

# Calibrated Surrogate Maximization of Linear-fractional Utility

07<sup>th</sup> Feb.

Han Bao (The University of Tokyo / RIKEN AIP)



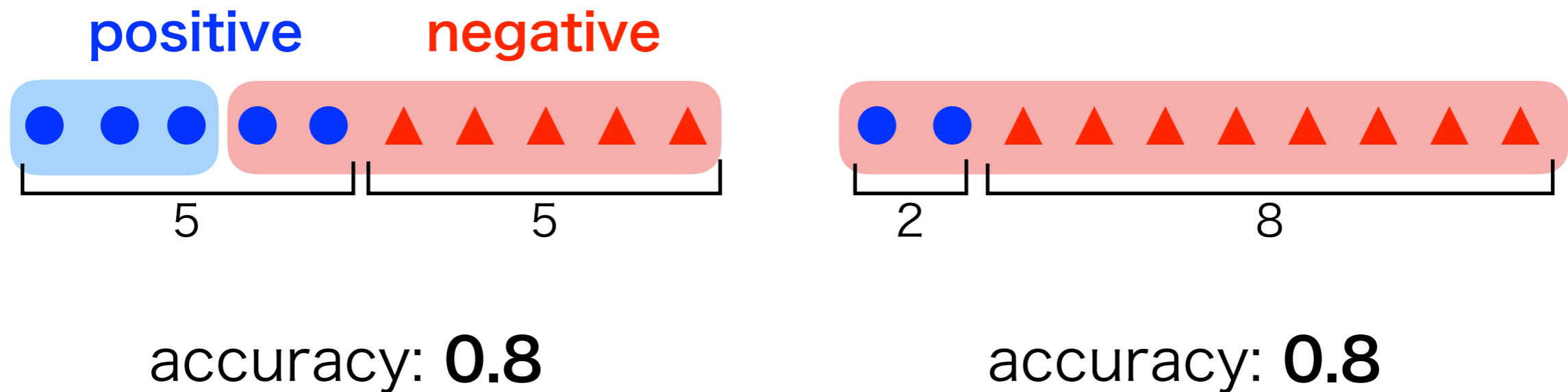
東京大学  
THE UNIVERSITY OF TOKYO



RIKEN

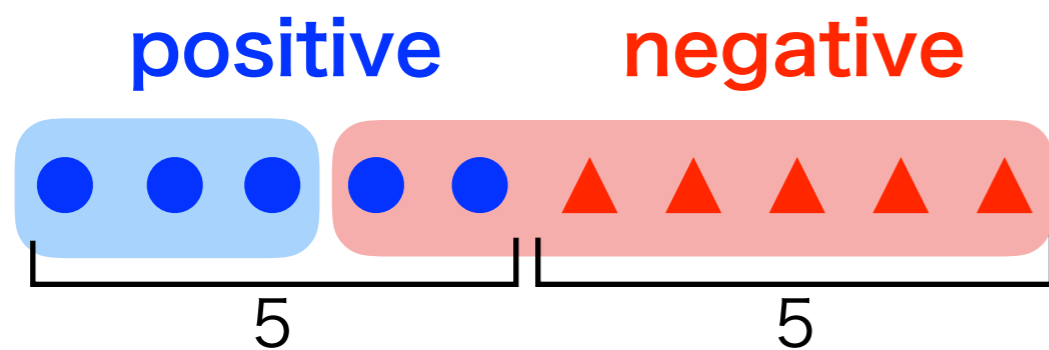
# Is accuracy appropriate?

- Our focus: **binary classification**



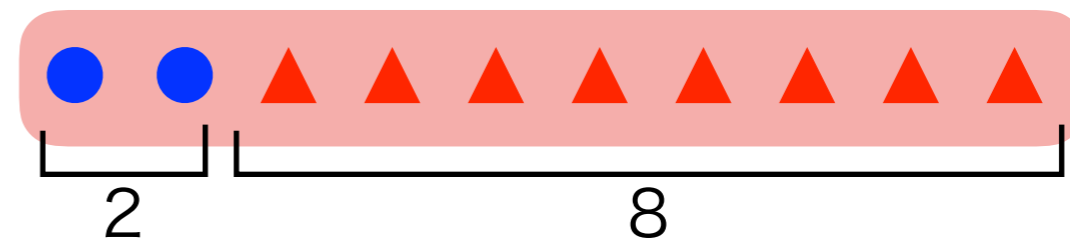
**May cause severe issues!**  
**(e.g. in medical diagnosis)**

