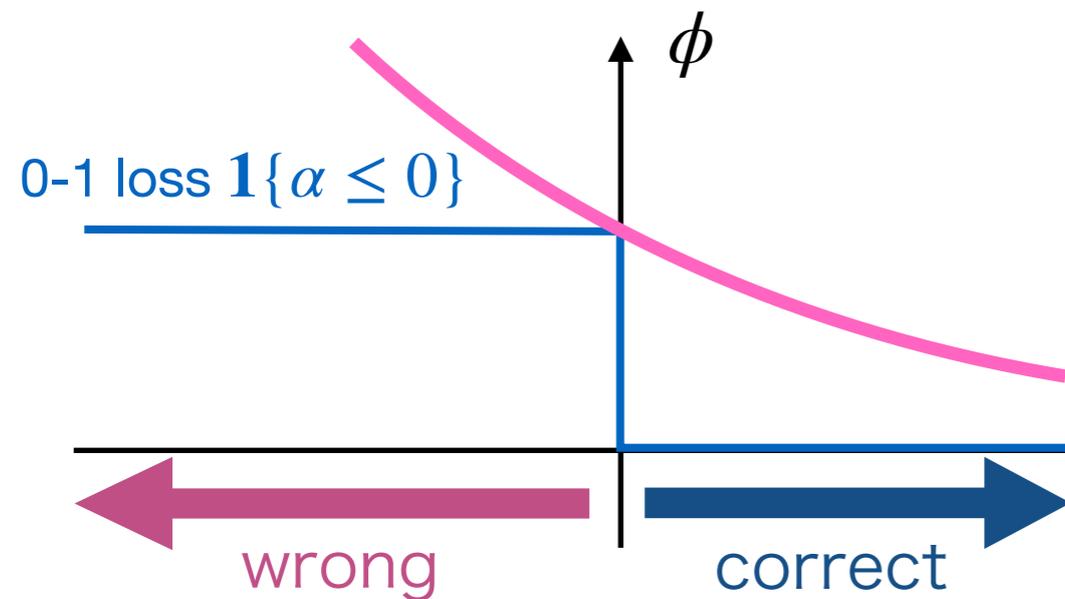
robust 0-1 loss

$$\max_{\Delta \in B_2(\gamma)} \phi_{01}(Yf(X + \Delta)) = \mathbf{1}\{Yf(X) \leq \gamma\} := \phi_\gamma(Yf(X))$$

- General case is hard to analyze

Isn't it a piece of cake?

● Standard learning



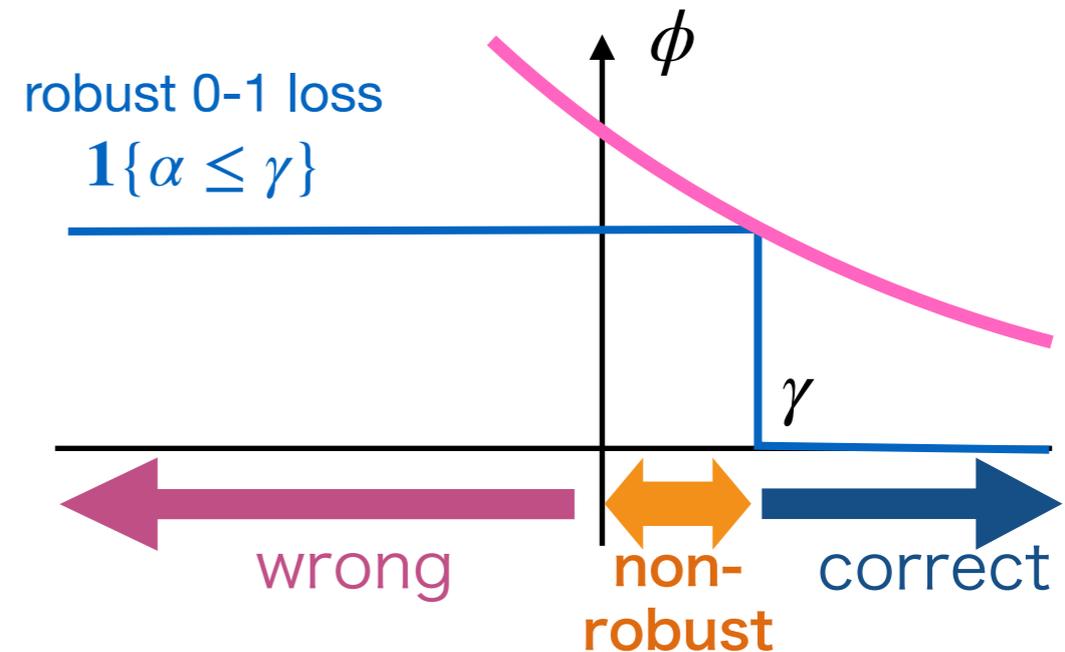
Theorem [Bartlett *et al.* 2006].

