

# Learning Theory Bridges

## Loss Functions

July 13<sup>rd</sup>, 2020

Han Bao (The University of Tokyo / RIKEN AIP)

# Han Bao (包 含)

<https://hermite.jp/>

- 2nd-year Ph.D. student @ Sugiyama-Honda-Yokoya Lab
- Research Interests:  
robustness and knowledge transfer via loss function

## robustness

Calibrated Surrogate Maximization of Linear-fractional Utility in Binary Classification. (AISTATS2020)

Calibrated Surrogate Losses for Adversarially Robust Classification. (COLT2020)

Calibrated surrogate maximization of Dice. (MICCAI2020)

## transfer learning

Unsupervised Domain Adaptation Based on Source-guided Discrepancy. (AAAI2019)

## similarity learning

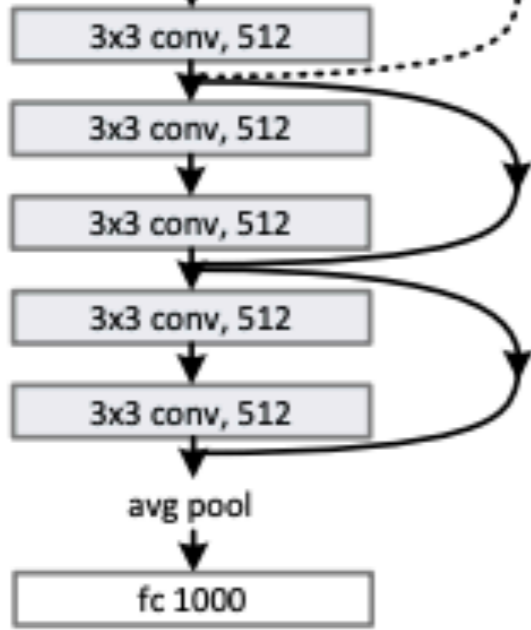
Similarity-based Classification: Connecting Similarity Learning to Binary Classification. (preprint)

Classification from Pairwise Similarity and Unlabeled Data. (ICML2018)

## knowledge transfer

# Deep Residual Learning for Image Recognition

Zhang Shaoqing Ren Jian Sun  
 Microsoft Research  
 {kaimingh, shren, jiansun}@microsoft.com

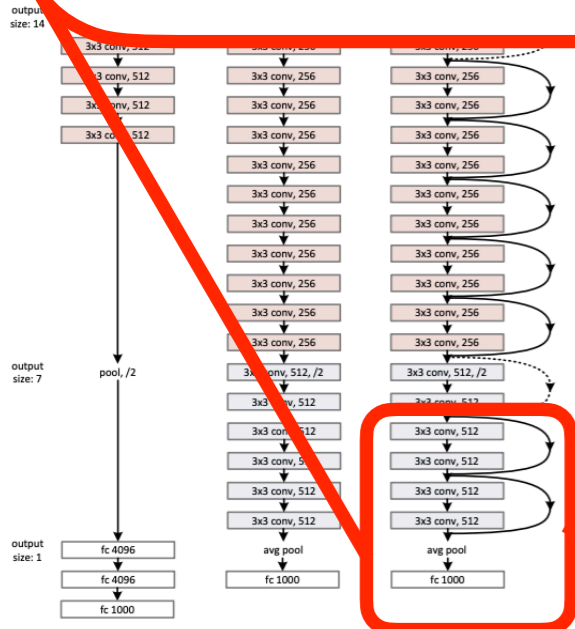
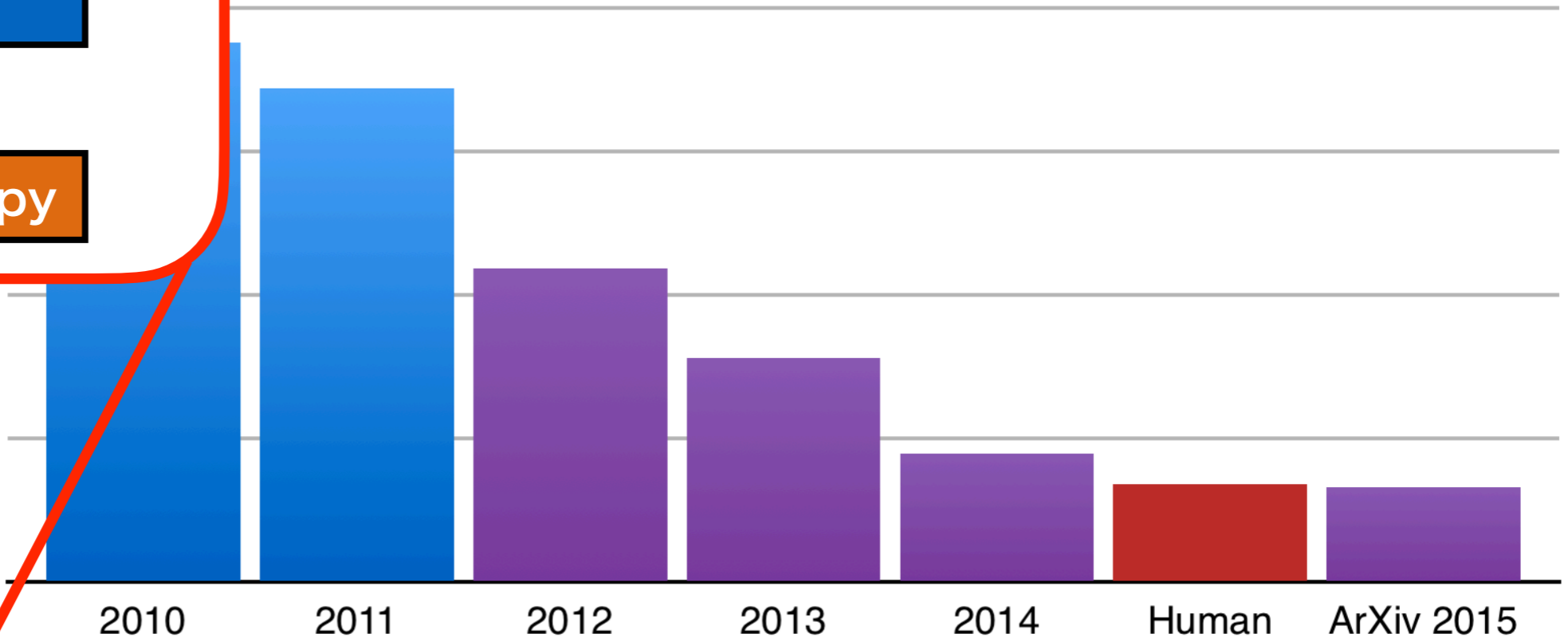


softmax



cross-entropy

## ILSVRC top-5 error on ImageNet



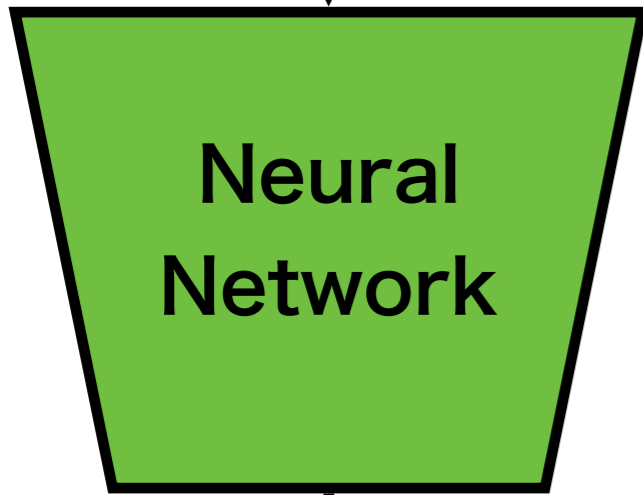
<https://devblogs.nvidia.com/mocha-jl-deep-learning-julia/image1/>

# Training

feature ( $x$ ) label ( $y$ )



traffic  
light



softmax

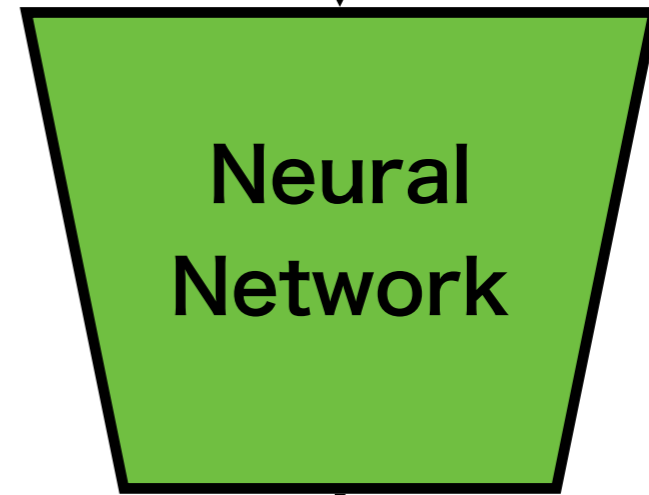
cross-entropy

→ minimize

Distance of  
label and prediction

# Prediction

feature ( $x$ )

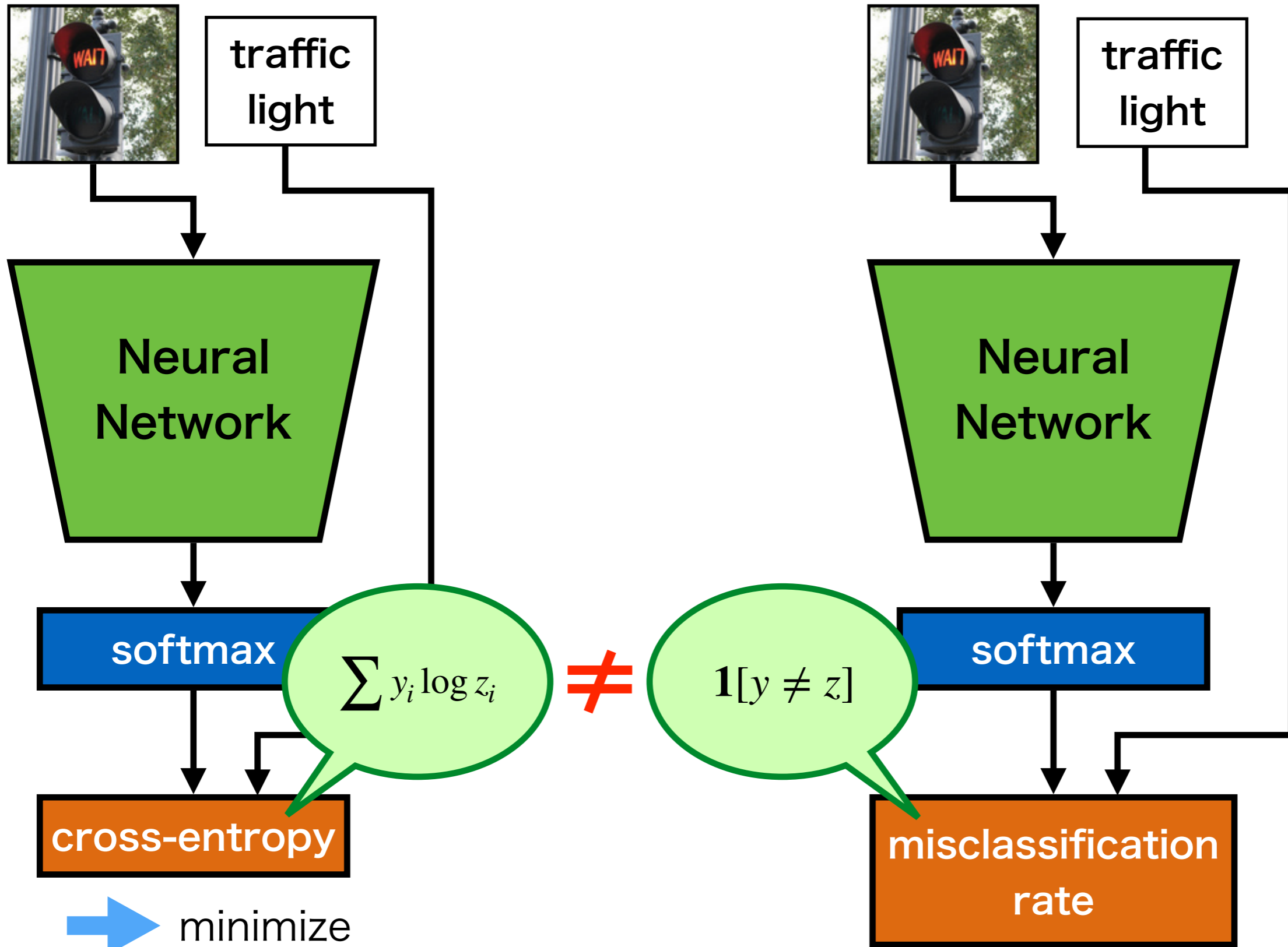


softmax

traffic  
light?

# Training

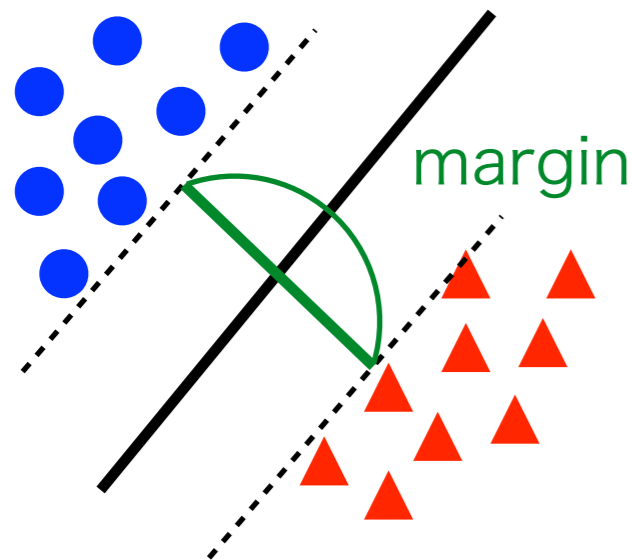
# Evaluation



# Support-Vector Networks

CORINNA CORTES  
VLADIMIR VAPNIK  
*AT&T Bell Labs., Holmdel, NJ 07733, USA*

corinna@neural.att.com  
vlad@neural.att.com



margin maximization

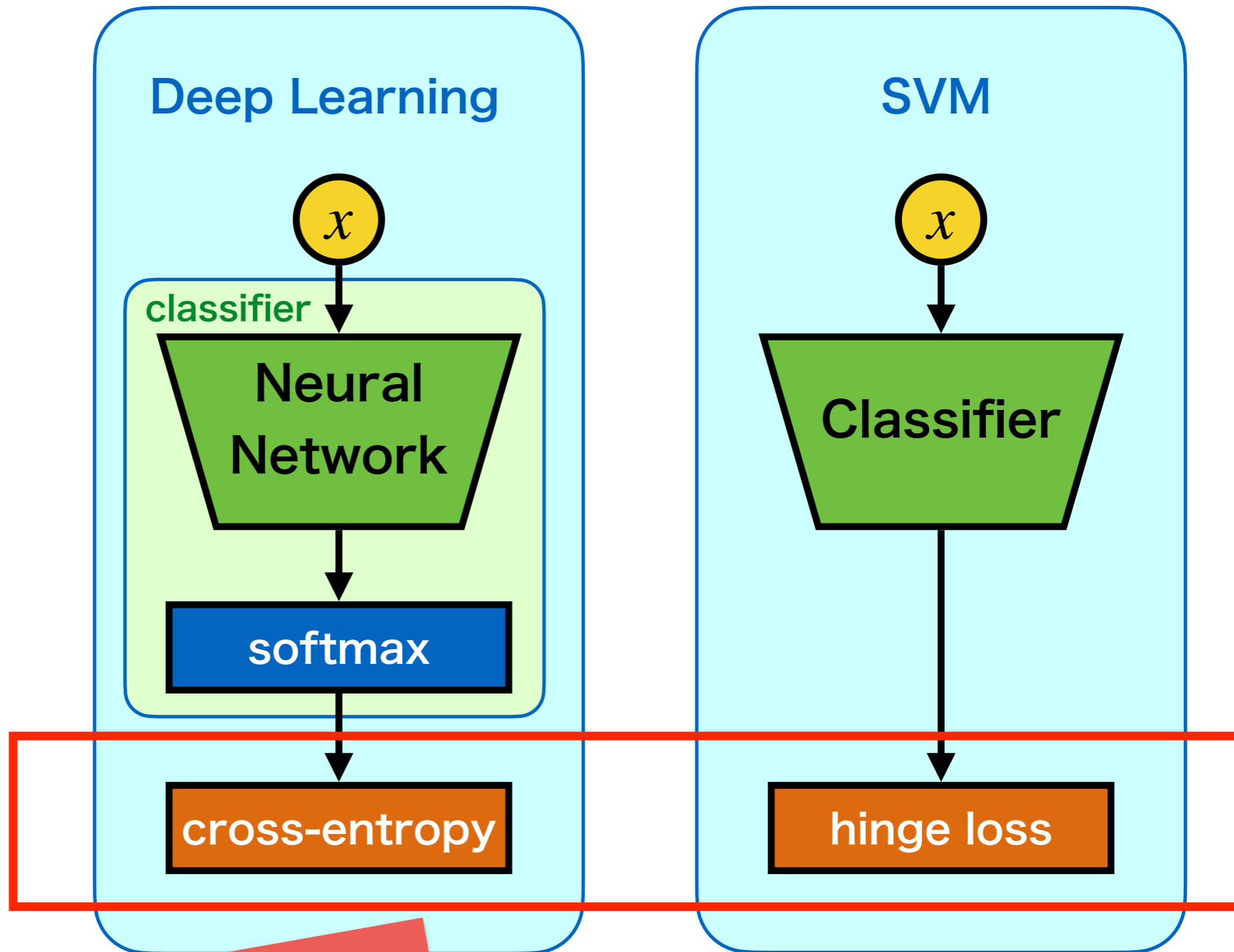
=

$$\min_{w,b} \sum_i \max \{ 0, 1 - y_i(w^\top x_i + b) \}$$

hinge loss minimization



misclassification rate



learning = minimize loss

Does it work?