# Is accuracy appropriate?



accuracy: **0.8**

F-measure: **0.75**



accuracy: **0.8**

F-measure: **0**

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

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

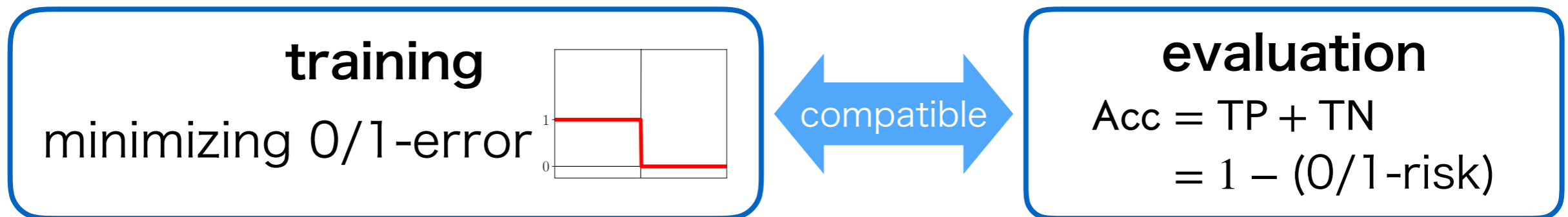
$$TN = \mathbb{E}_{X, Y=-1} [1_{\{f(X) < 0\}}]$$

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

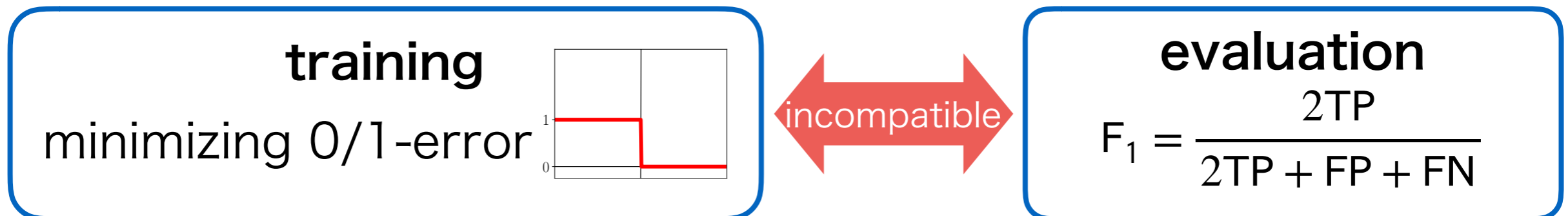
$$FN = \mathbb{E}_{X, Y=+1} [1_{\{f(X) < 0\}}]$$

# Training and Evaluation

- Usual empirical risk minimization (ERM)

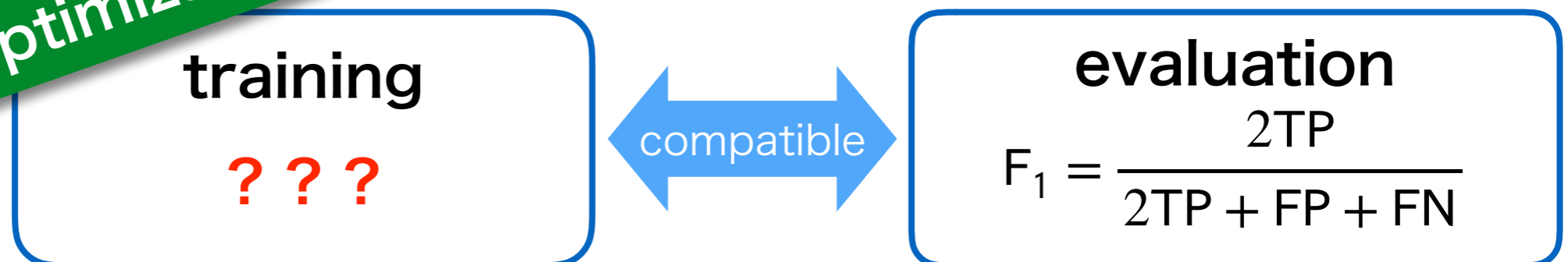


- Training with accuracy but evaluate with  $F_1$



- Why not?

