# On the Calibration of Multiclass Classification with Rejection
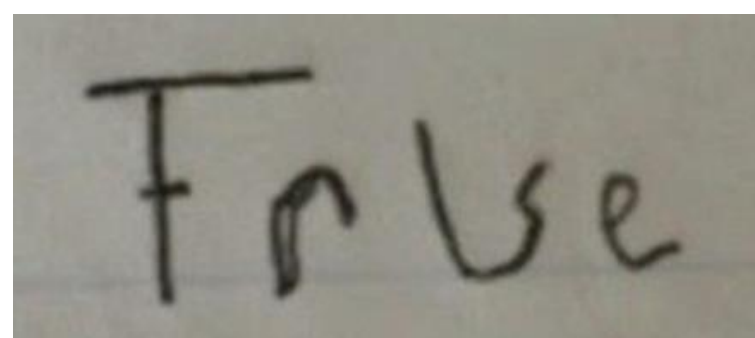
Chenri Ni[1] Nontawat Charoenphakdee[1,2] Junya Honda[1,2] Masashi Sugiyama[2,1]

1: The University of Tokyo    2: RIKEN AIP

## Introduction

Q: True or False!?

Q: 1,2 or 7 !?

Saying **"I don't know"** can **prevent misclassification**.
**Most theoretical works in this problem focused on binary case**.
Only Ramaswamy+ 2018 considered confidence-based approach in multiclass case.
**Contributions:**
- **An analysis of a recent classifier-rejector approach in multiclass case.**
- **Theoretical guarantee for well-known surrogate losses for confidence-based approach.**

## Multiclass classification with rejection

(Chow 1970, Ramaswamy+ 2018)

**Given:** Labeled data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n \overset{\text{i.i.d.}}{\sim} p(\boldsymbol{x}, y)$   $x \in \mathcal{X} \subseteq \mathbb{R}^d$

Rejection cost $c \in (0, 0.5)$   $y \in \mathcal{Y} = \{1, \dots, K\}$

**Find:** Classifier $f(\boldsymbol{x}) = \underset{y \in \mathcal{Y}}{\arg\max}\, g_y(\boldsymbol{x})$   $g_i(\boldsymbol{x}): \mathcal{X} \to \mathbb{R}$

Rejector $r: \mathcal{X} \to \mathbb{R}$

that minimizes the following risk: $R_{0\text{-}1\text{-}c}(r, f) = \underset{p(\boldsymbol{x}, y)}{\mathbb{E}}[\mathcal{L}_{0\text{-}1\text{-}c}(r, f; \boldsymbol{x}, y)]$

where $\mathcal{L}_{0\text{-}1\text{-}c}(r, f; \boldsymbol{x}, y) = \underbrace{\mathbb{1}_{[f(\boldsymbol{x}) \neq y]} \mathbb{1}_{[r(\boldsymbol{x}) > 0]}}_{\text{misclassification loss}} + \underbrace{c\, \mathbb{1}_{[r(\boldsymbol{x}) \leq 0]}}_{\text{rejection loss}}$

$\mathcal{L}_{0\text{-}1\text{-}c}(r, f; \boldsymbol{x}, y)$ is hard to directly optimize.

(Yuan+, 2010, Cortes+ (2015, 2016), Ramaswamy+ 2018)

**A computationally-efficient and theoretically justified surrogate loss is needed.**

**Optimal solution of classification with rejection:**

$f^*(\boldsymbol{x}) = \underset{y \in \mathcal{Y}}{\arg\max}\, \eta_y(\boldsymbol{x})$   (Chow 1970)
$r^*(\boldsymbol{x}) = \underset{y \in \mathcal{Y}}{\max}\, \eta_y(\boldsymbol{x}) - (1 - c)$   $\eta_y(\boldsymbol{x}) = p(y|\boldsymbol{x})$