If surrogate ϕ is convex,

1. ϕ is differentiable at 0
2. $\phi'(0) < 0$

are necessary and sufficient for calibration.

● Robust learning



Conjecture.

If surrogate ϕ is convex,

1. ϕ is differentiable at $\alpha = \gamma$
2. $\phi'(\gamma) < 0$

are necessary and sufficient?

Main Result

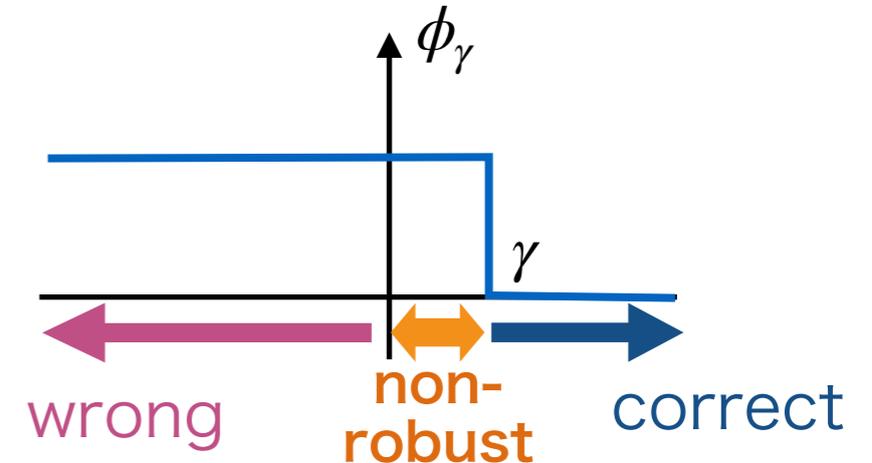
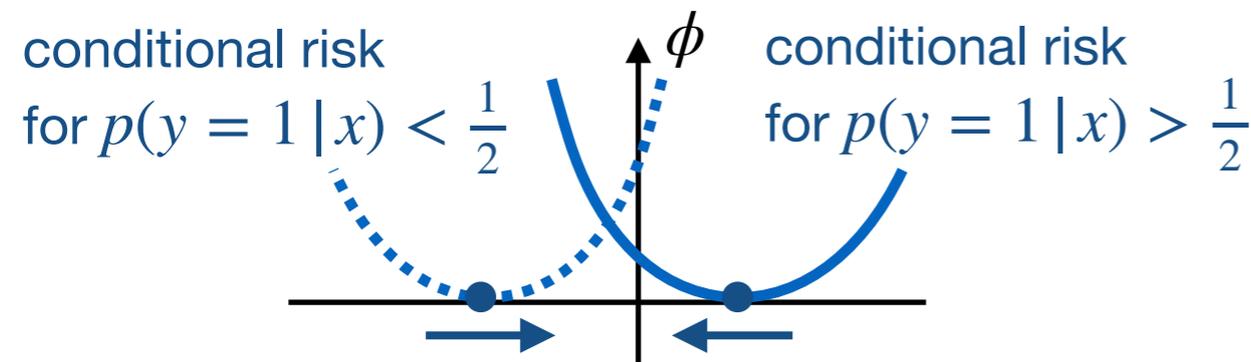
Theorem [Bao *et al.* 2020].

Any convex surrogates are not calibrated to robust 0-1 loss under linear models + L_2 attack.

- Intuition (Note: proven by checking $\delta(\varepsilon) = 0$ for some ε to be precise)

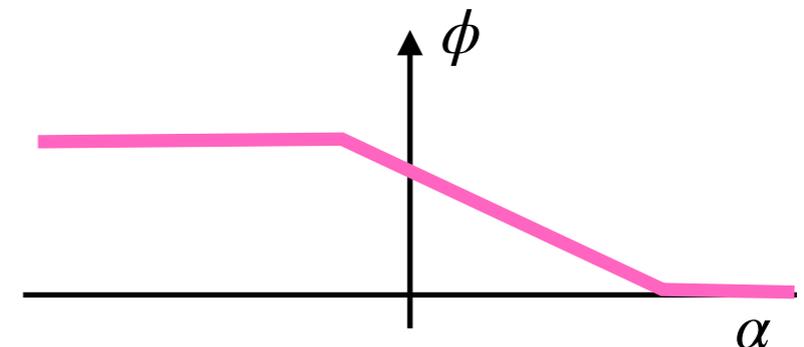
1. predictions becomes close to 0
as $p(y = 1 | x) \rightarrow \frac{1}{2}$