~~||~~

misclassification rate

# Background: Binary Classification

8

## Input

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

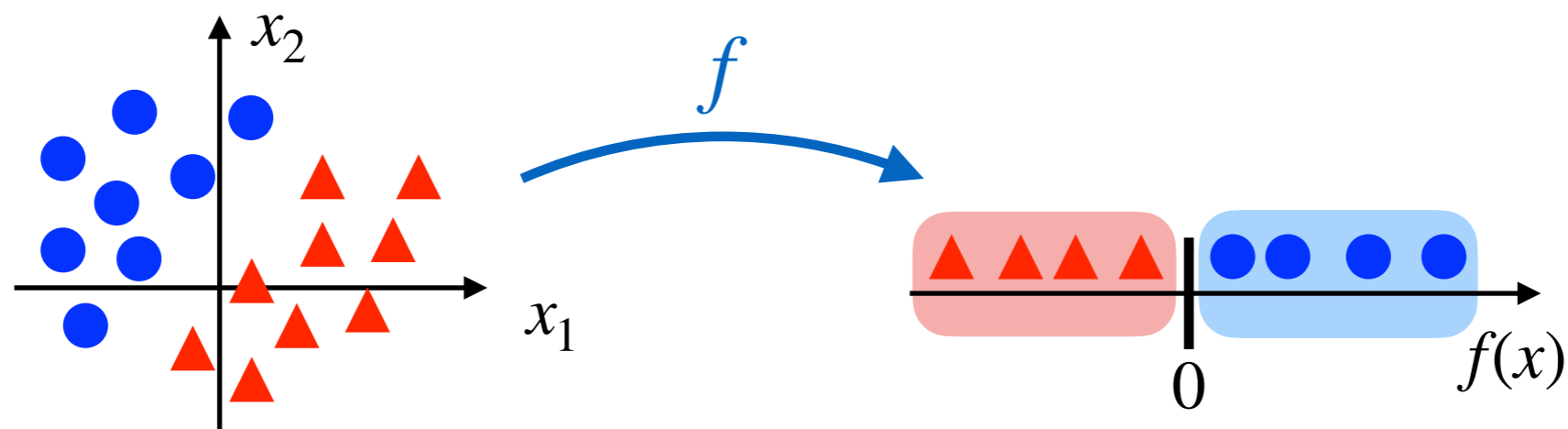
## Output

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

- ▶ predict class by  $\text{sign}(f(\cdot))$

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

1 if  $Y \neq \text{sign}(f(X))$ ,  
0 if  $Y = \text{sign}(f(X))$





# Loss function and Risk

- Goal of classification: minimize misclassification rate

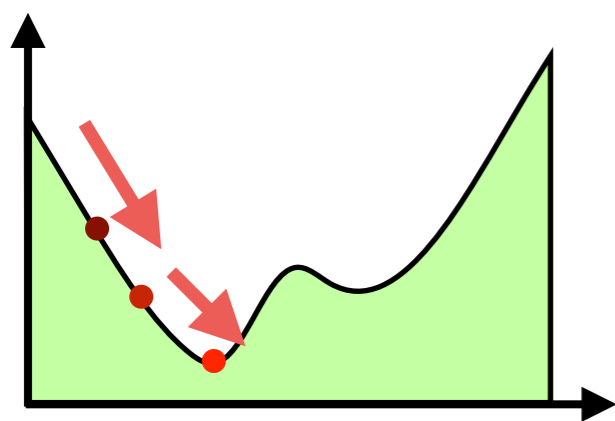
$$R_{01}(f) = \mathbb{E} [\mathbf{1}[Y \neq \text{sign}(f(X))]]$$

0-1 risk

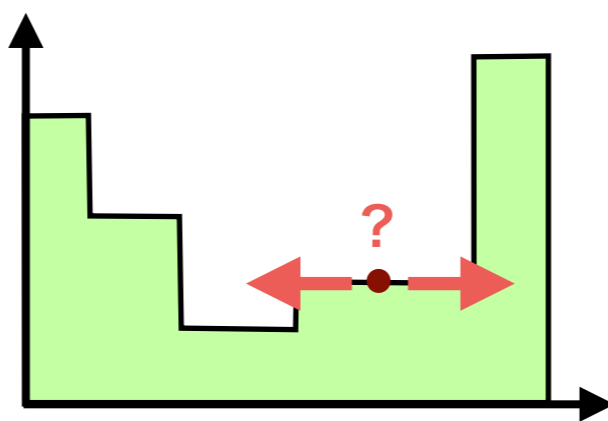
- Misclassification rate = expectation of 0-1 loss

$$\mathbf{1}[Y \neq \text{sign}(f(X))] = \phi_{01}(Yf(X))$$

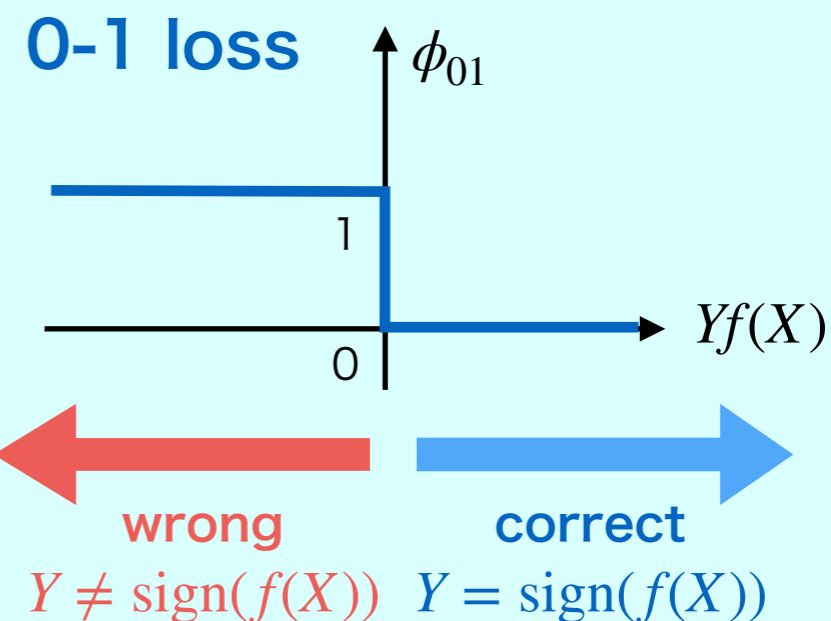
- Minimizing  $R_{01}$  is NP-hard [Feldman+ 2012]



minimization by  
gradient descent



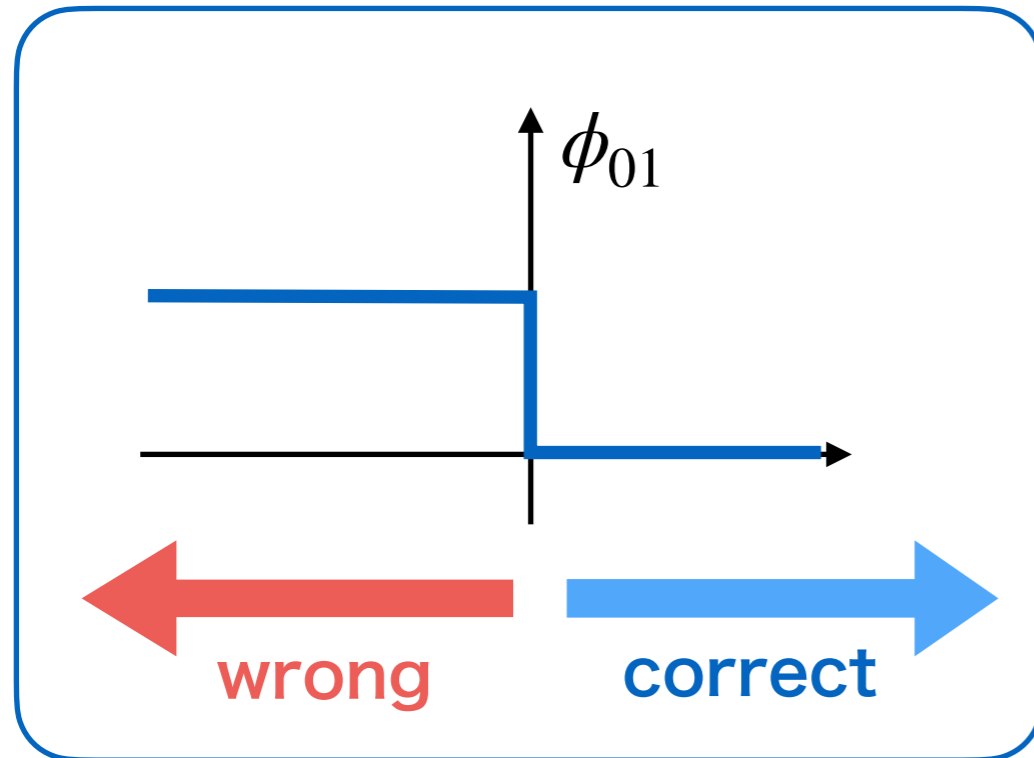
no gradient for  
discrete function



# Target Loss vs. Surrogate Loss

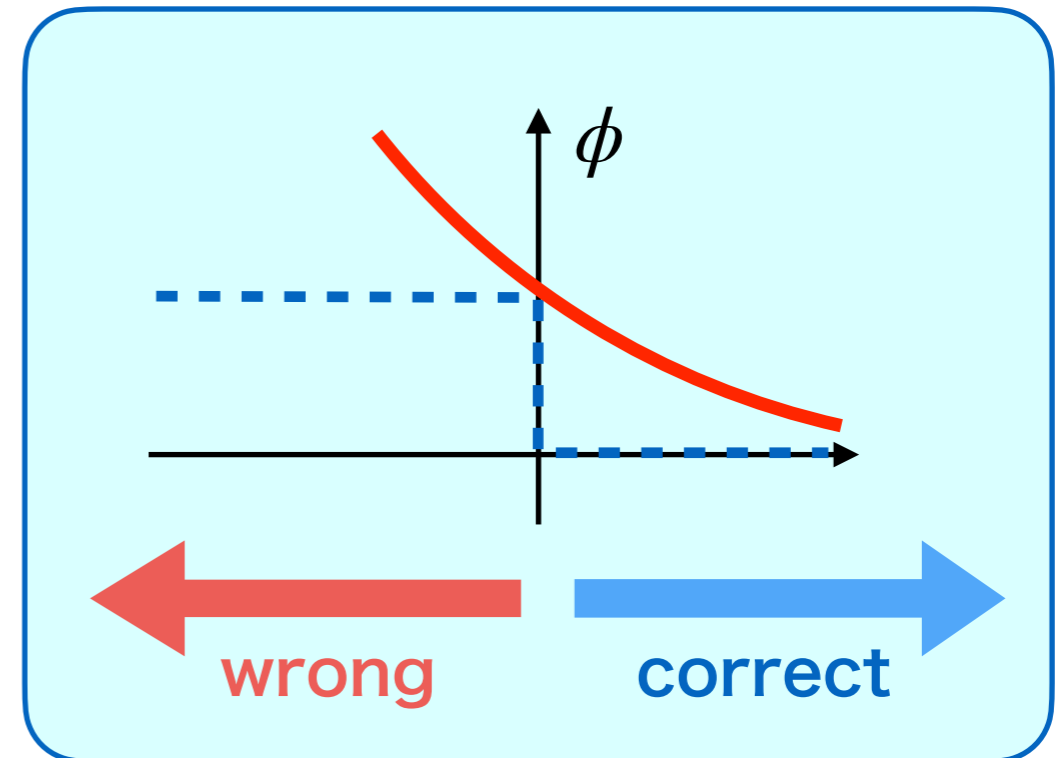
10

Target loss (0-1 loss)



- Final learning criterion
- Hard to optimize
  - ▶ nonconvex, no gradient

Surrogate loss



- Different from target loss
- Easily-optimizable criterion
  - ▶ usually convex, smooth

# Elements of Learning Theory <sup>11</sup>

(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**

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

**Generalization theory:**

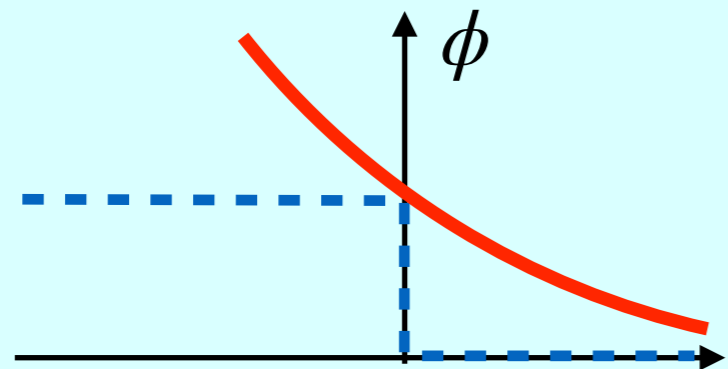
If model is not too complicated,  
then converges (roughly speaking)

Key ingredient:

**Calibration theory** for  
loss functions

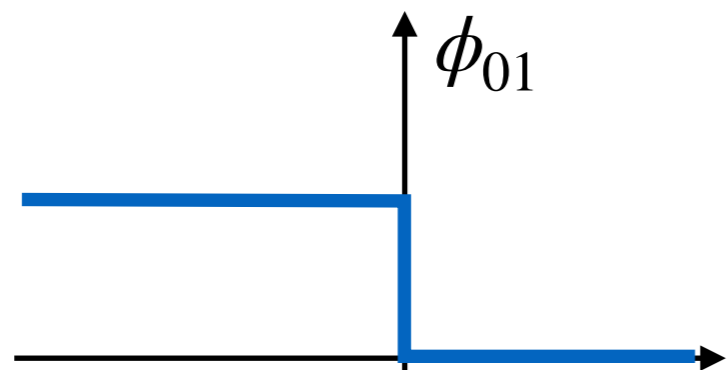
# What surrogate is desirable? <sup>12</sup>

## Surrogate loss



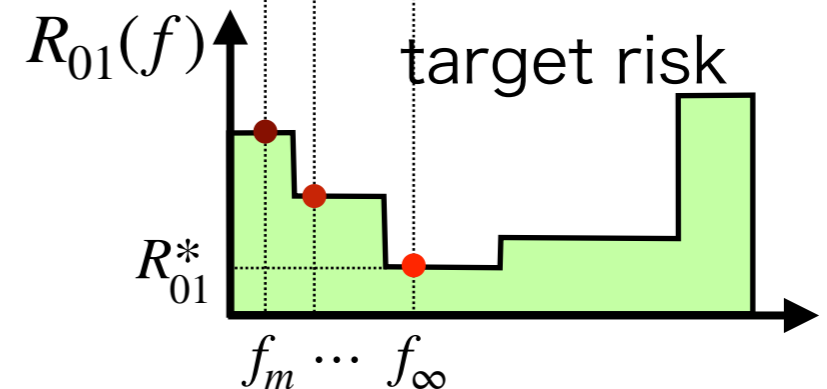
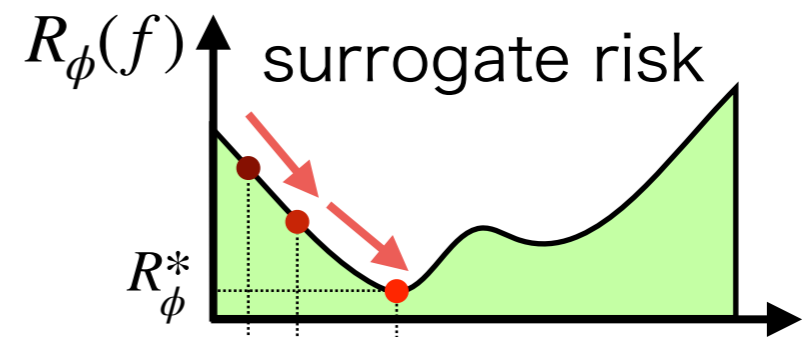
easily optimizable

## Target loss (0-1 loss)



final learning criterion

## Calibrated surrogate



$$R_\phi(f_m) \xrightarrow{m \rightarrow \infty} R_\phi^* \implies R_{01}(f_m) \xrightarrow{m \rightarrow \infty} R_{01}^*$$

# How to check risk convergence? <sup>13</sup>

[Steinwart 2007]

**Definition.**  $\phi$  is  $\psi$ -**calibrated** for a target loss  $\psi$

if for any  $\varepsilon > 0$ , there exists  $\delta > 0$  such that for all  $f$ ,

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

surrogate (excess) risk

target (excess) risk



**Idea:** write  $\delta$  as function of  $\varepsilon$  (by using contraposition)

**Definition. (calibration function)**

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

surrogate (excess) risk

target (excess) risk

If  $\delta(\varepsilon) > 0$  for all  $\varepsilon > 0$ , surrogate is calibrated!