## Calibration

**Calibration ensures that minimizing a surrogate loss will lead to an optimal solution**

- $(r, f)$ is **calibrated** if $R_{0\text{-}1\text{-}c}(r, f) = R_{0\text{-}1\text{-}c}(r^*, f^*)$
- $f$ is **classification-calibrated** if $f(\boldsymbol{x}) = f^*(\boldsymbol{x})$
- $r$ is **rejection-calibrated** if $\text{sign}[r(\boldsymbol{x})] = \text{sign}[r^*(\boldsymbol{x})]$
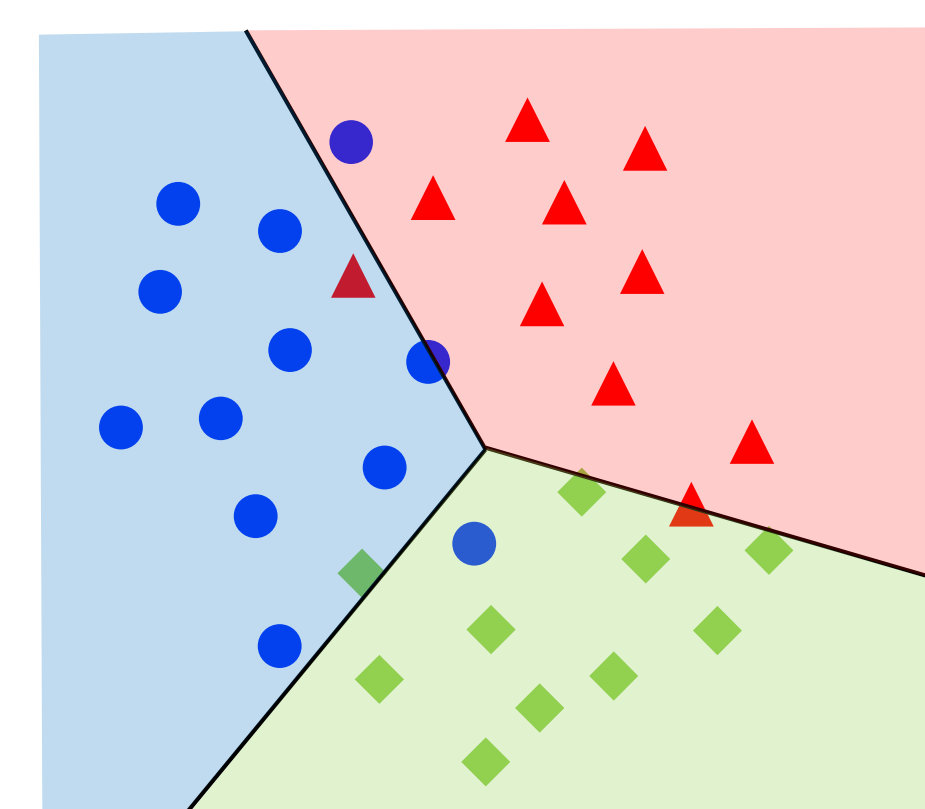
If $(r, f)$ is calibrated, $r$ must be rejection-calibrated.