2. predictions close to 0 are regarded as non-robust



- Nonconvex calibrated surrogates exist

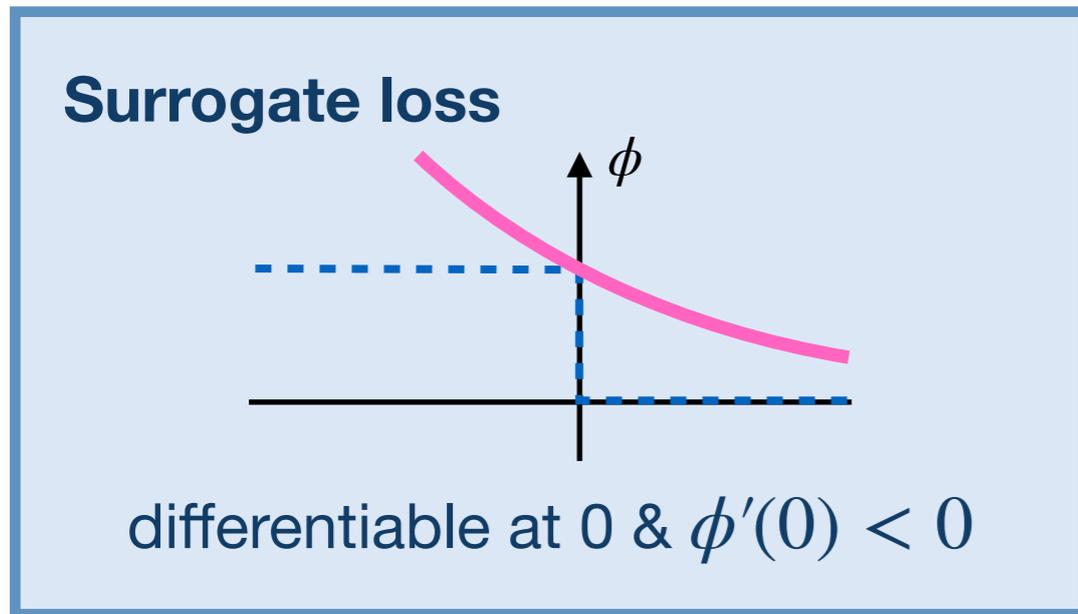
❖ e.g. ramp loss



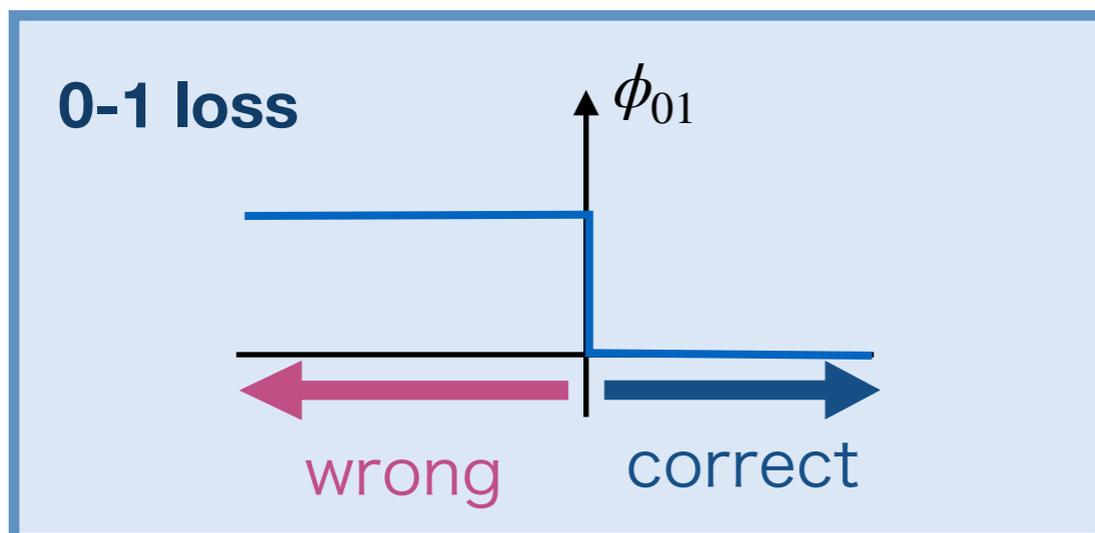
Summary |

Calibrated Surrogates and Robust Learning

● Standard learning



calibrated
[Bartlett *et al.* 2006]