# Main Tool: Calibration Function

14

## Definition. (calibration function)

$$\delta(\varepsilon) = \inf_f \underbrace{R_\phi(f) - R_\phi^*}_{\text{surrogate (excess) risk}} \quad \text{s.t.} \quad \underbrace{R_\psi(f) - R_\psi^*}_{\text{target (excess) risk}} \geq \varepsilon$$

### ■ Provides **iff condition**

▶  $\psi$ -calibrated  $\iff \delta(\varepsilon) > 0$  for all  $\varepsilon > 0$

### ■ Provides **excess risk bound** monotonically increasing

▶  $\psi$ -calibrated  $\implies \underbrace{R_\psi(f) - R_\psi^*}_{\text{target excess risk}} \leq \underbrace{(\delta^{**})^{-1}}_{\text{monotonically increasing}} \left( \underbrace{R_\phi(f) - R_\phi^*}_{\text{surrogate excess risk}} \right)$

$\delta^{**}$ : biconjugate of  $\delta$

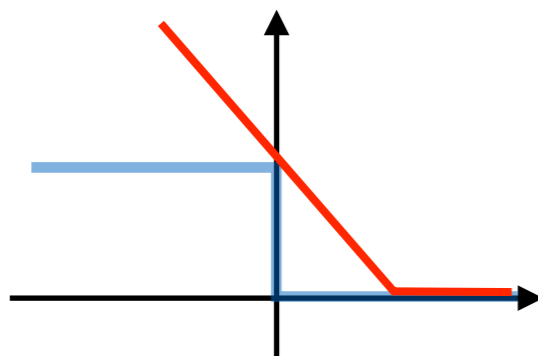
# Example: Binary Classification ( $\phi_{01}$ ) <sup>15</sup>

[Bartlett+ 2006]

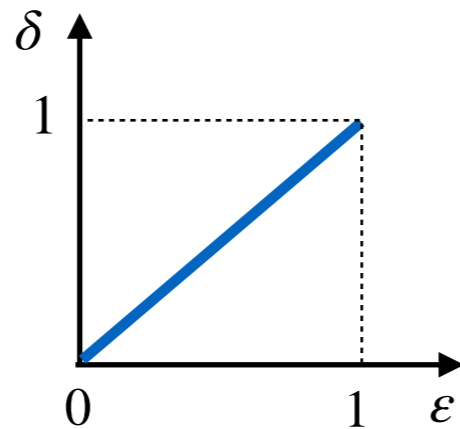
**Theorem.** If surrogate  $\phi$  is convex, it is  $\phi_{01}$ -calibrated iff

- ▶ differentiable at 0
- ▶  $\phi'(0) < 0$

hinge loss

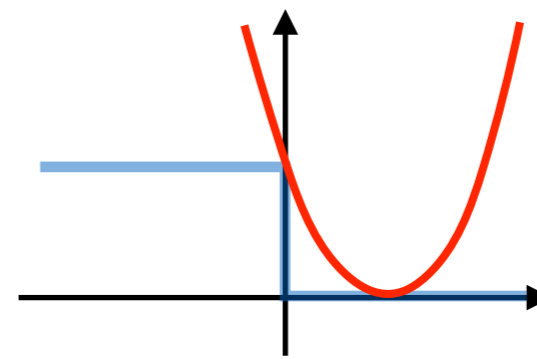


$$\phi(\alpha) = [1 - \alpha]_+$$

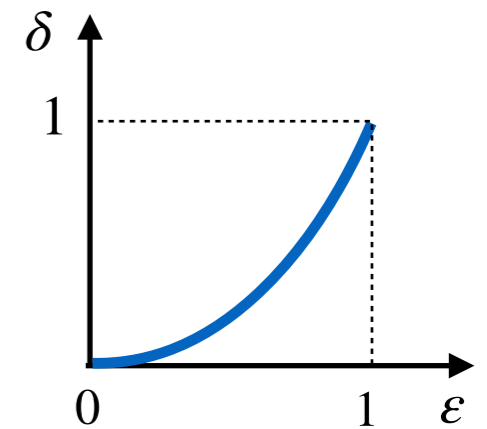


$$\delta(\epsilon) = \epsilon$$

squared loss



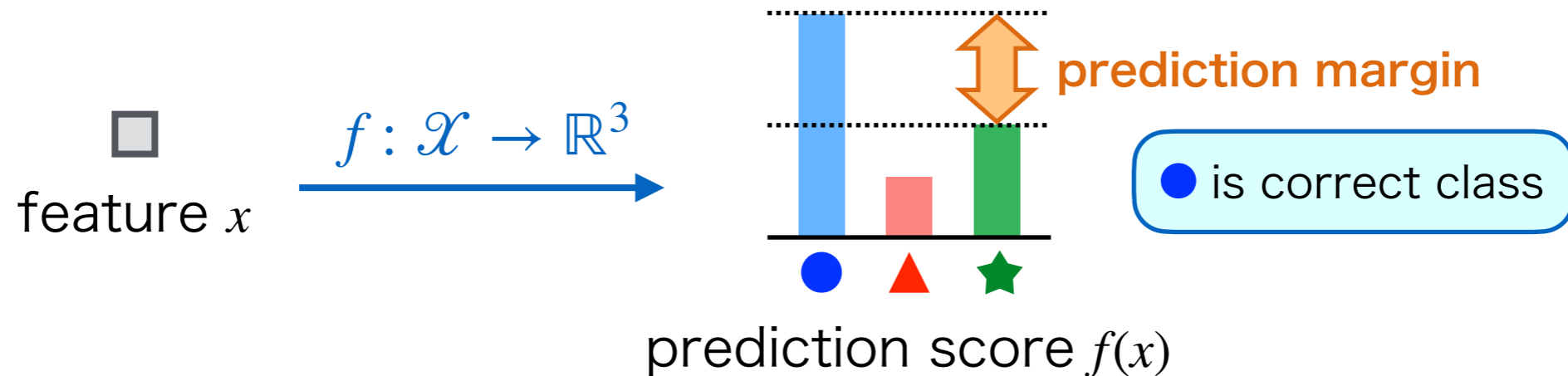
$$\phi(\alpha) = (1 - \alpha)^2$$



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

# Counterintuitive Result

- e.g. multi-class classification  $\Rightarrow$  maximize prediction margin



## Crammer-Singer loss

$$\max\{0, 1 - \text{pred. margin}\}$$

[Crammer & Singer 2001]

one of multi-class extensions  
of hinge loss

**Crammer-Singer loss is not calibrated to 0-1 loss !**

(similar extension of logistic loss is calibrated)

[Zhang 2004]

Crammer, K., & Singer, Y. (2001). On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of machine learning research*, 2(Dec), 265-292

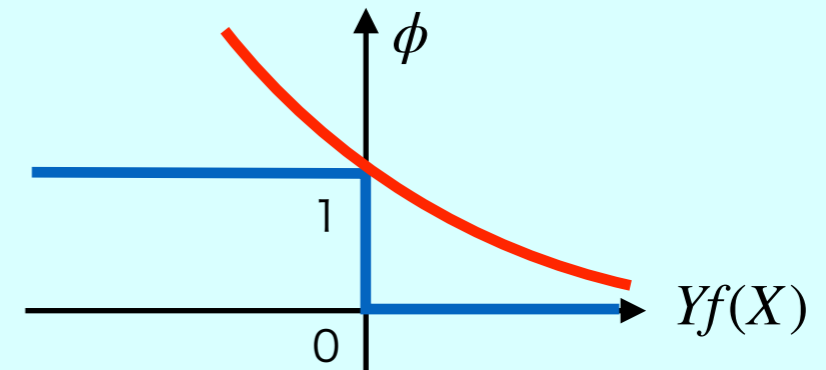
Zhang, T. (2004). Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5(Oct), 1225-1251.



# Summary: Calibration Theory 17

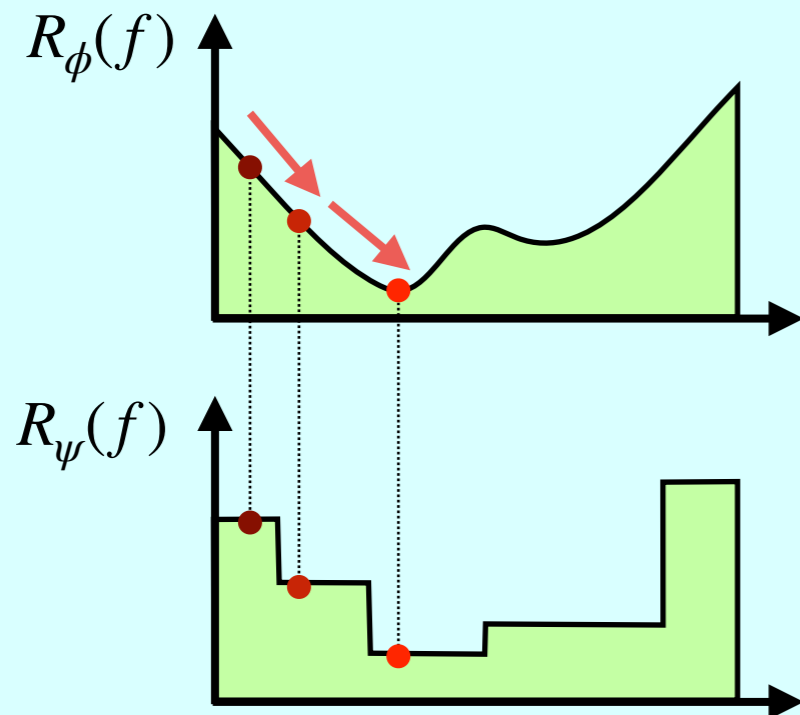
## Surrogate vs. Target loss

**Target** loss is often hard to optimize  
⇒ replace with **surrogate** loss



## Calibrated Surrogate

leading to minimization of target



## Binary Classification

Hinge, logistic is calibrated  
Calibrated iff  $\phi'(0) < 0$

## Multi-class Classification

CS-loss (MC-hinge loss) is  
not calibrated!

cross-entropy is calibrated  
(omitted)

Stringent justification of surrogate loss!

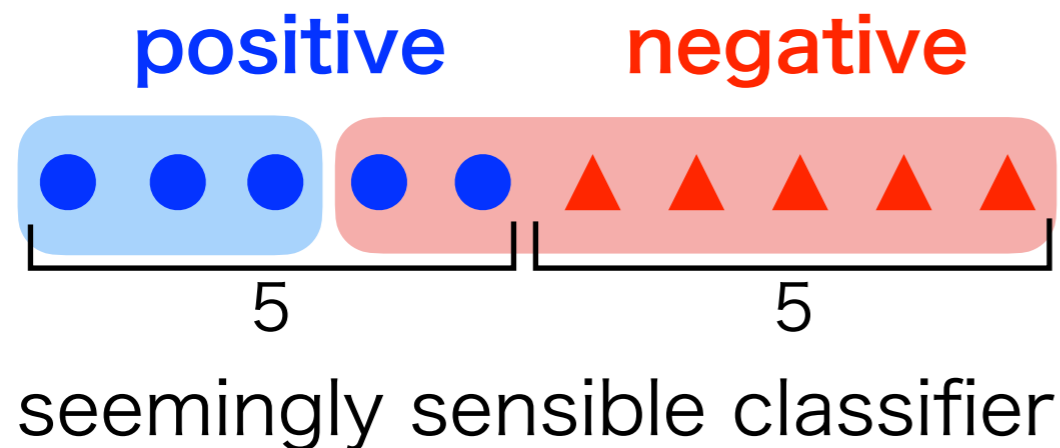
# When target is not 0-1 loss

**H. Bao** and M. Sugiyama.

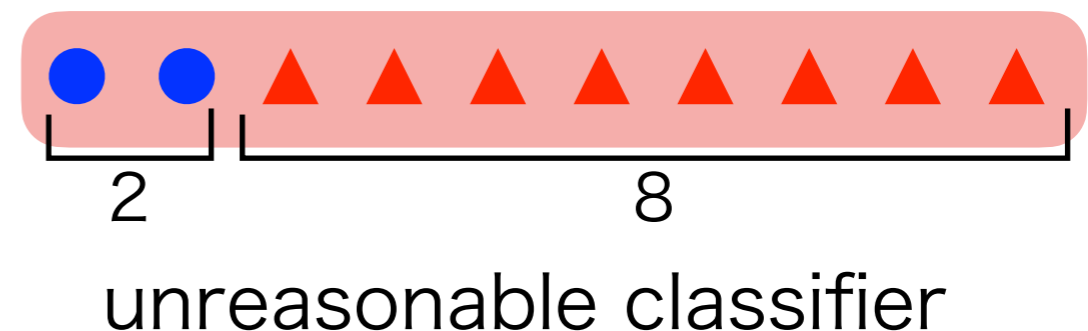
Calibrated Surrogate Maximization of Linear-fractional Utility in Binary Classification. In *AISTATS*, 2020.

# Is accuracy appropriate?

- Our focus: **binary classification**



accuracy: **0.8**

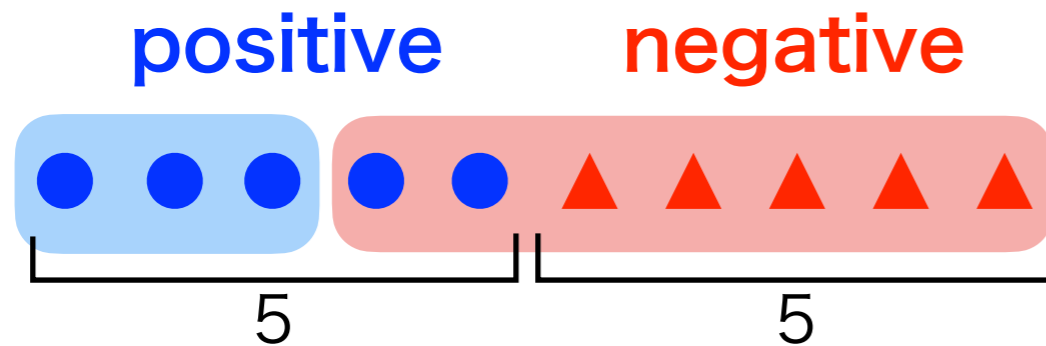


accuracy: **0.8**

Accuracy can't detect unreasonable classifiers  
under **class imbalance!**

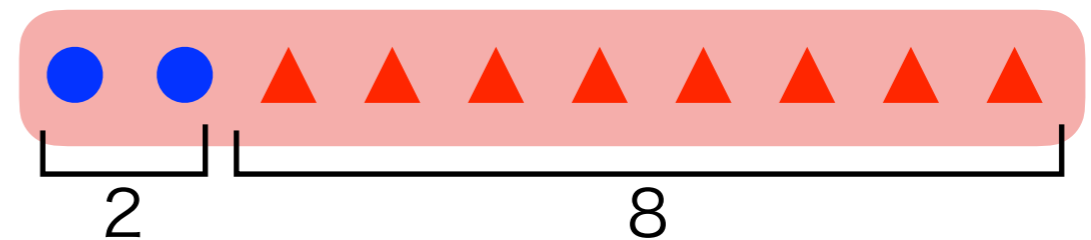
# Is accuracy appropriate?

- F-measure is more appropriate under **class imbalance**



accuracy: **0.8**

F-measure: **0.75**



accuracy: **0.8**

F-measure: **0**

$$\text{F-measure } F_1 = \frac{2TP}{2TP + FP + FN}$$

$$TP = \mathbb{E}_{X, Y=+1} [1_{\{f(X) > 0\}}]$$

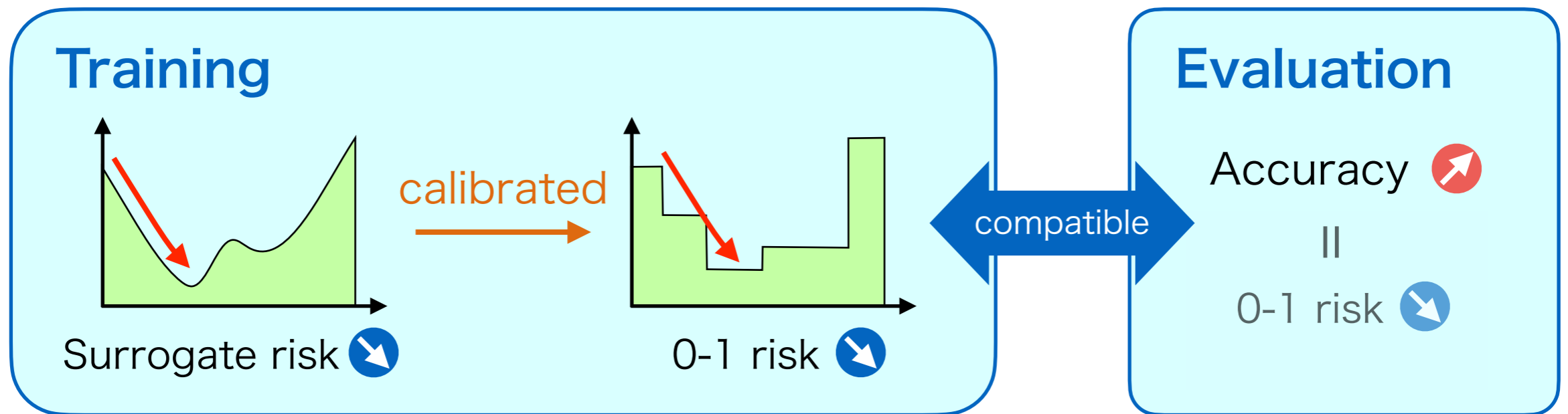
$$TN = \mathbb{E}_{X, Y=-1} [1_{\{f(X) < 0\}}]$$

$$FP = \mathbb{E}_{X, Y=-1} [1_{\{f(X) > 0\}}]$$

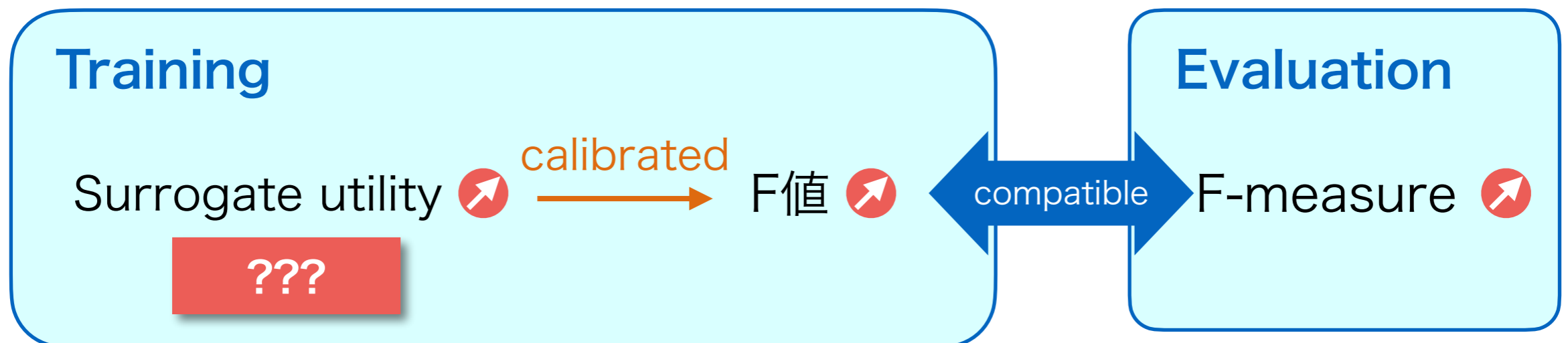
$$FN = \mathbb{E}_{X, Y=+1} [1_{\{f(X) < 0\}}]$$