**Direct Optimization**



Fowlkes-Mallows index

$$\text{FMI} = \frac{\text{TP}}{\pi} \sqrt{\frac{1}{\text{TP} + \text{FP}}}$$

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

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

Accuracy

Balanced Error Rate

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

Jaccard index

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

Matthews Correlation Coefficient

$$\text{MCC} = \frac{\text{TP} \cdot \text{TN} - \text{FP} \cdot \text{FN}}{\sqrt{\pi(1 - \pi)(\text{TP} + \text{FP})(\text{TN} + \text{FN})}}$$

Gower-Legendre index

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

Wanna Unify!!

# 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

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

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

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

- TP, FP = expectation of 0/1-loss

▶ e.g.  $\text{TP} = \mathbb{P}(Y = +1, f(X) > 0) = \mathbb{E}_{X, Y=+1}[1_{\{f(X) > 0\}}]$

# 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.  $U(f) = \sup_{f'} U(f')$



# Outline

- Introduction
- **Preliminary**
  - ▶ Convex Risk Minimization
  - ▶ Plug-in Principle vs. Cost-sensitive Learning
- Key Idea
  - ▶ Quasi-concave Surrogate
- Calibration Analysis & Experiments

# Formulation of Classification <sup>10</sup>

- Goal of classification: maximize accuracy  
= minimize mis-classification rate

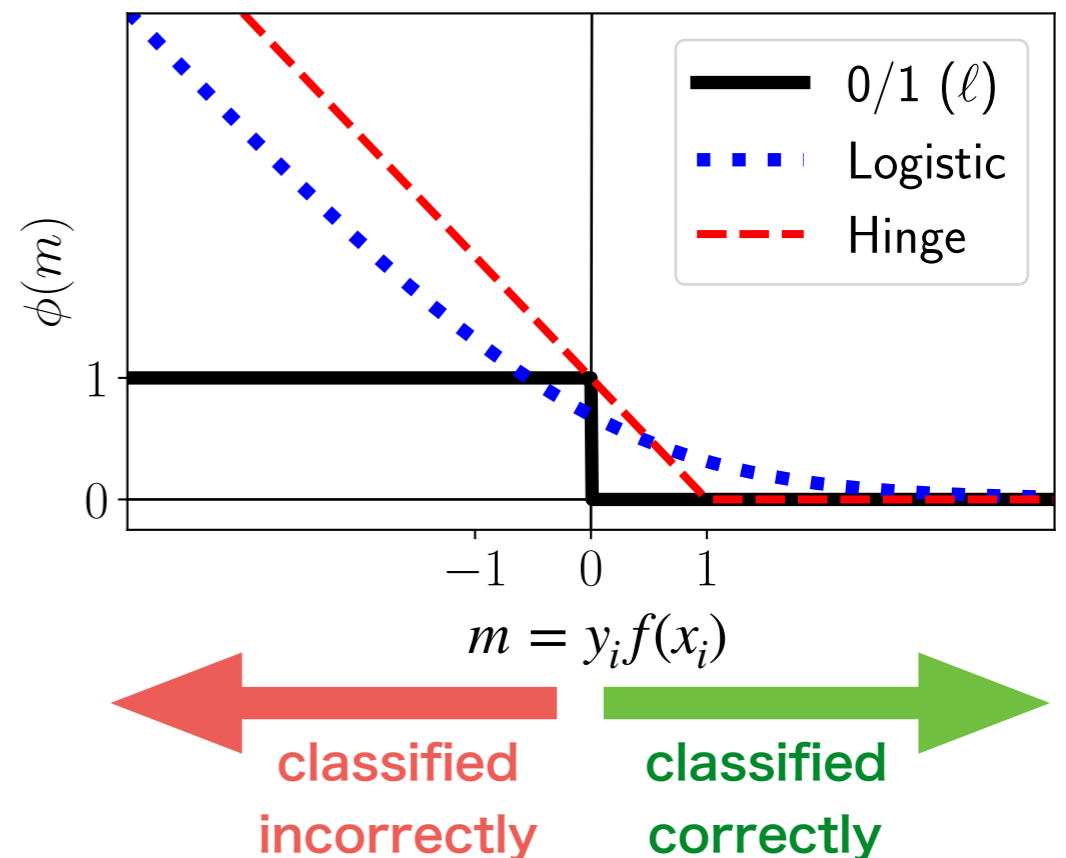
$$\begin{aligned}\hat{R}(f) &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}[y_i \neq \text{sign}(f(x_i))] \\ &= \frac{1}{n} \sum_{i=1}^n \ell(y_i f(x_i))\end{aligned}$$

↓ make 0/1 loss smoother

(Empirical) Surrogate Risk

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

**convex in  $f$ !**



Example of  $\phi$

- ▶ logistic loss
- ▶ hinge loss  $\Rightarrow$  SVM
- ▶ exponential loss  $\Rightarrow$  AdaBoost

