## On the Calibration of Multiclass Classification with Rejection Chenri Ni<sup>1</sup> Nontawat Charoenphakdee<sup>1,2</sup> Junya Honda<sup>1,2</sup> Masashi Sugiyama<sup>2,1</sup>



## Introduction: Learning with rejection





https://me.me/i/the-right-way-to-answer-true-and-false-questions-18781463 Saying "I don't know" can prevent misclassification. **Related work:** 

Approach	Binary	
<b>Confidence-base</b>	Bartlett+ (2008); Yuan+ (2010)	Ran
<b>Classifier-rejector</b>	Cortes+ (2015, 2016)	

**Ramaswamy+ (2018)** only focused on specific types of non-differentiable losses. **Contributions:** 

- Calibration condition for surrogate losses in the classifier-rejector approach, which suggests the difficulty especially in the multiclass case
- Excess risk bounds and estimation error bounds to guarantee the one-vs-all (OVA) and cross-entropy (CE) losses in the confidence-based approach

Multiclass classification with	rejec
Given: Labeled data: $\{(x_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(x, y)$ Rejection cost: $c \in (0, 0.5)$ Find: Classifier: $f(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} g_y(x) \in \mathcal{Y}$ Rejector: $r(x) \in \mathbb{R}$	Chow (197 $x \in \mathcal{X} \subseteq y \in \mathcal{Y} = g_i(x)$ : decision( $x$ )
Goal: Minimize $R_{0-1-c}(r, f) = \underset{p(\boldsymbol{x}, y)}{\mathbb{E}} [\mathcal{L}_{0-1-c}(r, f;$	[x, y)]
$\mathcal{L}_{0-1-c}(r, f; \boldsymbol{x}, \boldsymbol{y}) = \underbrace{\mathbb{I}[f(\boldsymbol{x}) \neq \boldsymbol{y}] \mathbb{I}[r(\boldsymbol{x}) > 0]}_{\text{misclassification loss}}$ $\mathcal{L}_{0-1-c}(r, f; \boldsymbol{x}, \boldsymbol{y}) \text{ is difficult to directly optimize}_{\text{Yuan+ (2010); Cortes+ (2015, 2)}}$ A computationally-efficient and theoretically justified s	rejectio
Calibration	
Calibration ensures that minimizing a surrogate loss will lea	ad to an
Optimal solution of classification with rejection	tion:
$f^*(\boldsymbol{x}) = \arg \max n_u(\boldsymbol{x})$ $n_u(\boldsymbol{x}) = n(u \boldsymbol{x})$	Chang
$r^{*}(\boldsymbol{x}) = \max_{\substack{y \in \mathcal{Y} \\ y \in \mathcal{Y}}} \eta_{y}(\boldsymbol{x}) - (1-c)$	c) Chow
$f(x) = \underset{y \in \mathcal{Y}}{\underset{y \in \mathcal{Y}}{\max}} \eta_y(x) - (1 - c)$ $f(x, f) \text{ is calibrated if } R_{0-1-c}(r, f) = R_{0-1-c}(r, f)$ $f(x) = r \text{ is classification-calibrated if } f(x) = r$ $f(x) = r \text{ is rejection-calibrated if } sign[r(x)]$	chow $r^*, f^*)$ $f^*(\boldsymbol{x})$ $=  ext{sign}$

If (r, f) is calibrated, r must be rejection-calibrated. A minimizer of a surrogate loss should give a calibrated (r, f).

## 1: The University of Tokyo 2: RIKEN AIP



$$(\boldsymbol{x}) - r(\boldsymbol{x}) \Big) + c\psi \big(\beta r(\boldsymbol{x}) \big)$$

Electronic Journal of Statistics, 2018.

[5] H.G. Ramaswamy, A. Tewari, and S. Agarwal. Consistent algorithms for multiclass classification with an abstain option.

[6] M. Yuan, M.H. Wegkamp. Classification methods with reject option based on convex risk minimization. JMLR, 2010. [7] Y. Lecun, The MNIST database of handwritten digits. <u>http://yann.lecun.com/exdb/mnist/</u>, 1998.