# Training and Evaluation

## ■ Usual training with accuracy



## ■ Training with accuracy but evaluating with F-measure



# Not only $F_1$ , but many others

Q. Can we handle in the same way?

Accuracy

$$\text{Acc} = \text{TP} + \text{TN}$$

Weighted Accuracy

$$\text{WAcc} = \frac{w_1 \text{TP} + w_2 \text{TN}}{w_1 \text{TP} + w_2 \text{TN} + w_3 \text{FP} + w_4 \text{FN}}$$

F-measure

$$F_1 = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

Balanced Error Rate

$$\text{BER} = \frac{1}{\pi} \text{FN} + \frac{1}{1 - \pi} \text{FP}$$

Gower-Legendre index

$$\text{GLI} = \frac{\text{TP} + \text{TN}}{\text{TP} + \alpha(\text{FP} + \text{FN}) + \text{TN}}$$

Jaccard index

$$\text{Jac} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}}$$

# Unification of Metrics

Actual Metrics

$$F_1 = \frac{2TP}{2TP + FP + FN}$$

$$\text{Jac} = \frac{TP}{TP + FP + FN}$$

Note:

$$TN = \mathbb{P}(Y = -1) - FP$$

$$FN = \mathbb{P}(Y = +1) - TP$$

linear-fraction

$$U(f) = \frac{a_0 TP + b_0 FP + c_0}{a_1 TP + b_1 FP + c_1}$$

$a_k, b_k, c_k$  : constants

# Unification of Metrics

linear-fraction

$$U(f) = \frac{a_0 \text{TP} + b_0 \text{FP} + c_0}{a_1 \text{TP} + b_1 \text{FP} + c_1}$$

expectation divided by expectation

$$= \frac{a_0 \mathbb{E}_P + b_0 \mathbb{E}_N + c_0}{a_1 \mathbb{E}_P + b_1 \mathbb{E}_N + c_1} = \frac{\mathbb{E}_X[W_0(f(X))]}{\mathbb{E}_X[W_1(f(X))]}$$

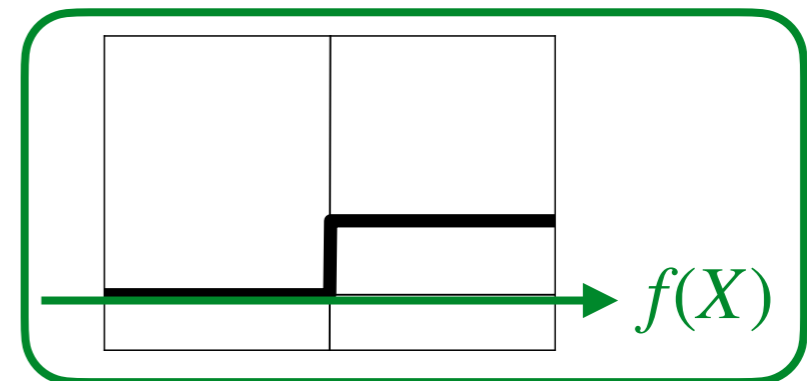
■ TP, FP = expectation of 0/1-loss

▶ TP =  $\mathbb{E}_{X, Y=+1} [\mathbf{1}[f(X) > 0]]$

positive data && positive prediction

▶ FP =  $\mathbb{E}_{X, Y=-1} [\mathbf{1}[f(X) > 0]]$

negative data && positive prediction





# Goal of This Talk

Given a metric (utility)  $U(f) = \frac{a_0 \text{TP} + b_0 \text{FP} + c_0}{a_1 \text{TP} + b_1 \text{FP} + c_1}$

## Q. How to optimize $U(f)$ directly?

- ▶ without estimating class-posterior probability

labeled sample  $\{(x_i, y_i)\}_{i=1}^n$  i.i.d.  $\mathbb{P}$   
metric  $U$



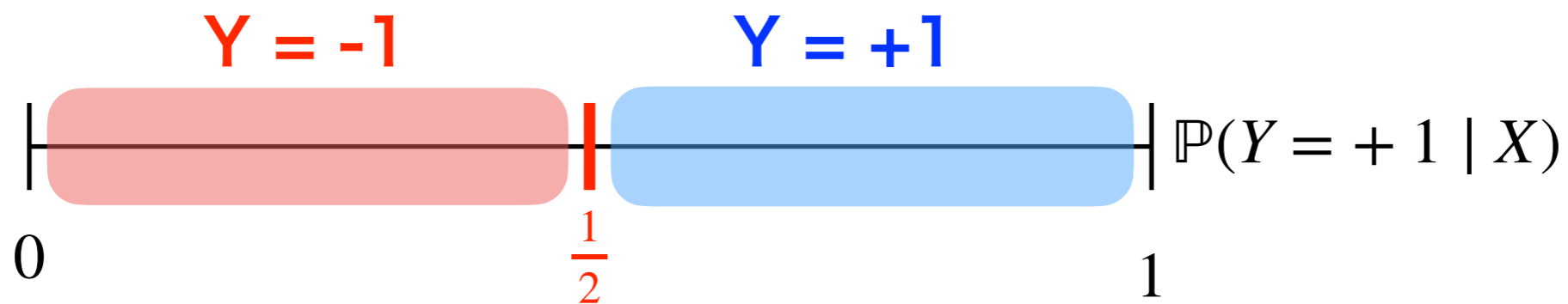
classifier  $f: \mathcal{X} \rightarrow \mathbb{R}$   
s.t.  $U(f) = \sup_{f'} U(f')$

# Related: Plug-in Classifier

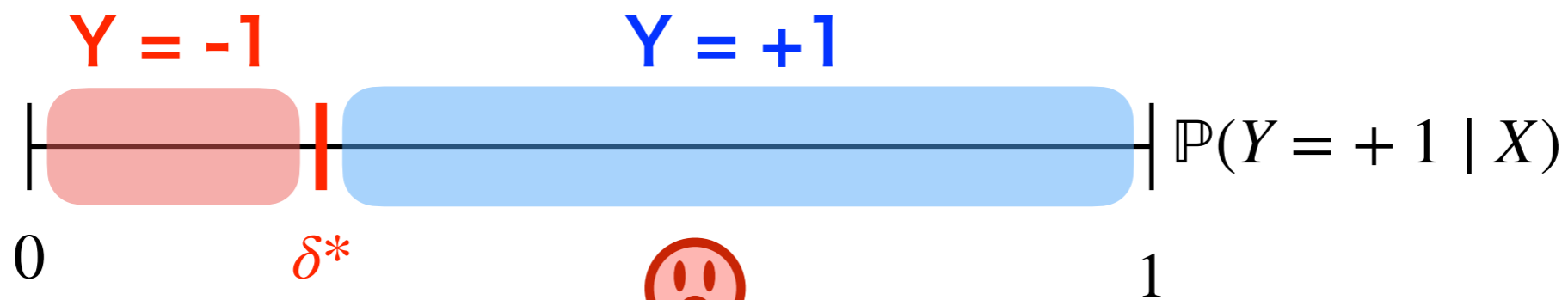
[Koyejo+ NIPS2014; Yan+ ICML2018]

- Estimating class-posterior probability is costly!

Bayes-optimal classifier (accuracy):  $\mathbb{P}(Y = +1 | x) - \frac{1}{2}$



Bayes-optimal classifier (general case):  $\mathbb{P}(Y = +1 | x) - \delta^*$



$\Rightarrow$  estimate  $\mathbb{P}(Y = +1 | x)$  and  $\delta^*$  independently

O. O. Koyejo, N. Natarajan, P. K. Ravikumar, & I. S. Dhillon.  
Consistent binary classification with generalized performance metrics. In *NIPS*, 2014.

B. Yan, O. Koyejo, K. Zhong, & P. Ravikumar.  
Binary classification with Karmic, threshold-quasi-concave metrics. In *ICML*, 2018.

# Convexity & Statistical Property <sup>27</sup>

Q. How to make tractable surrogate?

## Accuracy

tractable (convex)

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



calibrated

intractable

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

## Linear-fractional Metrics

① tractable?



② calibrated?

intractable

$$U(f) = \frac{\mathbb{E}_X[W_0(f(X))]}{\mathbb{E}_X[W_1(f(X))]}$$

# Non-concave, but quasi-concave 28

Idea:  $\frac{\text{concave}}{\text{convex}} = \underline{\text{quasi-concave}}$

$\frac{f(x)}{g(x)}$  is quasi-concave

if  $f$  : concave,  $g$  : convex,

$f(x) \geq 0$  and  $g(x) > 0$  for  $\forall x$

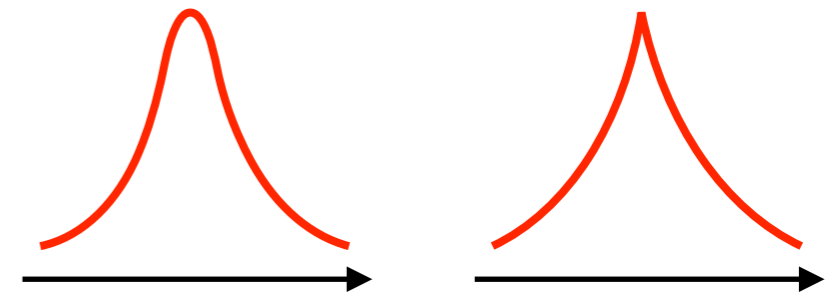
(proof) Show  $\{x | f/g \geq \alpha\}$  is convex.

$$\frac{f(x)}{g(x)} \geq \alpha \iff \underbrace{f(x) - \alpha g(x)}_{\text{concave}} \geq 0$$

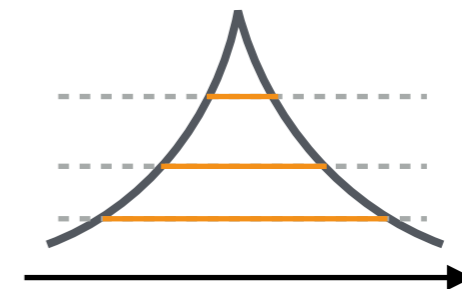
NB: super-level set of concave func.  
is convex

$\therefore \{x | f/g \geq \alpha\}$  is convex for  $\forall \alpha \geq 0$

non-concave, but unimodal  
 $\Rightarrow$  efficiently optimized

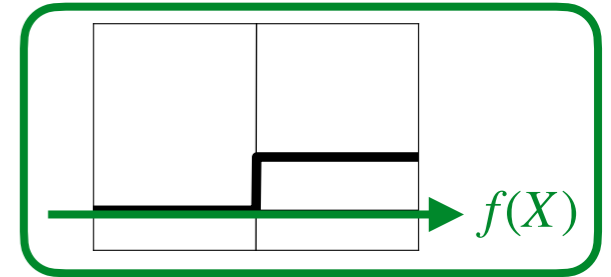


- quasi-concave  $\supseteq$  concave
- super-levels are convex



# Surrogate Utility

- Idea: bound true utility from below



linear-fraction

$$U(f) = \frac{a_0 TP + b_0 FP + c_0}{a_1 TP + b_1 FP + c_1}$$

$$= \frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_1}$$

numerator from below

non-negative sum of concave  
 $\Rightarrow$  concave

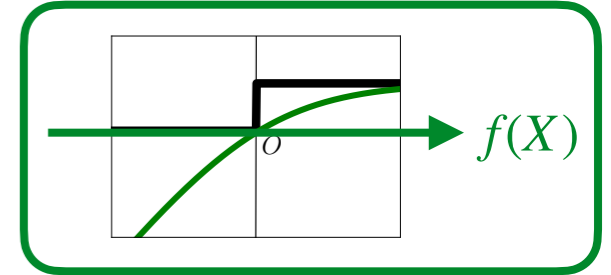
$$\geq \frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_1}$$

non-negative sum of convex  
 $\Rightarrow$  convex

denominator from above

# Surrogate Utility

- Idea: bound true utility from below

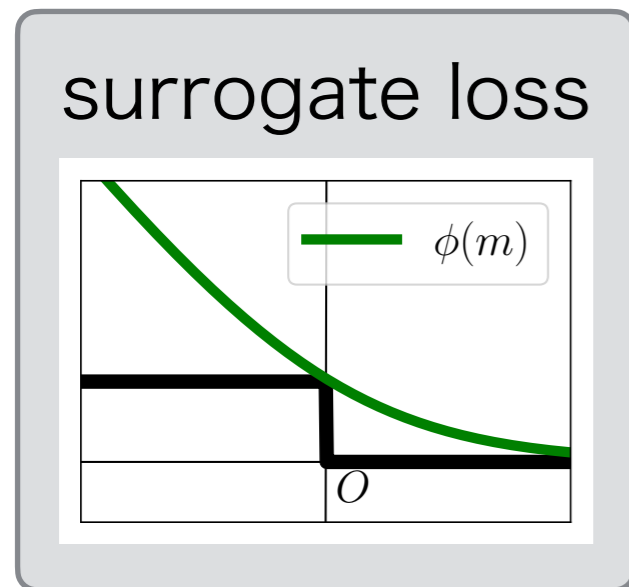


linear-fraction

$$U(f) = \frac{a_0 TP + b_0 FP + c_0}{a_1 TP + b_1 FP + c_1}$$

$$\geq \frac{a_0 \mathbb{E}_P \left[ \text{graph} \right] + b_0 \mathbb{E}_N \left[ \text{graph} \right] + c_0}{a_1 \mathbb{E}_P \left[ \text{graph} \right] + b_1 \mathbb{E}_N \left[ \text{graph} \right] + c_1}$$

||



$$U_\phi(f) = \frac{a_0 \mathbb{E}_P [1 - \phi(f(X))] + b_0 \mathbb{E}_N [-\phi(-f(X))] + c_0}{a_1 \mathbb{E}_P [1 + \phi(f(X))] + b_1 \mathbb{E}_N [\phi(-f(X))] + c_1}$$

$$:= \frac{\mathbb{E}[W_{0,\phi}]}{\mathbb{E}[W_{1,\phi}]} : \text{Surrogate Utility}$$

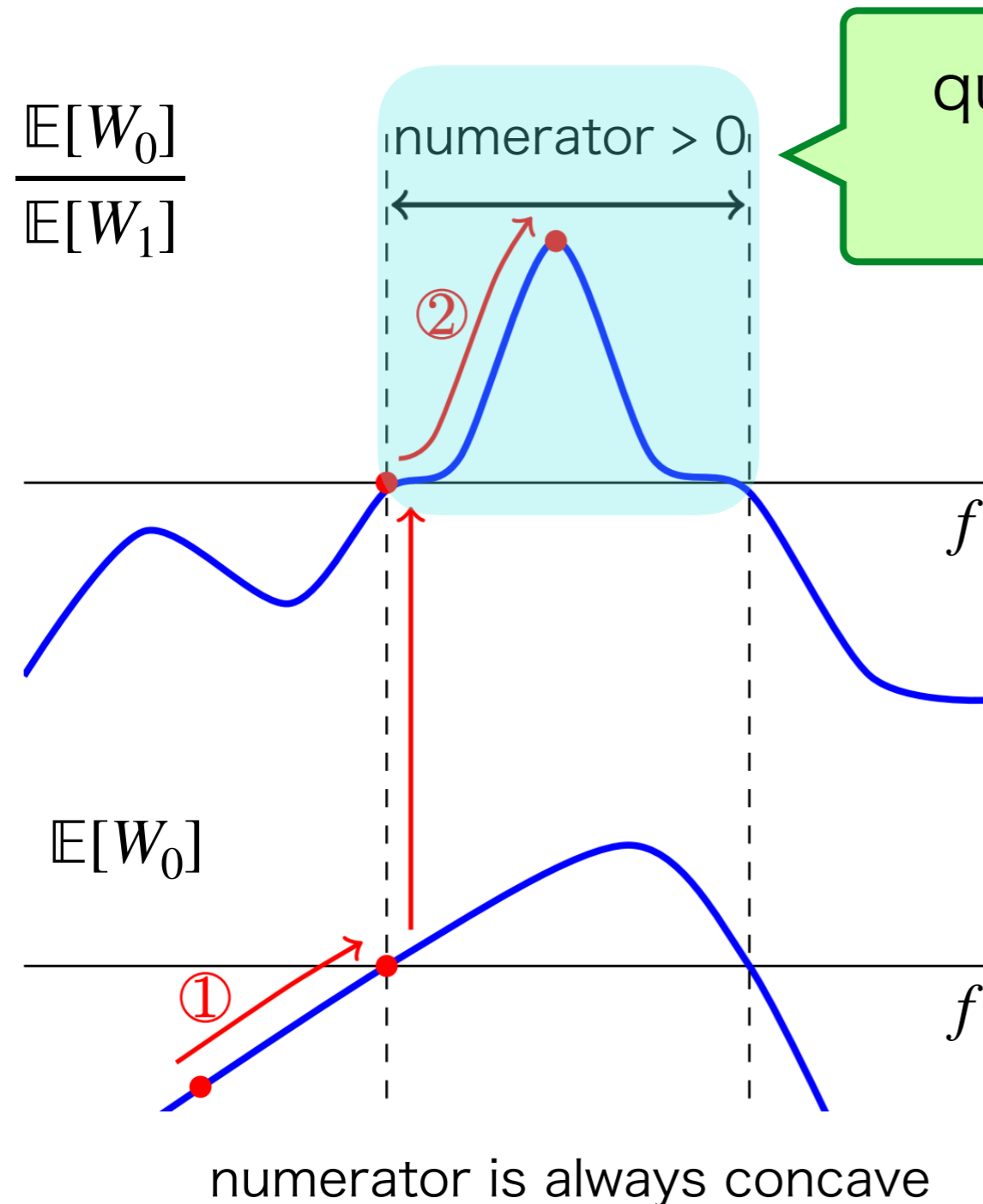
# Hybrid Optimization Strategy 31

$$U_\phi(f) = \frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_1} = \frac{\text{Concave Curve}}{\text{Convex Curve}}$$

The equation shows the utility function  $U_\phi(f)$  as a ratio of two expected values. The numerator consists of  $a_0 \mathbb{E}_P$  (with a plot of a concave curve and a step function),  $+ b_0 \mathbb{E}_N$  (with a plot of a concave curve and a step function), and  $+ c_0$ . The denominator consists of  $a_1 \mathbb{E}_P$  (with a plot of a convex curve and a step function),  $+ b_1 \mathbb{E}_N$  (with a plot of a convex curve and a step function), and  $+ c_1$ . The result is shown as a fraction of two curves: a concave curve on top and a convex curve on the bottom, both plotted against a horizontal axis.