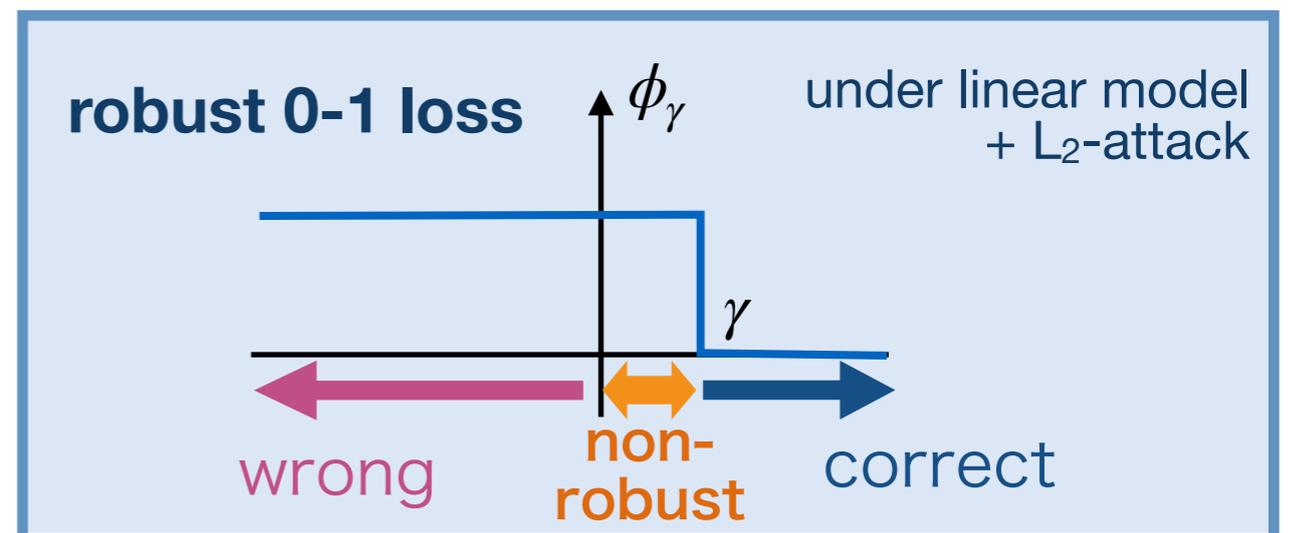


● Robust learning

Result

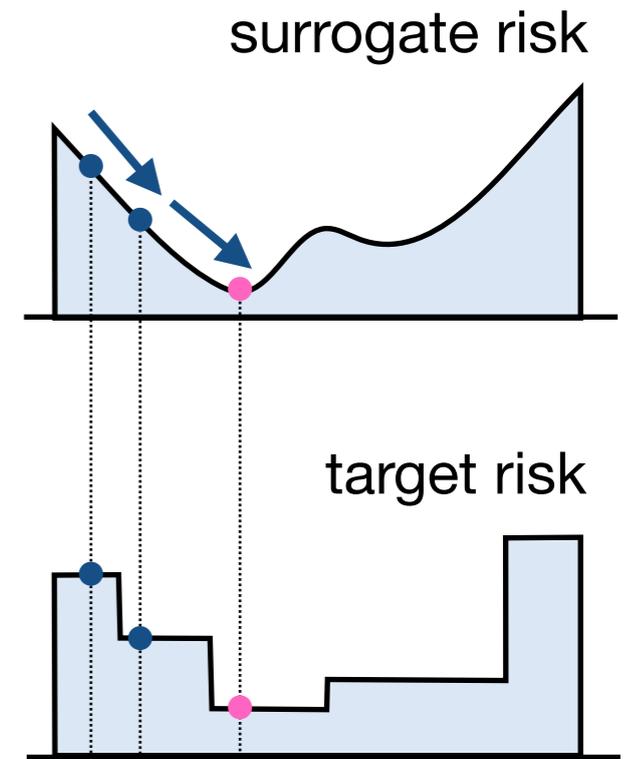
- no convex calibrated surrogates
- nonconvex calibrated surrogates exist (e.g. ramp loss)

calibrated
[Bao *et al.* 2020]



Summary

- **Calibrated surrogate losses:**
surrogate risk minimizer = target risk minimizer
 - ❖ can be confirmed via calibration function
- **Robust learning** from calibration perspective:
⇒ no convex calibrated losses
 - ❖ future: how about minimax surrogates?
- Take home:



calibration is interesting not only for minimizer consistency
but also for robust loss design!

- ❖ similar idea adopted to analyze robustness
to symmetric label noise [Reid & Williamson 2010]