			Conf	fider	ice-ba	ase	d	appi	roac	h		
					Bartle	tt+ (20	08);	Yuan+ (	(2010); F	Ramaswai	my+ (2018	;)
		Reject	tor de	pends	solelv	on cl	as	sifier'	<mark>s</mark> conf	idence		
				•	• Cross-e	entropy	/ (C	E) loss:				
					$\mathcal{L}_{ ext{CE}}(j$	$f; oldsymbol{x}, y)$	=	$-g_y(\boldsymbol{x})$	$) + \log \Sigma$	$\sum_{y'\in\mathcal{Y}} \mathbf{e}_{\mathbf{\mathcal{Y}}}$	$\exp\left(g_{y'}(oldsymbol{x}$	))
•					• One-ve	ersus-a	II (C	DVA) los	s:	<i>5</i> - <b>0</b>	·	ŗ
$\mathcal{L}_{\mathrm{OVA}}(f;;$							$f(\boldsymbol{x}, y) = \phi(g_{y}(\boldsymbol{x})) + \sum_{u' \neq u} \phi(-g_{u'}(\boldsymbol{x}))$					
					• Rejector: $g(x) = [g_1(x), \dots, g_K(x)]^{ op}$							
			•		$r_f(\boldsymbol{x}) = \max_{\boldsymbol{u} \in \mathcal{V}} \Psi^{-1}(\boldsymbol{g}(\boldsymbol{x})) - (1-c)$							
					$\Psi^{-1} \colon \mathbb{R}^K \to [0,1]^K$ Inverse link function							
					<b>π</b> -1 (		φ'(	$(-a_n)$	_ 1		$\exp(a)$	
					$\Psi_{y,{ m OVA}}^{-1}(g)$ See our pap	$({m g})=rac{1}{\phi'({m g})}$ , where ${m for}$ is the component of the	$\frac{\varphi}{(-g_y)}$	$(g_y)$ $(g_y) + \phi'(g_y)$ $(g_y)$ $(g_y)$ $(g_y)$	$\Psi_{y,\mathrm{CE}}^{-1}$	$\mathbf{g}(oldsymbol{g}) = rac{1}{\sum_y}$ Softmax fun	$\frac{\exp(g_y)}{f \in \mathcal{Y}} \exp(g_{y'})$ ction	
Ve pr	ovi	de exc	ess ris	k bour	nds to g	uarar	nte	e OVA	and Cl	E losses	•	
Exces	s r	isk:										
	4	$\Delta R_{0-1-}$	$_{c}(r_{f},f)$	$) = R_{0}$	$1-c(r_f, f$	r) —	j	inf	$R_{0-1-c}$	$(r_f, f)$		
			$\Lambda R_{c}(f)$	$) - R_{ab}$	(f)	J' inf	:me	$\frac{1}{R} \left( f' \right)$	)			
		2	$\Delta n_{\ell}(J)$	) = Ill	(J) - f':	measura	able		)			
Exces	s ri	sk bou	ind of	OVA lo	DSS:			Loss N	Vame	$\phi(z)$	C	8
(2C	')-	$^{s}\Delta R_{0}$	-1-c(r)	$(f,f)^s$	$\leq \Delta R$	COVA	(f)	Expone	ential	$\exp(1 + \exp(-z))$	$(-z)) = \frac{1}{2}$	$\frac{2}{2}$
Extensio	on o	f the res	ult by <mark>Yu</mark> a	n+ (2010)	to the mu	Ilticlass	case	e. Squa	red	$(1-z)^2$	$\frac{1}{2}$	2
Exces	ss r	isk bo	und of	CE los	S:			Squared	Hinge	$(1-z)_{+}^{-}$	- 2	Ζ
	$\frac{1}{2}$	$\Delta R_0$	-1-c(r)	$(f, f)^2$	$\leq \Delta I$	$R_{\rm CE}($	f)					
Needs	ے anal	ysis spec	ific to the	e multicla	ass case wh	nere pre	vio	us techni	ques can	not be ap	plied.	
Mi	inir	nizers	of OV	A and	CE losse	s also	) m	ninimiz	e the	0-1-c lo	SS.	
				See	our paper fo	or estima	atior	n error bo	und using	Rademacl	ner complex	ity.
				E	Exper	ime	n	ts				
lassifie	er-r	ejector	MPC+l	og (MPC	c with log	istic los	ss),	APC+lo	g (APC \	with logi	stic loss)	
onfide	ence	e-based	: OVA+h	nin by <mark>R</mark> a	amaswamy	y+ (201	8),	OVA+lo	g (OVA v	with logi	stic loss),	CE
0.225	erro	vehicle		-		atimage			0 225	letter		
0.200	— MPC — OVA — CE	C+log A+log	1	0.12	OVA+log CE		1		0.200 0 0.175 0	IPC+log VA+log E		
0.175 	— OVA	A+hin			OVA+hin				0.175 O	VA+hin		
0.125 0.100				0.00					9 0.125 0.100	1/		
0.075				0.04	4				0.075			
0.025	0.1	0.2	0.3 0.	4	0.1	0.2 0.	3	0.4	0.025	1 0.2	0.3 0.4	
ccura	cy c	of non-	rejecte	d data:	"- (-)" ind	dicates	s al	l data v	vere rej	ected.		
dataset	C O O O	APC+log	MPC+log	OVA+log	CE							
vehicle	0.05	- ( - ) 98.4 (1.9)	96.6 (2.3) 92.4 (3.0)	100 (0.0) 97.9 (0.7)	100 (0.0) 97.4 (0.1)	dataset	<i>c</i>	APC+log	MPC+log	OVA+log	CE	
	0.4	89.1 (2.9)	85.3 (4.2)	90.2 (1.6)	<b>91.7</b> (0.9)	covtype	0.05	7 <b>9.5 (2.1)</b> 74.0 (1.8)	79.8 (1.7) 73.8 (1.0)	<b>82.1</b> (2.7) 74.9 (1.4)	82.0 (3.2) 77.1 (0.3)	
satimage	0.05	<b>99.1 (0.2)</b> 95.0 (1.0)	97.2 (1.4) 92.6 (1.2)	98.7 (0.1) 96.2 (0.2)	98.3 (0.1) 95.7 (0.1)		0.4	<b>69.8</b> (1.3)	64.9 (3.4)	<b>68.7</b> (1.1)	<b>69.4</b> ( <b>1.8</b> )	
	0.4	91.5 (0.7)	89.0 (1.1)	92.2 (0.3)	91.8 (0.2)	letter	0.05	<b>99.8 (0.1)</b> 97.9 (0.3)	98.6 (0.2) 96.9 (0.5)	99.6 ( 0.2 ) 98.3 (0.2)	99 8 (0.0) 98.4 (0.1)	
veast	0.05 0.2	- ( - ) - ( - )	- ( - ) - ( - )	- ( - ) - ( - )	- ( - ) 80.6 (6.2)		0.4	95.2 (0.5)	94.6 (3.8)	94.6 (0.2)	94.9 (0.3)	
jeast	0.4	- ( - )	- ( - )	75.0 (3.9)	76.6 (1.7)							