- Note: numerator can be negative
  - ▶  $U_\phi$  isn't quasi-concave only if numerator  $< 0$
  - ▶ make numerator positive first (concave), then maximize fractional form (quasi-concave)

# Hybrid Optimization Strategy <sup>32</sup>



## Strategy

- ① update gradient-ascent direction while  $\mathbb{E}[W_0] < 0$
- ② maximize fraction by normalized-gradient ascent [Hazan+ NeurIPS2015]



# Convexity & Statistical Property

Q. How to make surrogate calibrated?

## Accuracy

tractable (convex)

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

calibrated

intractable

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

## Linear-fractional Metrics

① tractable

$$\frac{a_0 \mathbb{E}_P \cdot \left[ \begin{array}{|c|} \hline \text{graph} \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|} \hline \text{graph} \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \cdot \left[ \begin{array}{|c|} \hline \text{graph} \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|} \hline \text{graph} \\ \hline \end{array} \right] + c_1}$$

② calibrated?

intractable

$$U(f) = \frac{\mathbb{E}_X[W_0(f(X))]}{\mathbb{E}_X[W_1(f(X))]}$$

# Special Case: F<sub>1</sub>-measure

## Theorem

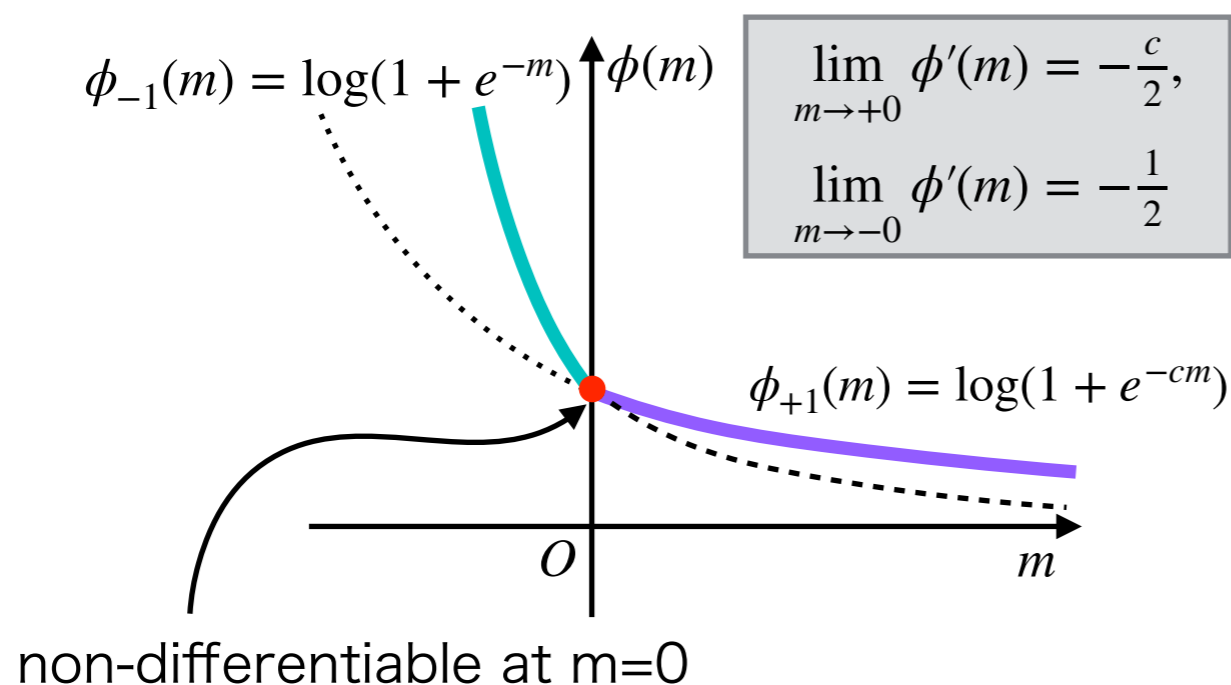
$$U_{\phi}(f_n) \xrightarrow{n \rightarrow \infty} 1 \implies U(f_n) \xrightarrow{n \rightarrow \infty} 1 \quad \forall \{f_n\}$$

if  $\phi$  satisfies

- ▶  $\exists c \in (0,1)$  s.t.  $\sup_f U_{\phi}(f) \geq \frac{2c}{1-c}$ ,  $\lim_{m \rightarrow +0} \phi'(m) \geq c \lim_{m \rightarrow -0} \phi'(m)$
- ▶  $\phi$  is non-increasing
- ▶  $\phi$  is convex

Note: informal

## Example



## Intuition:

trade off **TP** and **FP**  
by gradient steepness

# Experiment: F<sub>1</sub>-measure

(F <sub>1</sub> -measure)	Proposed		Baselines		
Dataset	U-GD	U-BFGS	ERM	W-ERM	Plug-in
adult	0.617 (101)	0.660 (11)	0.639 (51)	0.676 (18)	<b>0.681 (9)</b>
australian	<b>0.843 (41)</b>	<b>0.844 (45)</b>	0.820 (123)	0.814 (116)	0.827 (51)
breast-cancer	<b>0.963 (31)</b>	<b>0.960 (32)</b>	0.950 (37)	0.948 (44)	0.953 (40)
cod-rna	0.802 (231)	0.594 (4)	0.927 (7)	0.927 (6)	<b>0.930 (2)</b>
diabetes	<b>0.834 (32)</b>	<b>0.828 (31)</b>	0.817 (50)	0.821 (40)	0.820 (42)
fourclass	<b>0.638 (70)</b>	<b>0.638 (64)</b>	0.601 (124)	0.591 (212)	0.618 (64)
german.numer	0.561 (102)	<b>0.580 (74)</b>	0.492 (188)	0.560 (107)	<b>0.589 (73)</b>
heart	<b>0.796 (101)</b>	<b>0.802 (99)</b>	<b>0.792 (80)</b>	0.764 (151)	0.764 (137)
ionosphere	<b>0.908 (49)</b>	<b>0.901 (43)</b>	0.883 (104)	0.842 (217)	<b>0.897 (54)</b>
madelon	<b>0.666 (19)</b>	0.632 (67)	0.491 (293)	0.639 (110)	<b>0.663 (24)</b>
mushrooms	1.000 (1)	0.997 (7)	<b>1.000 (1)</b>	1.000 (2)	0.999 (4)
phishing	0.937 (29)	<b>0.943 (7)</b>	<b>0.944 (8)</b>	0.940 (12)	<b>0.944 (8)</b>
phoneme	<b>0.648 (27)</b>	0.559 (22)	0.530 (201)	0.616 (135)	0.633 (35)
skin_nonskin	0.870 (3)	0.856 (4)	0.854 (7)	<b>0.877 (8)</b>	0.838 (5)
sonar	<b>0.735 (95)</b>	<b>0.740 (91)</b>	0.706 (121)	0.655 (189)	<b>0.721 (113)</b>
spambase	0.876 (27)	0.756 (61)	0.887 (42)	0.881 (58)	<b>0.903 (18)</b>
splice	0.785 (49)	<b>0.799 (46)</b>	0.785 (55)	0.771 (67)	<b>0.801 (45)</b>
w8a	0.297 (80)	0.284 (96)	0.735 (35)	<b>0.742 (29)</b>	<b>0.745 (26)</b>

(F<sub>1</sub>-measure is shown)

$$\text{model: } f_{\theta}(x) = \theta^{\top} x$$

$$\text{surrogate loss: } \phi(m) = \max\{\log(1 + e^{-m}), \log(1 + e^{-\frac{m}{3}})\}$$

# Loss for Complicated Metrics

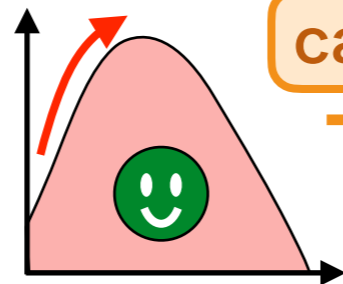
## Linear-fractional metrics

contains F-measure, Jaccard  
often used with imbalanced data

$$U(f) = \frac{a_0 TP + b_0 FP + c_0}{a_1 TP + b_1 FP + c_1}$$

### surrogate utility

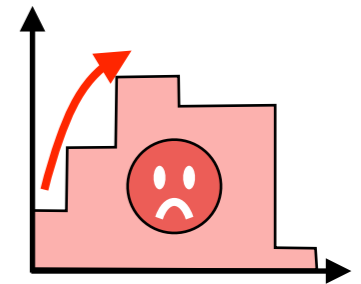
$$\frac{a_0 \mathbb{E}_P + c_0}{a_1 \mathbb{E}_P + c_1} + \frac{b_0 \mathbb{E}_N}{b_1 \mathbb{E}_N + c_1}$$



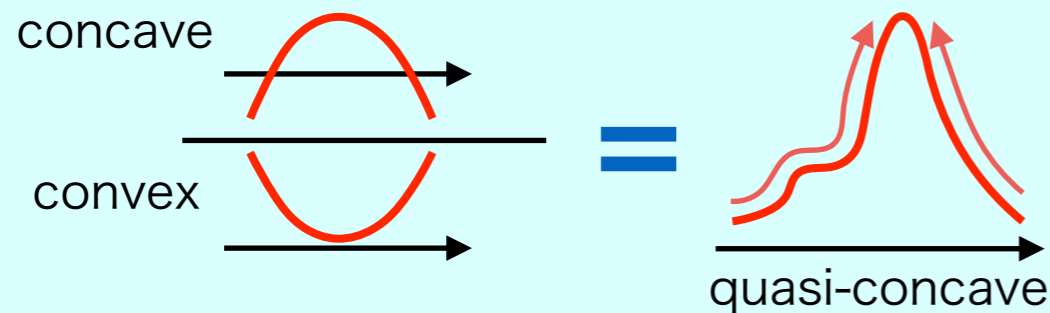
calibrated

### target utility

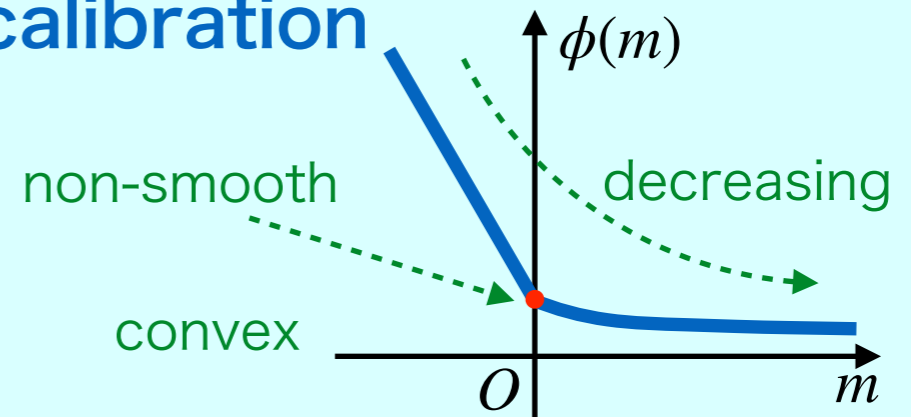
$$U(f) = \frac{a_0 TP + b_0 FP + c_0}{a_1 TP + b_1 FP + c_1}$$



### ① tractability (quasi-concave)



### ② calibration



Provides guideline of designing loss for complicated metrics!

# When adversary presents

**H. Bao**, C. Scott, and M. Sugiyama.

Calibrated Surrogate Losses for Adversarially Robust Classification.

In *COLT*, 2020.

# Adversarial Attacks

[Goodfellow+ 2015]

**Adding imperceptible small noise can fool classifiers!**

original data



$\mathbf{x}$

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$

“nematode”

8.2% confidence

=

perturbed data



$\mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$

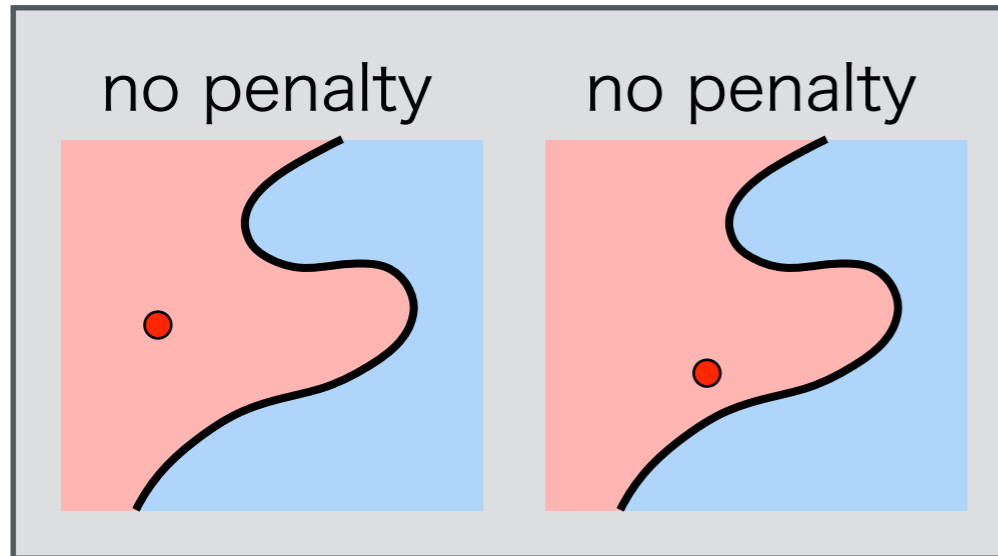
“gibbon”

99.3 % confidence



# Penalize Vulnerable Prediction

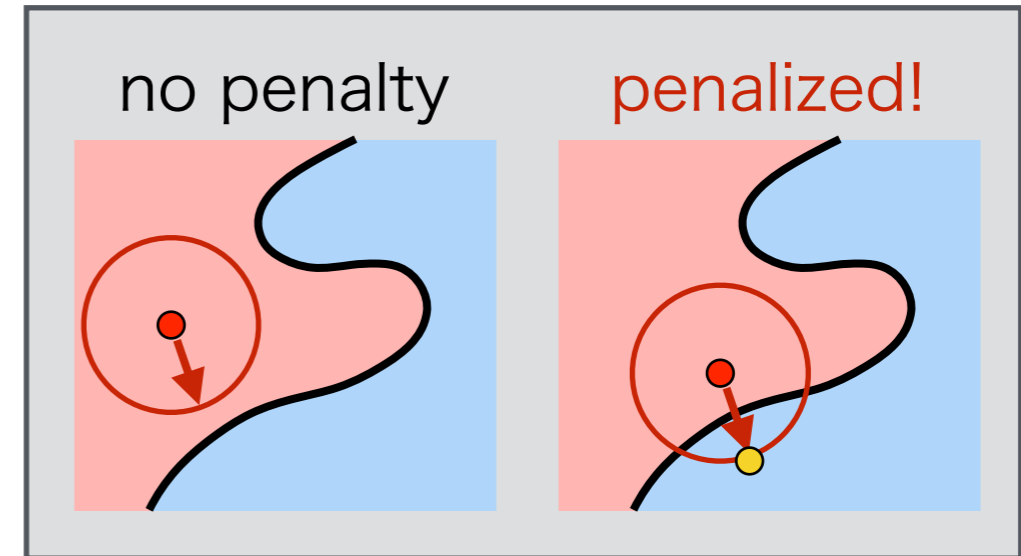
## Usual Classification



usual 0-1 loss

$$\ell_{01}(x, y, f) = \begin{cases} 1 & \text{if } yf(x) \leq 0 \\ 0 & \text{otherwise} \end{cases}$$