# 3 Actors in Risk Minimization <sup>11</sup>

- Minimize classification risk (= 1 - Accuracy)

$$R(f) = \mathbb{E} [ \underbrace{\ell}_{0/1\text{-loss}} ( \underbrace{Yf(X)}_{\text{prediction margin}} ) ]$$

0/1-loss represents if  $X$  is correctly classified by  $f$

- Surrogate loss makes tractable  
differentiable upper bound of 0/1-loss  
(surrogate risk)

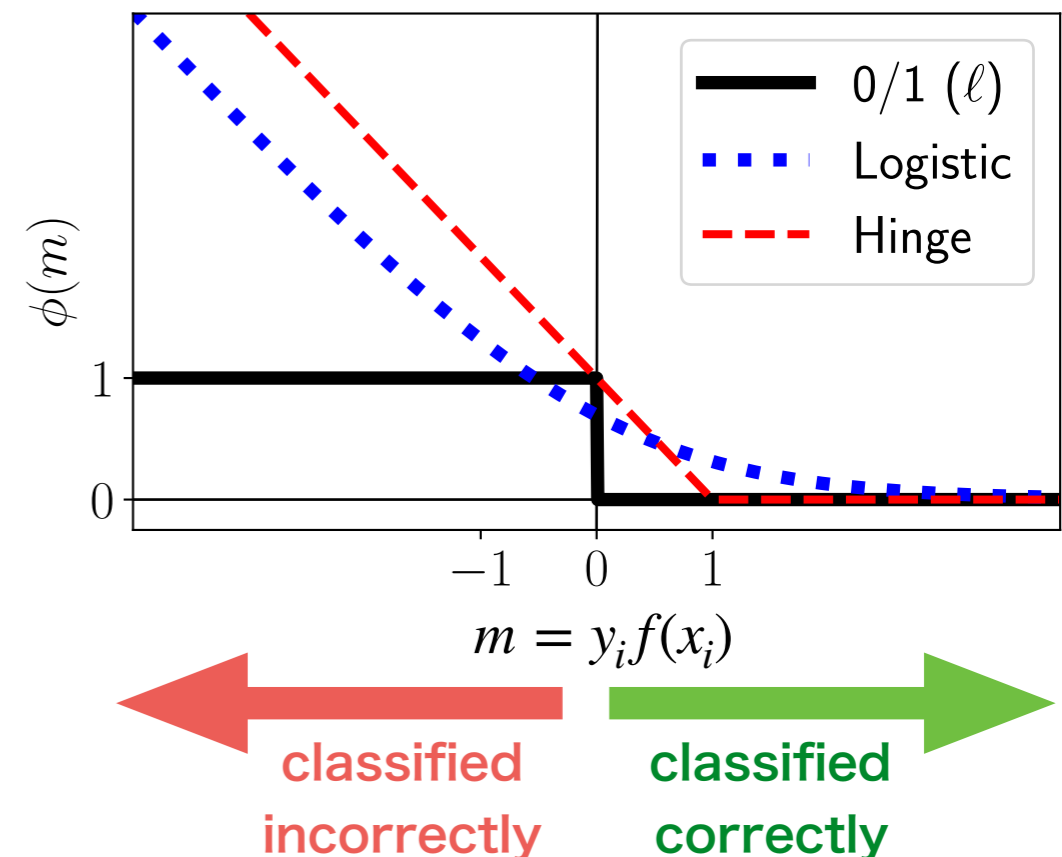
$$R_\phi(f) = \mathbb{E} [ \underbrace{\phi}_{\text{surrogate loss}} ( Yf(X) ) ]$$

- Sample approximation (M-estimation)

(empirical (surrogate) risk)

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

what we actually minimize



# Convexity & Statistical Property <sup>12</sup>

tractable (convex)

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

generalize

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

$$R(f) = \mathbb{E}[\ell(Yf(X))]$$

intractable

Q.  $\operatorname{argmin} R_\phi = \operatorname{argmin} R$  ?

A. Yes, w/ calibrated surrogate

**Theorem.** [Bartlett+ 2006]

Assume  $\phi$ : convex.

Then,  $\operatorname{argmin}_f R_\phi(f) = \operatorname{argmin}_f R(f)$   
iff  $\phi'(0) < 0$ .