## Robust Classification



robust 0-1 loss

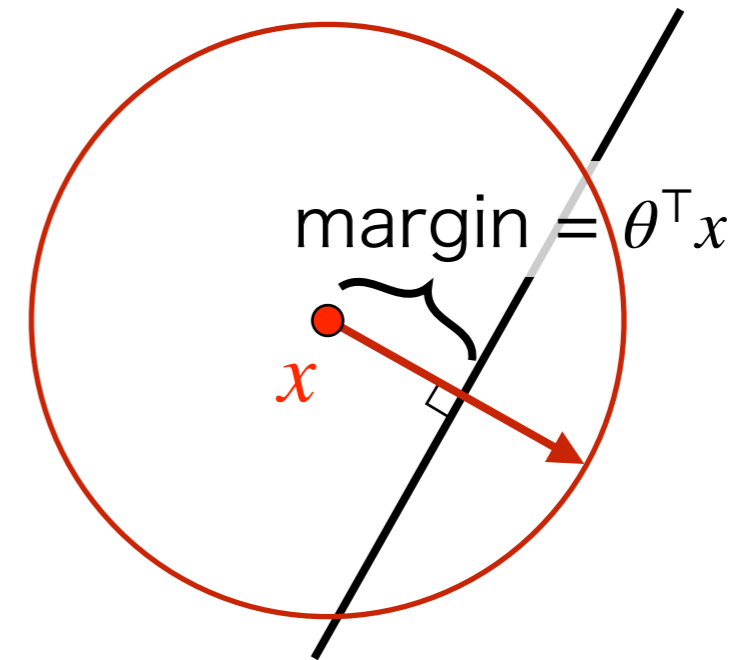
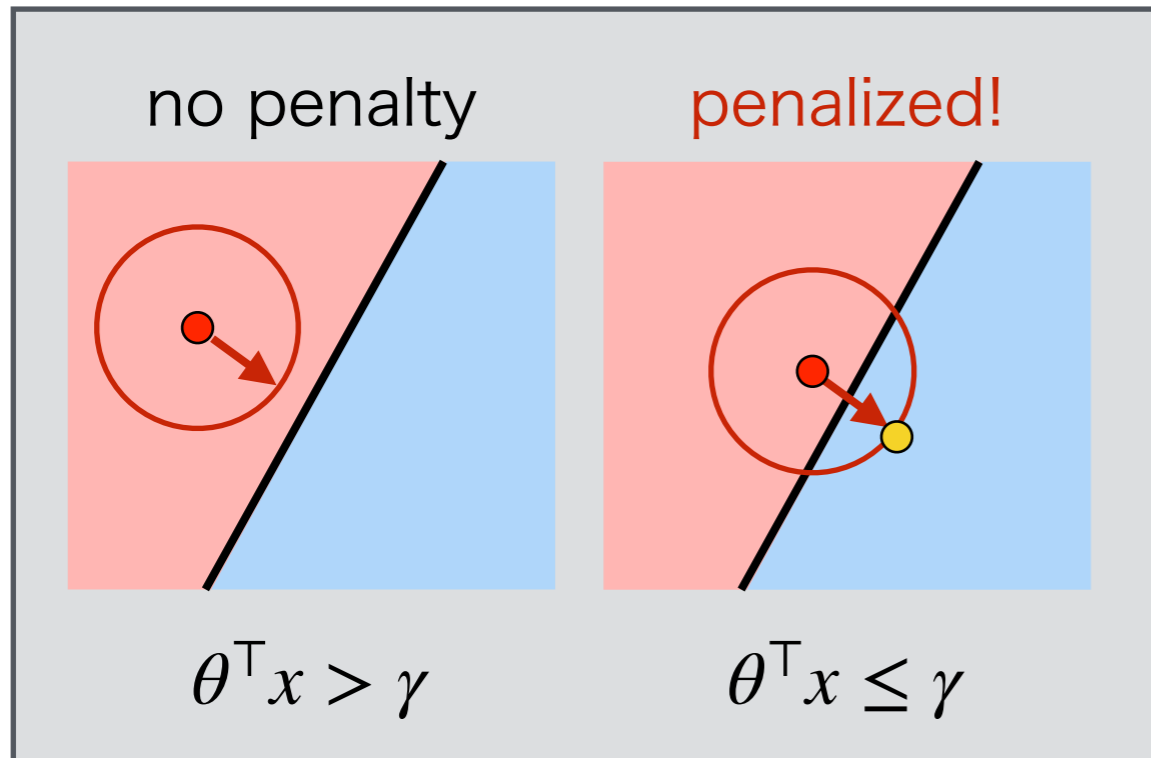
$$\ell_{\gamma}(x, y, f) = \begin{cases} 1 & \text{if } \exists \Delta \in \mathbb{B}_2(\gamma) . yf(x + \Delta) \leq 0 \\ 0 & \text{otherwise} \end{cases}$$

prediction too close to boundary  
should be penalized

$$\mathbb{B}_2(\gamma) = \{x \in \mathbb{R}^d \mid \|x\|_2 \leq \gamma\}: \gamma\text{-ball}$$

# In Case of Linear Predictors

linear predictors  $\mathcal{F}_{\text{lin}} = \{x \mapsto \theta^\top x \mid \|\theta\|_2 = 1\}$



**robust 0-1 loss**

$$\ell_\gamma(x, y, f) = \begin{cases} 1 & \text{if } \exists \Delta \in \mathbb{B}_2(\gamma) . yf(x + \Delta) \leq 0 \\ 0 & \text{otherwise} \end{cases} = \mathbf{1}\{yf(x) \leq \gamma\} := \phi_\gamma(yf(x))$$



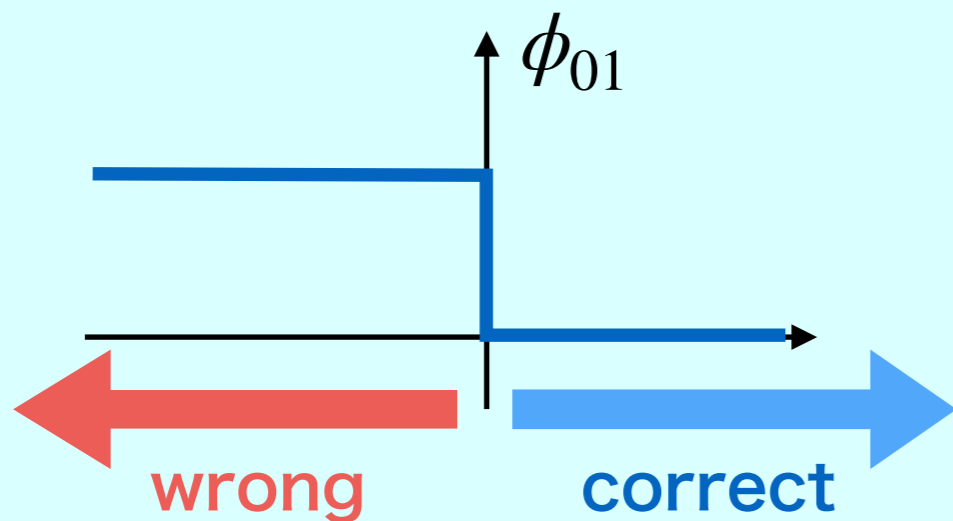
# Formulation of Classification <sup>41</sup>

## Usual Classification

minimize 0-1 risk

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

0-1 loss  $\phi_{01}(\alpha) = \mathbf{1}\{\alpha \leq 0\}$



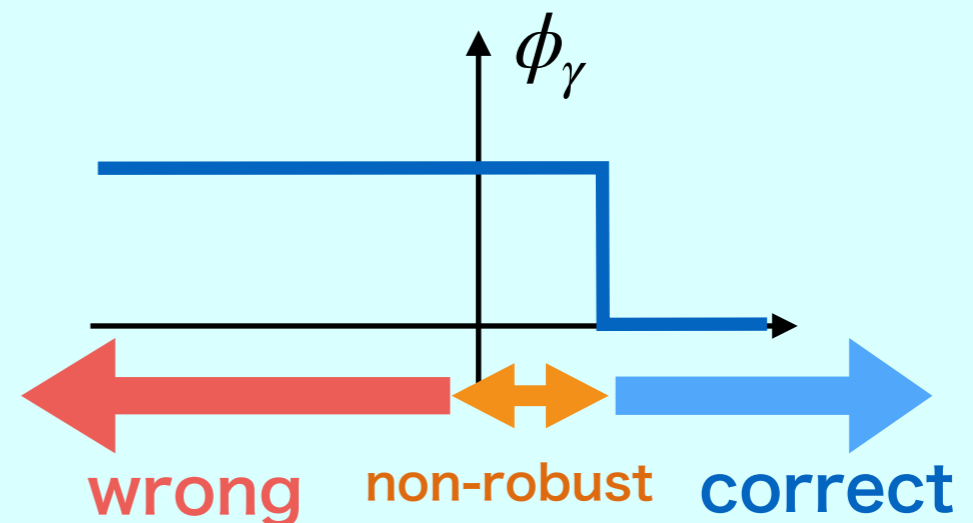
## Robust Classification

minimize  $\gamma$ -robust 0-1 risk

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

(restricted to linear predictors)

robust 0-1 loss  $\phi_{\gamma}(\alpha) = \mathbf{1}\{\alpha \leq \gamma\}$

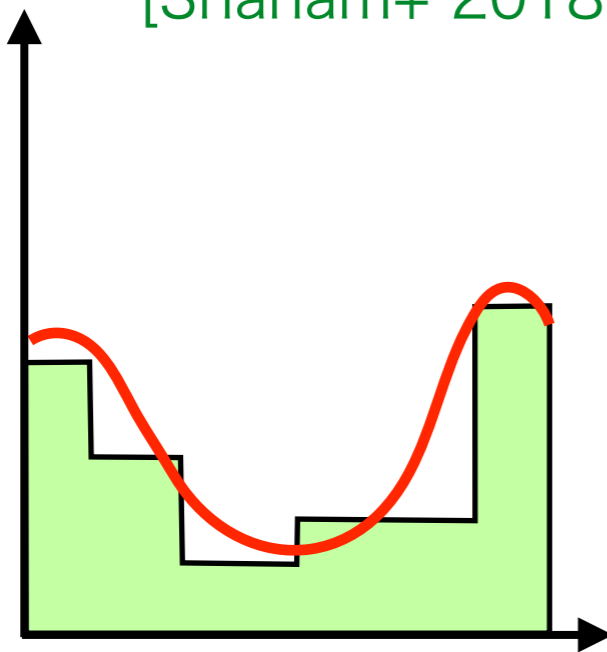


# Existing Approaches

Direct optimization of robust risk  $R_{\phi_\gamma}(f)$  is intractable

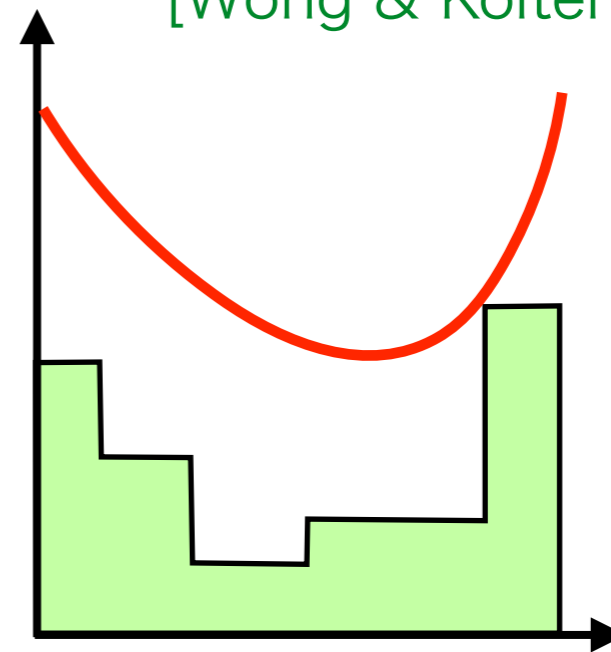
## Taylor approximation

[Shaham+ 2018; etc.]



## Convex upper bound

[Wong & Kolter 2018; etc.]



Both do not necessarily lead to true minimizer!

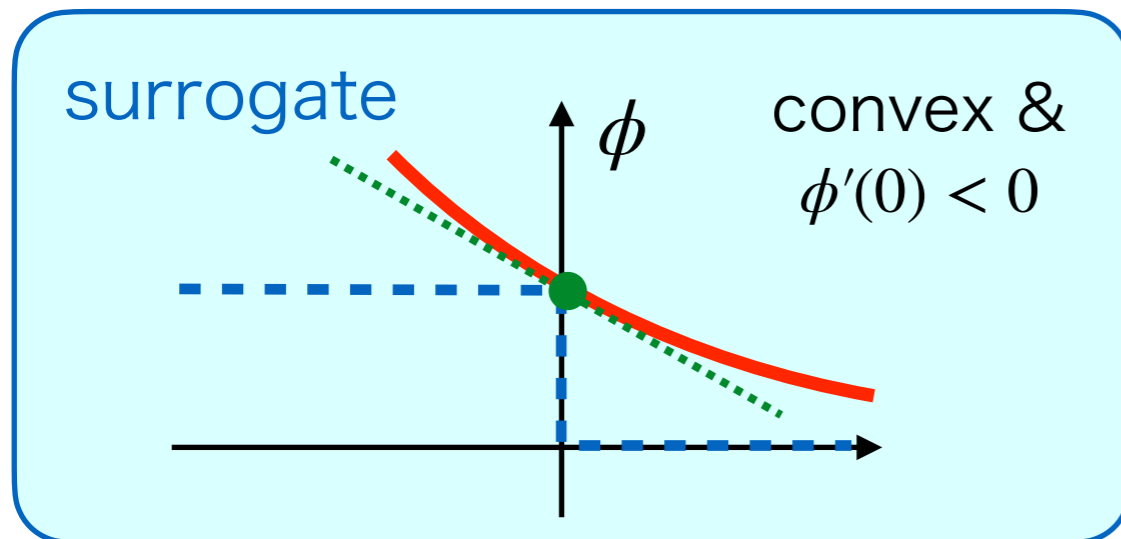
Shaham, U., Yamada, Y., & Negahban, S. (2018).

Understanding adversarial training: Increasing local stability of supervised models through robust optimization. *Neurocomputing*, 195-204.

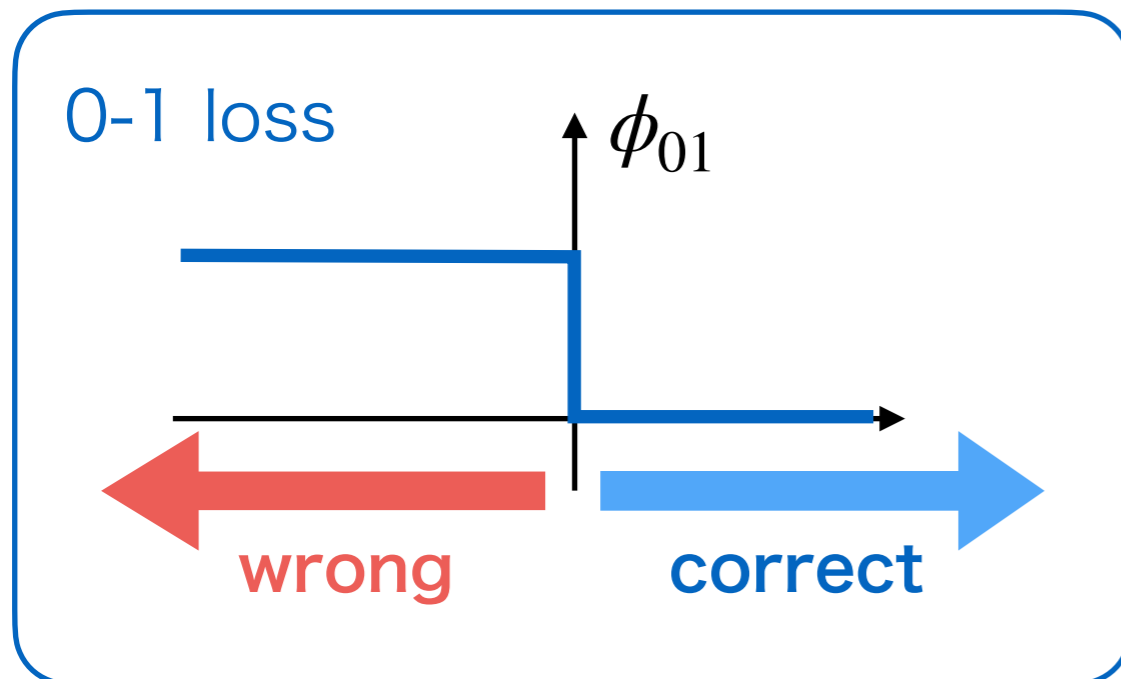
Wong, E., & Kolter, Z. (2018,). Provable Defenses against Adversarial Examples via the Convex Outer Adversarial Polytope. In *International Conference on Machine Learning* (pp. 5286-5295).

# What surrogate is calibrated? <sup>43</sup>

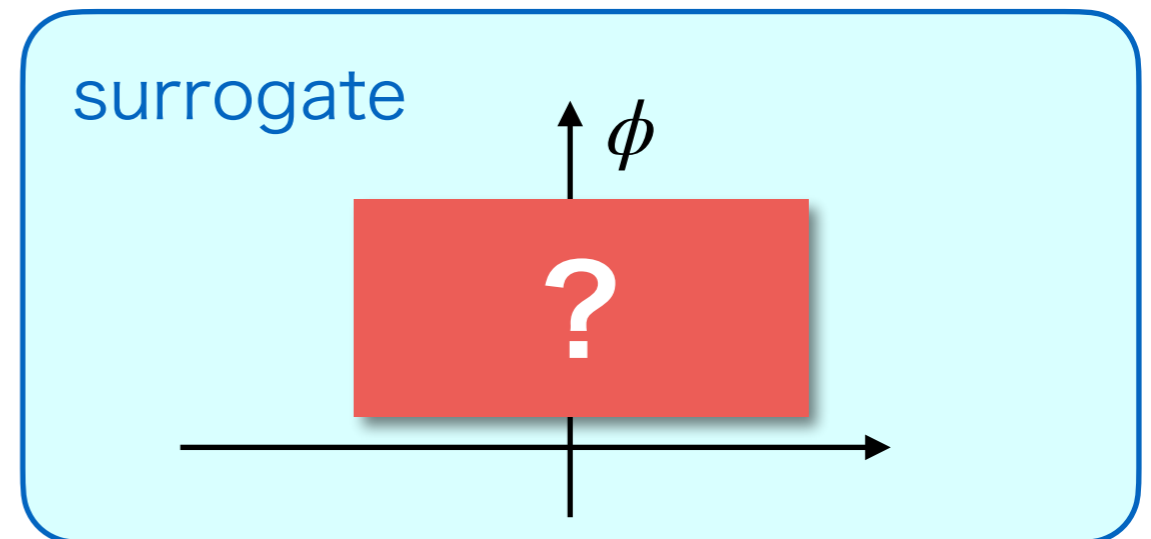
## Usual Classification



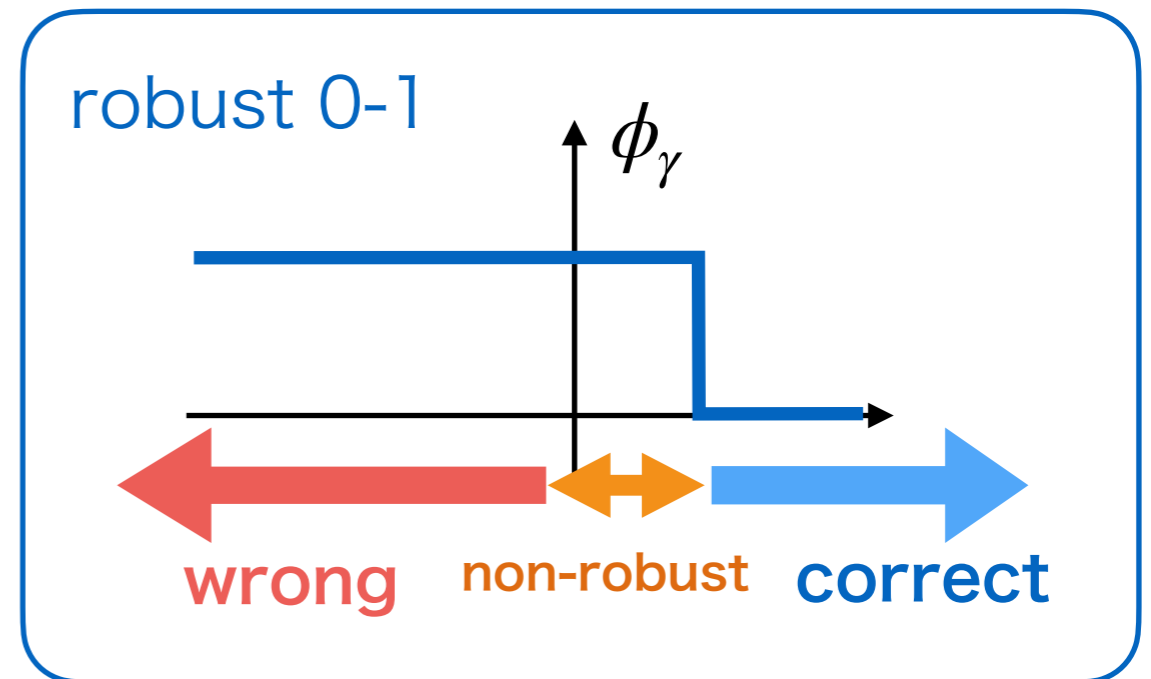
calibrated  
[Bartlett+ 2006]



## Robust Classification



calibrated

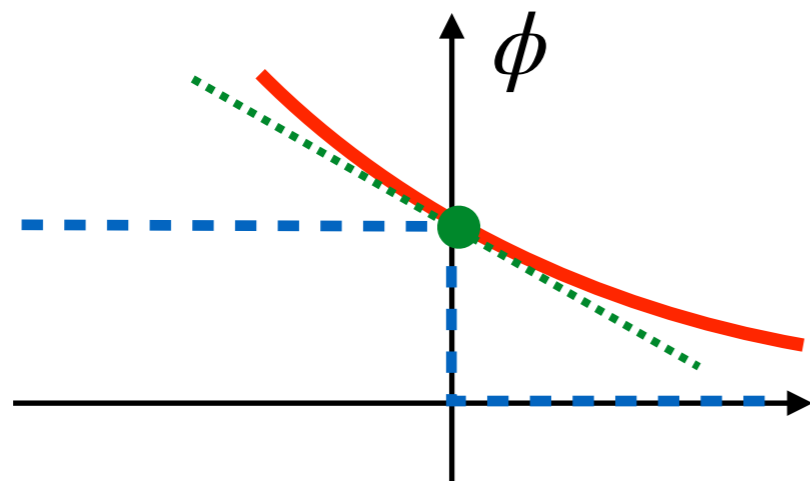


# Isn't it a piece of cake?

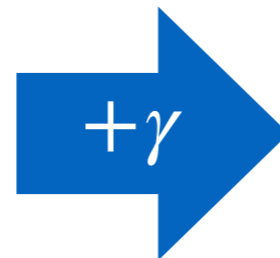
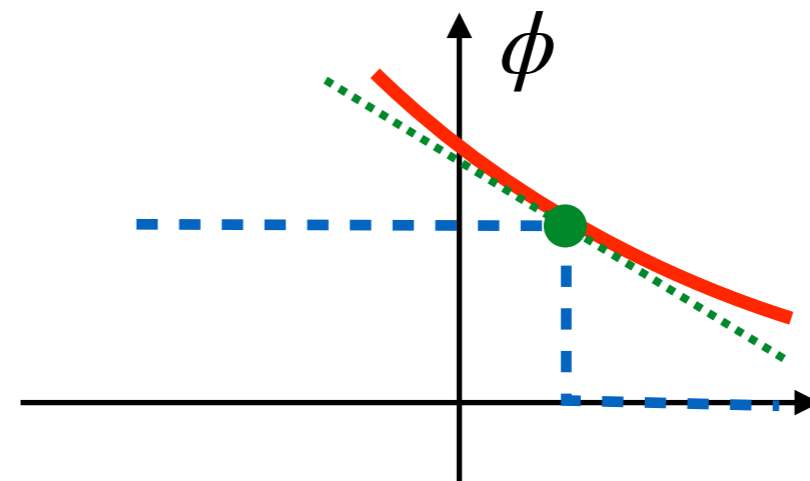
**Theorem.** If surrogate  $\phi$  is convex, it is  $\phi_{01}$ -calibrated iff

- ▶ differentiable at 0
- ▶  $\phi'(0) < 0$

Usual 0-1 loss



Robust 0-1 loss



If  $\phi'(\gamma) < 0$ , then calibrated to robust 0-1 loss?

# No convex calibrated surrogate

**Theorem.** Any convex surrogate is not  $\phi_\gamma$ -calibrated.

(under linear predictors)

## Proof Sketch

**Idea:** to show  $\delta(\varepsilon) = 0$  for some  $\varepsilon > 0$

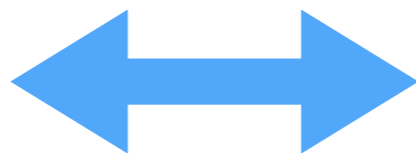
calibration function

convex in  $f$

$f$  is non-robust ( $|f(x)| < \gamma$ )

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

$$\delta(\varepsilon) = 0$$



$$\inf_{f: R_{\phi_\gamma}(f) - R_{\phi_\gamma}^* \geq \varepsilon} R_\phi(f) = \inf_f R_\phi(f)$$

“non-robust”  
minimizer  
( $|f(x)| < \gamma$ )

optimal  
minimizer

# No convex calibrated surrogate

**Theorem.** Any convex surrogate is not  $\phi_\gamma$ -calibrated.

(under linear predictors)

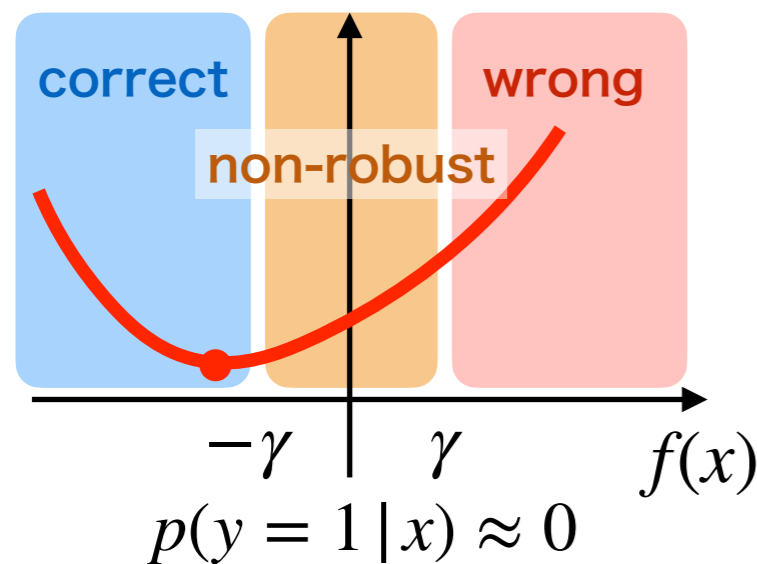
## Proof Sketch

calibration function

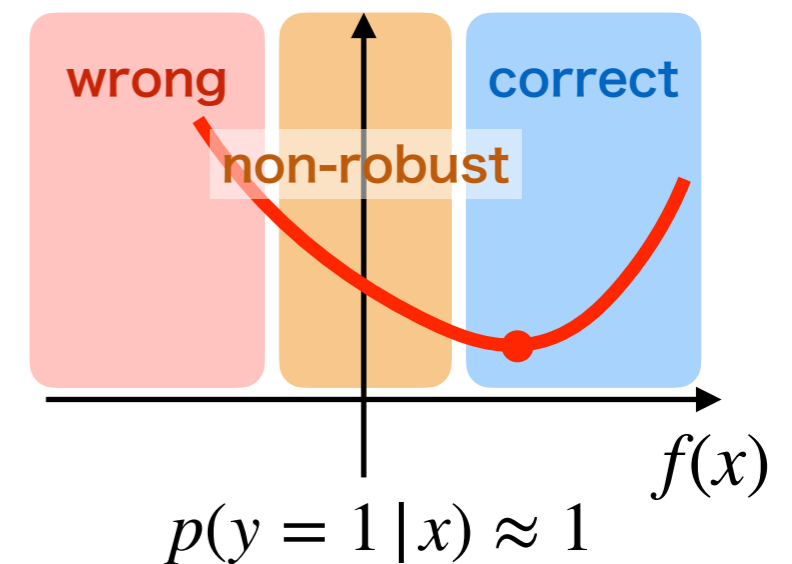
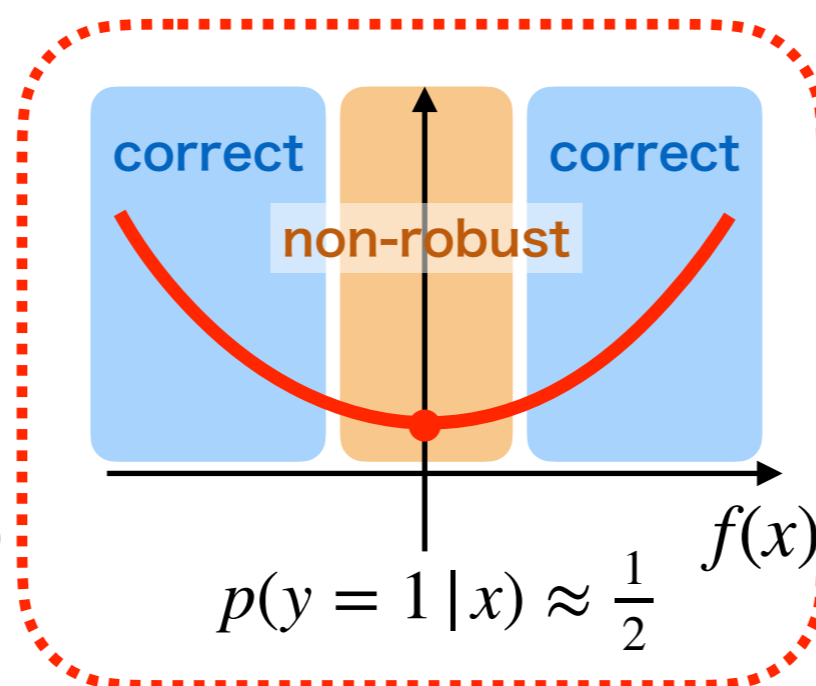
convex in  $f$

$f$  is non-robust ( $|f(x)| < \gamma$ )

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



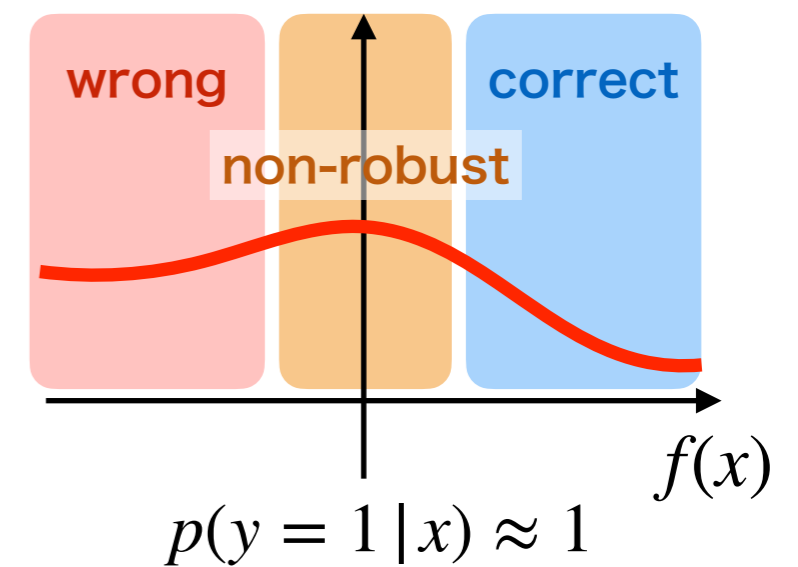
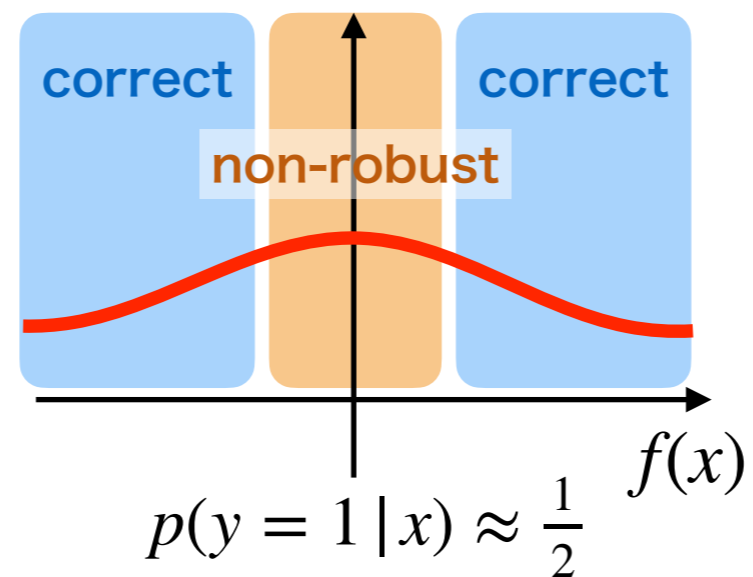
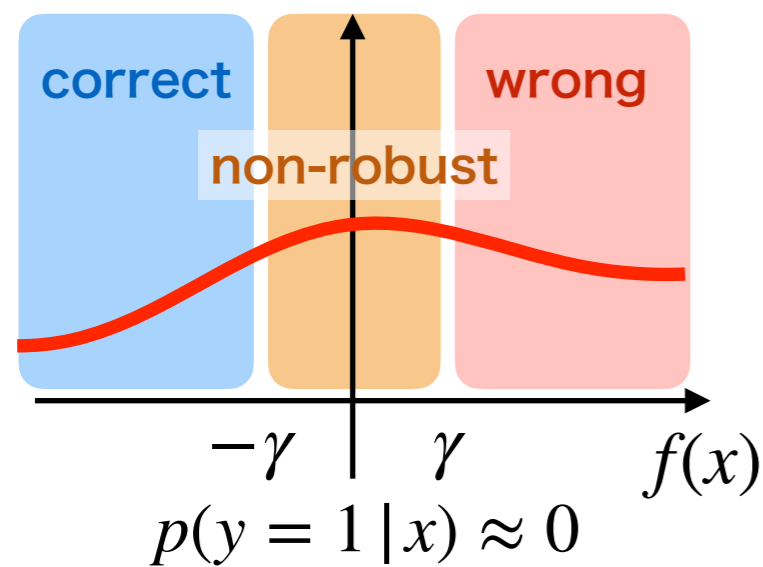
non-robust minimizer!