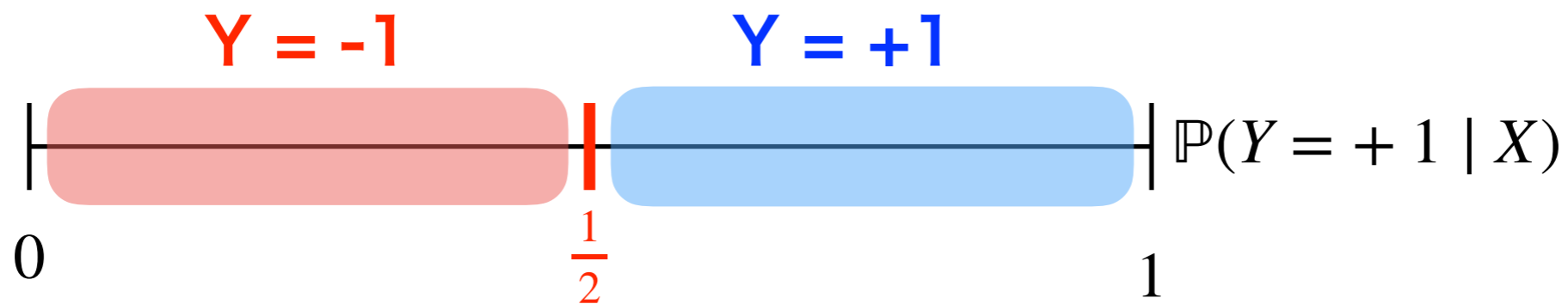
(informal)

# Related Work: Plug-in Rule

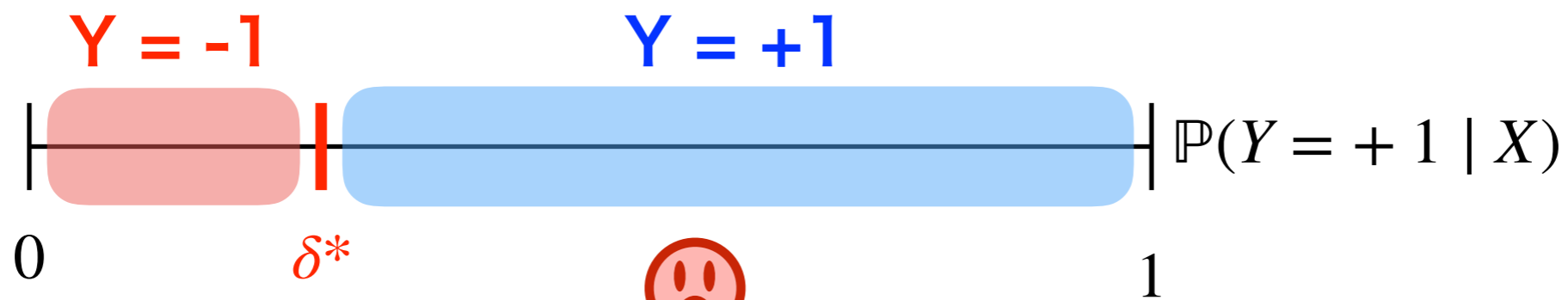
[Koyejo+ NIPS2014; Yan+ ICML2018]

## ■ Classifier based on class-posterior probability

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



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



⇒ 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.

# Outline

- Introduction
- Preliminary
  - ▶ Convex Risk Minimization
  - ▶ Plug-in Principle vs. Cost-sensitive Learning
- **Key Idea**
  - ▶ Quasi-concave Surrogate
- Calibration Analysis & Experiments

# Convexity & Statistical Property 15

tractable (convex)

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

generalize

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

calibration

$$R(f) = \mathbb{E}[\ell(Yf(X))]$$

**intractable**

①  
Q. tractable & calibrated  
objective?  
②

calibration

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

**intractable**

$$\operatorname{argmin} R_\phi = \operatorname{argmin} R$$

# Non-concave, but quasi-concave <sup>16</sup>

Idea: concave / convex = quasi-concave

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

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

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

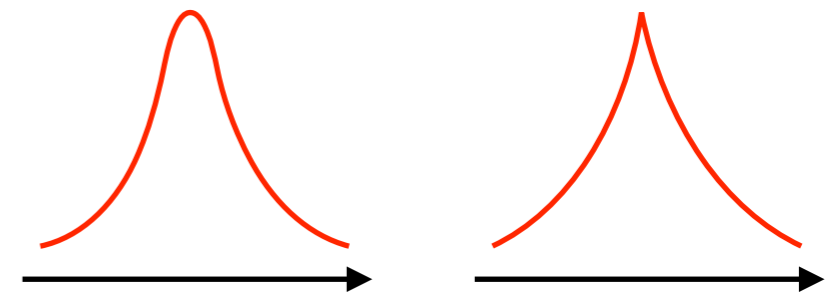
(proof) Show  $\{x | f/g \geq \alpha\}$  is convex.