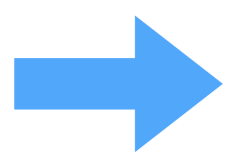
surrogate **conditional** risk is plotted

# How to find calibrated surrogate? <sup>47</sup>

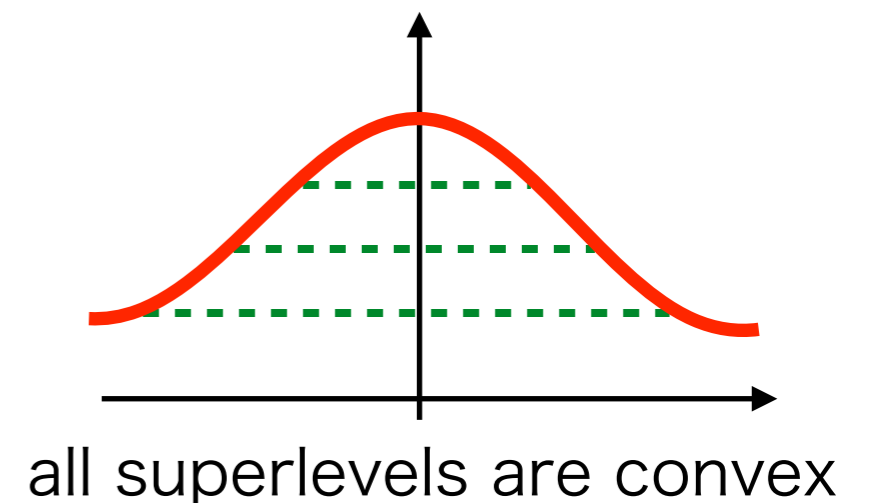
**Idea.** To make conditional risk not minimized in **non-robust area**



surrogate conditional risk is plotted



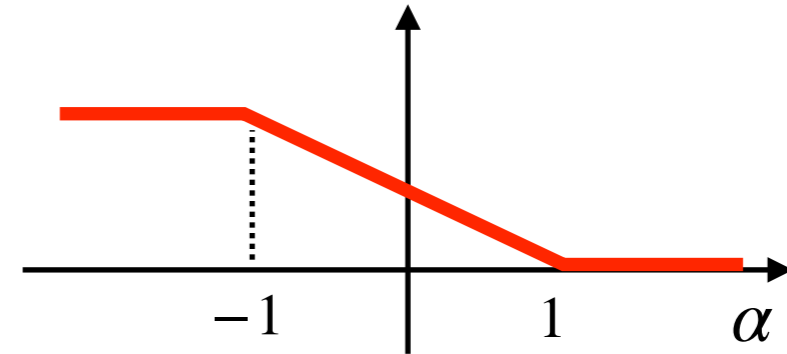
consider a surrogate  $\phi$  such that conditional risk is **quasiconcave**



# Example: Shifted Ramp Loss 48

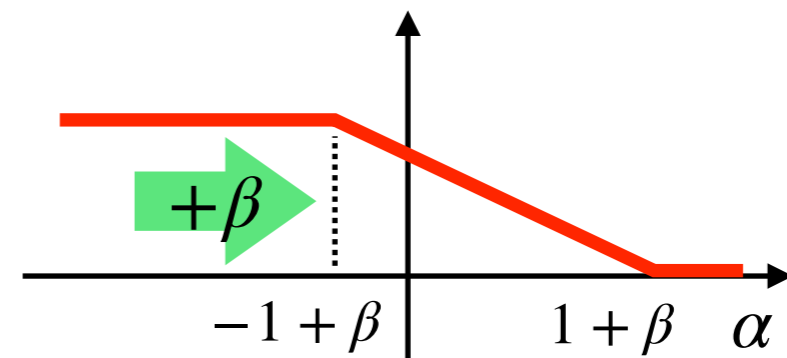
Ramp loss

$$\phi(\alpha) = \text{clip}_{[0,1]} \left( \frac{1 - \alpha}{2} \right)$$

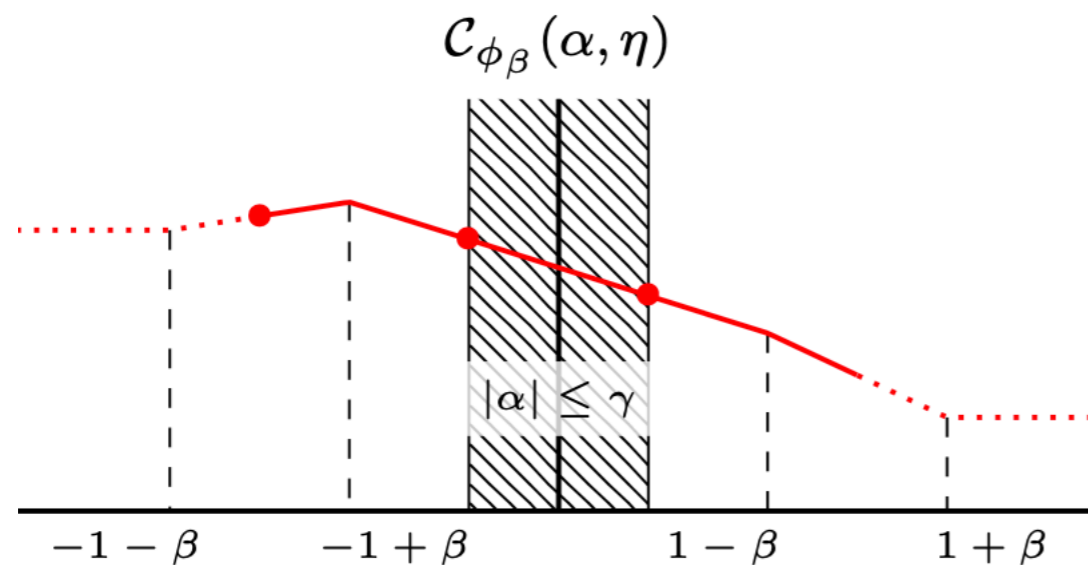


Shifted ramp loss

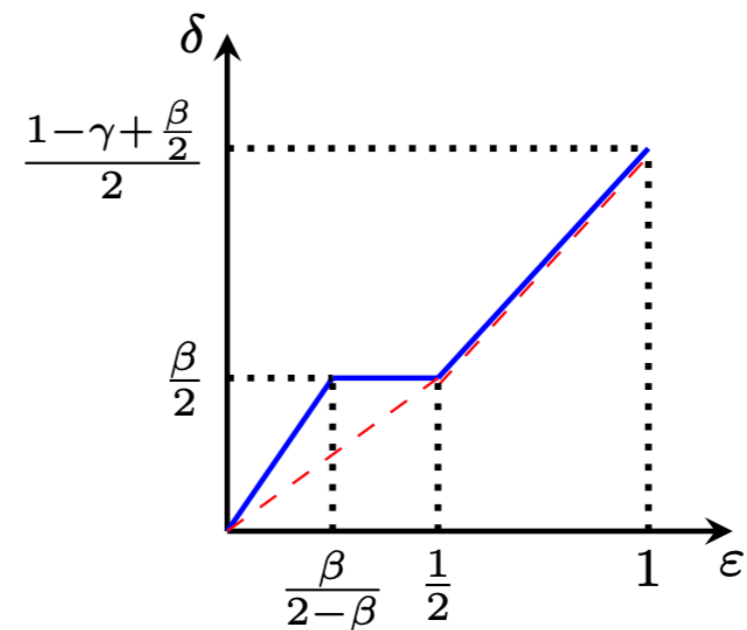
$$\phi_\beta(\alpha) = \text{clip}_{[0,1]} \left( \frac{1 - \alpha + \beta}{2} \right)$$



conditional risk ( $p(y = 1 | x) > \frac{1}{2}$ )



calibration function

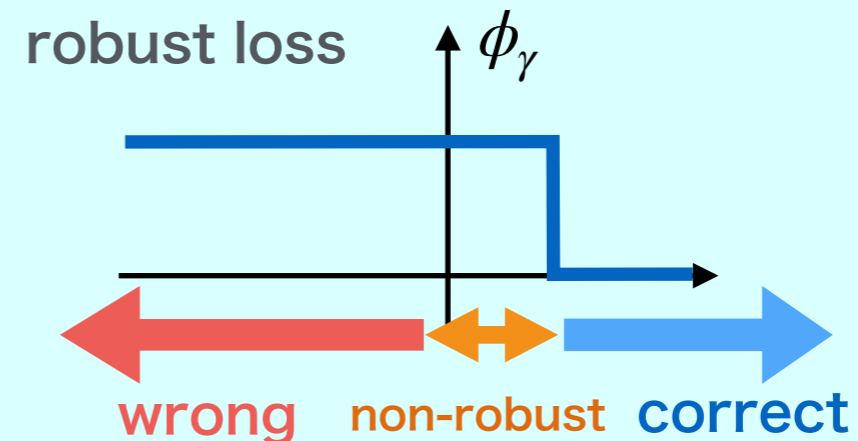
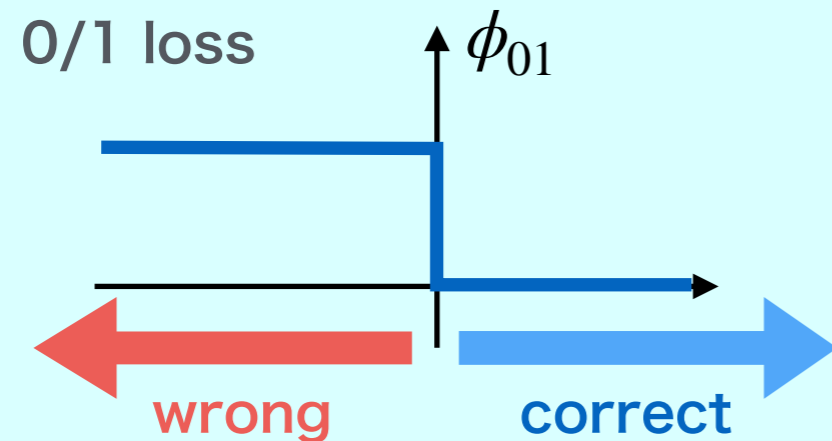


assume  $0 < \beta < 1 - \gamma$



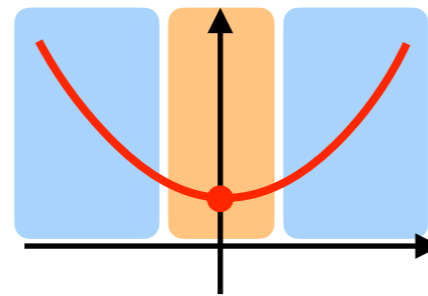
# Loss for Robust Learning

“Embed” robustness into loss function

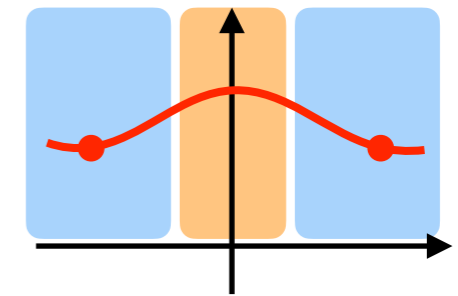


loss function can not only accommodate classification performance but also robustness

Inability of convex loss



convex loss is minimized in non-robust area

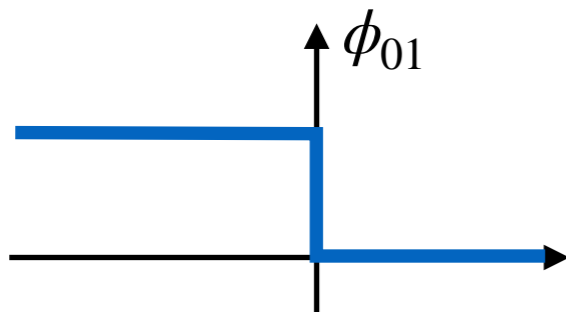
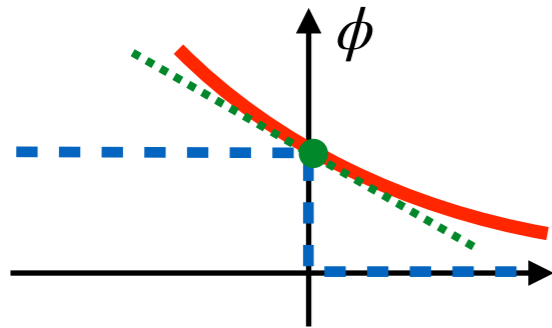


robust objective

Calibration theory helps to reveal classifiers' property!

# Summary

## Binary Classification



## Linear-fractional Metrics

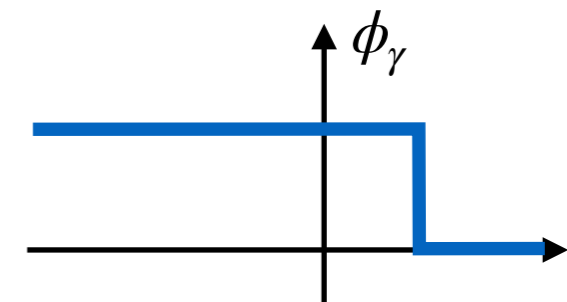
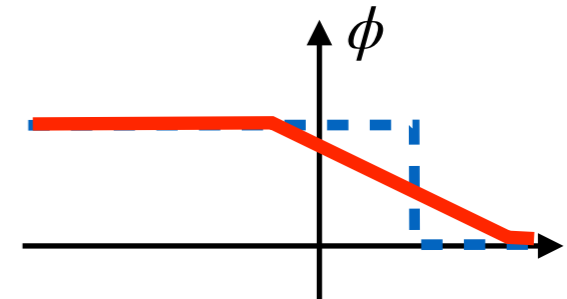
$$\frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + c_1}$$



$$\frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \text{step} \\ \hline \end{array} \right] + c_1}$$

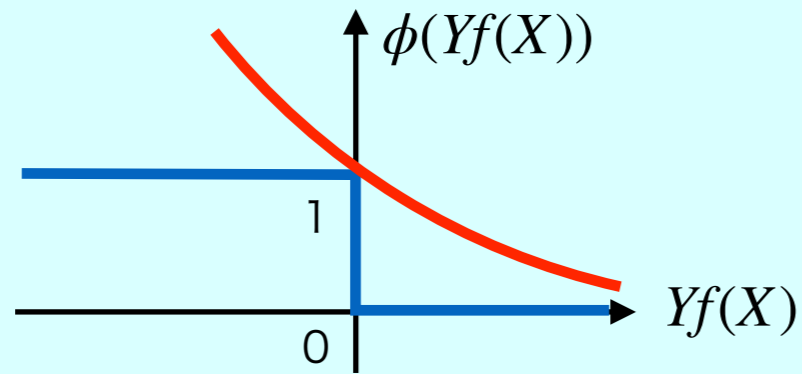
quasi-concave surrogate utility

## Adversarial Robustness

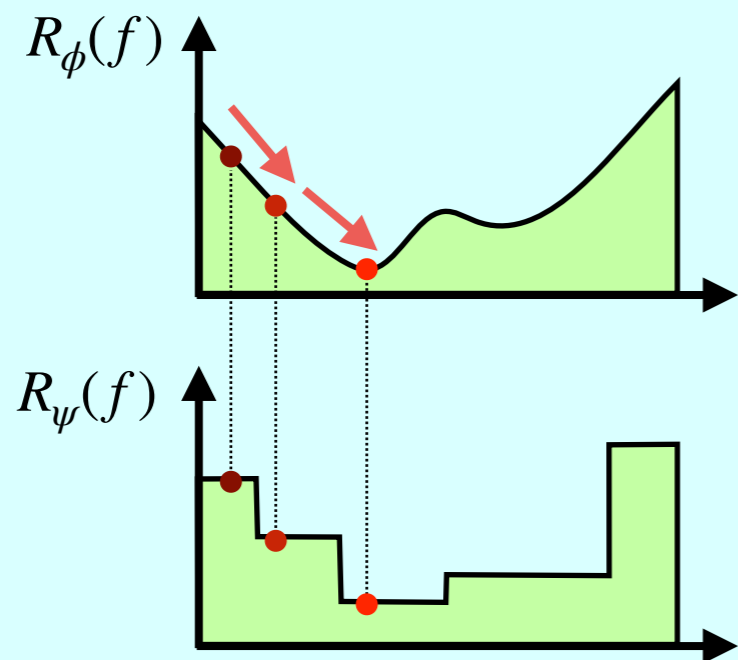


no convex & calibrated surrogate

## Surrogate vs. Target loss



## Calibrated loss



Binary Classification

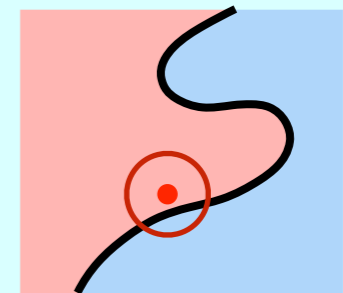
Multi-class Classification

Linear-fractional

$$\frac{a_0 \text{TP} + b_0 \text{FP} + c_0}{a_1 \text{TP} + b_1 \text{FP} + c_1}$$

more complicated metrics

Adversarial Robustness



not only  
classification performance!

More applications?

noise/distribution  
robustness

hypothesis testing

metric elicitation

representation  
learning