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

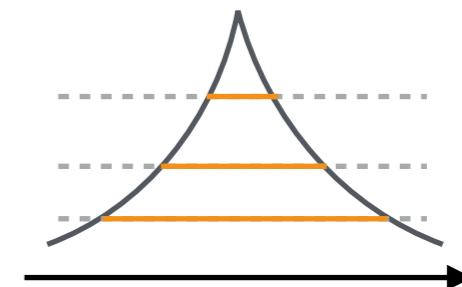
NB: super-level set of concave func.  
is convex

$\therefore \{x | f/g \geq \alpha\}$  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

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

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

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

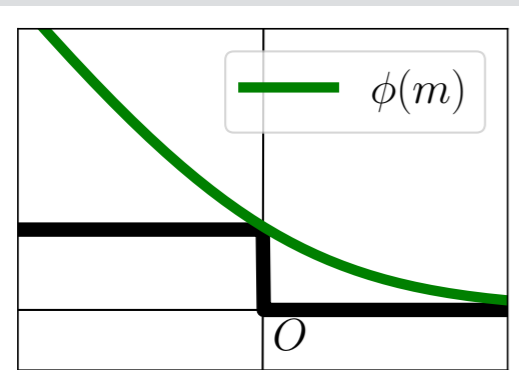
$$U(f) = \frac{a_0 \text{TP} + b_0 \text{FP} + c_0}{a_1 \text{TP} + b_1 \text{FP} + c_1}$$

$\geq$

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

$\equiv$

surrogate loss



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

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

# Hybrid Optimization Strategy <sup>19</sup>

$$U_\phi(f) = \frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + b_1 \mathbb{E}_N \left[ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \right] + c_1} = \frac{\text{Graph 1}}{\text{Graph 2}}$$

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 graph of a concave curve and a step function),  $+ b_0 \mathbb{E}_N$  (with a graph of a concave curve and a step function), and  $+ c_0$ . The denominator consists of  $a_1 \mathbb{E}_P$  (with a graph of a concave curve and a step function),  $+ b_1 \mathbb{E}_N$  (with a graph of a concave curve and a step function), and  $+ c_1$ . The result is shown as a fraction of two graphs: the top graph is a concave curve above a horizontal axis, and the bottom graph is a concave curve below a horizontal axis.

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

# Hybrid Optimization Strategy <sup>20</sup>

---

## Algorithm 1: Hybrid Optimization Algorithm

---

**Input** :  $\phi$  convex loss,  $\theta$  initial classifier  
parameter

**repeat**

|  $g^n \leftarrow \nabla_{\theta} \hat{U}_{\phi}^n(f_{\theta})$   
|  $\theta \leftarrow \text{gradient\_based\_update}(\theta, g^n)$

] maximize numerator

**until**  $\hat{U}_{\phi}^n(f_{\theta}) \leq 0$

**repeat**

|  $g \leftarrow \nabla_{\theta} \hat{U}_{\phi}(f_{\theta}), \hat{g} = g / \|g\|$   
|  $\theta \leftarrow \text{gradient\_based\_update}(\theta, \hat{g})$

] maximize fraction

**until** stopping criterion is satisfied

**Output:** maximizer  $f_{\theta}$

---

normalized gradient  
for quasi-concave optimization

[Hazan+ NeurIPS2015]

# Outline

- Introduction
- Preliminary
  - ▶ Convex Risk Minimization
  - ▶ Plug-in Principle vs. Cost-sensitive Learning
- Key Idea
  - ▶ Quasi-concave Surrogate
- **Calibration Analysis & Experiments**

# Justify Surrogate Optimization

## ■ For classification risk

surrogate risk

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

classification risk

$$R(f) = \mathbb{E}[\ell(Yf(X))]$$

If  $\phi$  is **classification-calibrated** loss, [Bartlett+ 2006]

$$R_\phi(f_n) \xrightarrow{n \rightarrow \infty} 0 \implies R(f_n) \xrightarrow{n \rightarrow \infty} 0 \quad \forall \{f_n\}$$

Note: informal

## ■ For fractional utility

surrogate utility

$$U_\phi(f) = \frac{\mathbb{E}_X[W_{0,\phi}(f(X))]}{\mathbb{E}_X[W_{1,\phi}(f(X))]}$$

true utility

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

**Q.** What kind of conditions are needed for  $\phi$  to satisfy

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

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

## Theorem

merely sufficient!

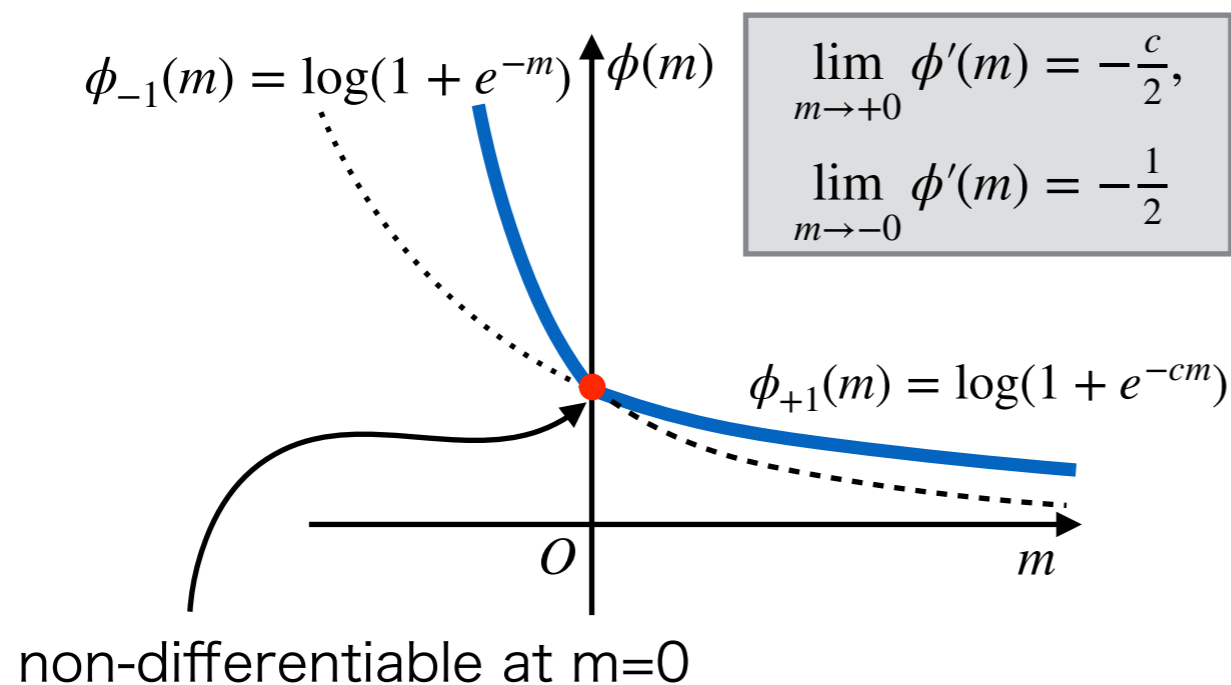
$$U_{\phi}(f_n) \xrightarrow{n \rightarrow \infty} 1 \implies U(f_n) \xrightarrow{n \rightarrow \infty} 1 \quad \forall \{f_n\}$$

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



# 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)

model: linear-in-parameter

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



# Experiment: Jaccard index

(Jaccard index)	Proposed		Baselines		
Dataset	U-GD	U-BFGS	ERM	W-ERM	Plug-in
adult	0.499 (44)	0.498 (11)	0.471 (51)	0.510 (20)	<b>0.516 (10)</b>
australian	<b>0.735 (63)</b>	<b>0.733 (59)</b>	0.702 (144)	0.693 (143)	0.707 (76)
breast-cancer	<b>0.921 (54)</b>	<b>0.918 (55)</b>	0.905 (66)	0.903 (78)	<b>0.913 (69)</b>
cod-rna	0.854 (3)	0.785 (8)	0.864 (11)	0.865 (9)	<b>0.869 (3)</b>
diabetes	<b>0.714 (44)</b>	0.702 (50)	0.692 (70)	0.698 (56)	0.695 (60)
fourclass	<b>0.469 (69)</b>	<b>0.457 (68)</b>	0.436 (112)	0.434 (171)	0.449 (66)
german.numer	<b>0.433 (64)</b>	<b>0.429 (69)</b>	0.335 (153)	0.391 (98)	<b>0.418 (71)</b>
heart	<b>0.665 (135)</b>	<b>0.675 (135)</b>	<b>0.664 (102)</b>	0.629 (178)	0.626 (163)
ionosphere	<b>0.826 (76)</b>	<b>0.829 (65)</b>	0.796 (134)	0.749 (245)	<b>0.815 (87)</b>
madelon	<b>0.495 (31)</b>	0.459 (69)	0.346 (225)	0.474 (100)	<b>0.496 (27)</b>
mushrooms	0.999 (2)	0.995 (4)	<b>1.000 (1)</b>	0.999 (4)	0.997 (7)
phishing	0.883 (43)	<b>0.893 (11)</b>	<b>0.894 (14)</b>	0.888 (22)	<b>0.894 (15)</b>
phoneme	0.435 (51)	0.436 (24)	0.371 (160)	<b>0.450 (104)</b>	<b>0.461 (34)</b>
skin_nonskin	0.744 (5)	0.751 (5)	0.746 (10)	<b>0.780 (13)</b>	0.722 (7)
sonar	<b>0.600 (125)</b>	<b>0.600 (111)</b>	0.552 (147)	0.495 (202)	<b>0.572 (134)</b>
spambase	<b>0.827 (22)</b>	0.708 (22)	0.798 (67)	0.790 (86)	<b>0.824 (31)</b>
splice	<b>0.670 (60)</b>	<b>0.672 (56)</b>	0.646 (71)	0.629 (84)	<b>0.672 (57)</b>
w8a	0.496 (151)	0.452 (28)	0.580 (44)	<b>0.590 (35)</b>	<b>0.595 (33)</b>

(Jaccard index is shown)

model: linear-in-parameter

surrogate loss:  $\phi(m) = \max\{\log(1 + e^{-m}), \log(1 + e^{-\frac{3m}{4}})\}$

## ■ Goal

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

maximize linear-fractional utility

## ■ Tractable Optimization

surrogate utility

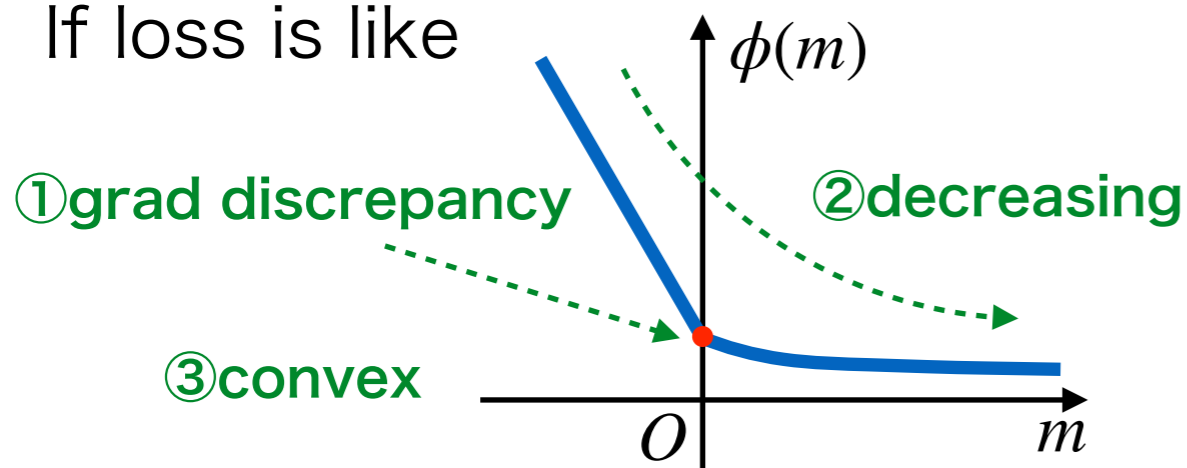
$$\frac{a_0 \mathbb{E}_P \left[ \begin{array}{|c|} \hline \text{step} \\ \hline \end{array} \right] + b_0 \mathbb{E}_N \left[ \begin{array}{|c|} \hline \text{step} \\ \hline \end{array} \right] + c_0}{a_1 \mathbb{E}_P \left[ \begin{array}{|c|} \hline \text{step} \\ \hline \end{array} \right] + b_1 \mathbb{E}_N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

quasi-concave optimization



## ■ Calibrated Surrogate

If loss is like



then

$$\operatorname{argmax}_f U_\phi(f) = \operatorname{argmax}_f U(f)$$

## ■ Open Problems

- ▶ necessary and sufficient condition of calibration
- ▶ explicit convergence rate
- ▶ theoretical comparison with probability estimation