**A minimizer of a surrogate loss should give a calibrated $(r, f)$**
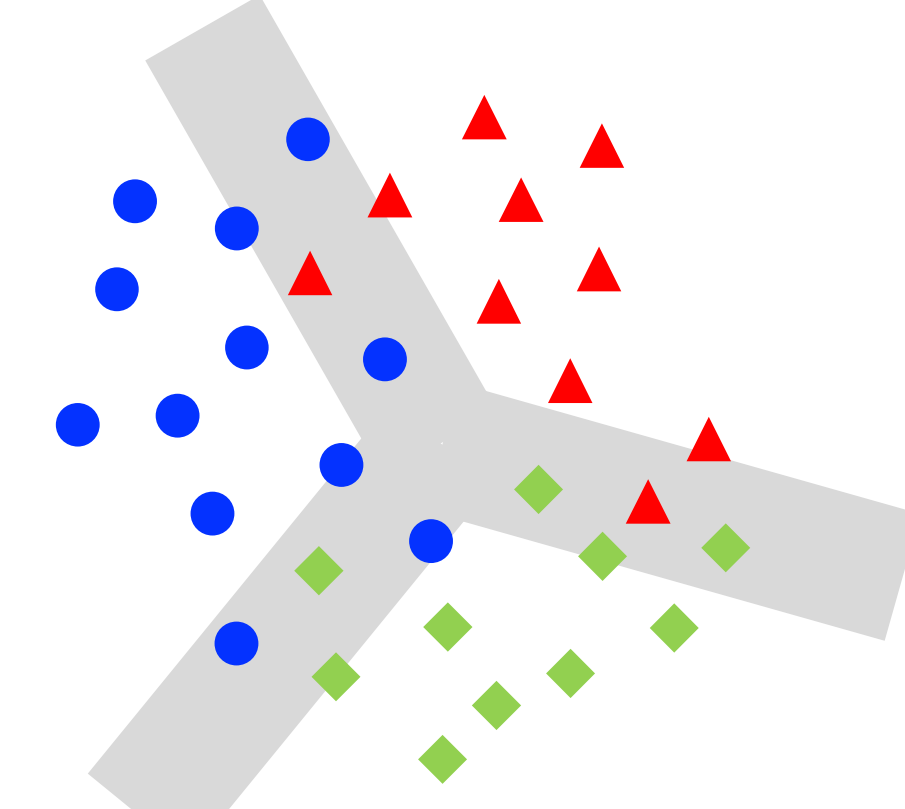
## Classifier-rejector approach

(Cortes+, 2015, 2016)

**Classifier** and **rejector** are trained simultaneously

**Classifier**
(colored area indicates classifier prediction)

**Rejector**
(gray area indicates rejection area)

$(r_{\boldsymbol{\eta}}^\dagger, f_{\boldsymbol{\eta}}^\dagger) = \underset{r \in \mathbb{R},\, \boldsymbol{g} \in \mathbb{R}^K}{\arg\min}\, W(r, f; \boldsymbol{\eta})$   $\boldsymbol{\eta}(\boldsymbol{x}) = [\eta_1(\boldsymbol{x}), \dots, \eta_K(\boldsymbol{x})]^\top$   $W(r(\boldsymbol{x}), f(\boldsymbol{x}); \boldsymbol{\eta}(\boldsymbol{x})) = \sum_{y \in \mathcal{Y}} \eta_y(\boldsymbol{x}) \mathcal{L}(r, f; \boldsymbol{x}, y)$

### Corollary 5: (Necessary condition for rejection calibration)

**For $\mathcal{L}(r, f; \boldsymbol{x}, y)$ that is convex with respect to $r$ and $\left.\frac{\partial^2 W(r, f_{\boldsymbol{\eta}}^\dagger; \boldsymbol{\eta})}{\partial r^2}\right|_{r=0} > 0$**

$r^\dagger$ is rejection-calibrated **only if** both conditions hold:

**Condition (1)**
$\underset{\boldsymbol{\eta}:\, \max_y \eta_y = 1 - c}{\sup} \left.\frac{\partial W(r, f_{\boldsymbol{\eta}}^\dagger; \boldsymbol{\eta})}{\partial r}\right|_{r=0} = 0$
Condition for false reject rate to be zero

**Condition (2)**
$\underset{\boldsymbol{\eta}:\, \max_y \eta_y = 1 - c}{\inf} \left.\frac{\partial W(r, f_{\boldsymbol{\eta}}^\dagger; \boldsymbol{\eta})}{\partial r}\right|_{r=0} = 0$
Condition for false accept rate to be zero

Necessary and sufficient condition is also provided in our paper (Theorem 4)

$\underset{\boldsymbol{\eta}:\, \max_y \eta_y = 1 - c}{\sup} \left.\frac{\partial W(r, f_{\boldsymbol{\eta}}^\dagger; \boldsymbol{\eta})}{\partial r}\right|_{r=0} = \underset{\boldsymbol{\eta}:\, \max_y \eta_y = 1 - c}{\inf} \left.\frac{\partial W(r, f_{\boldsymbol{\eta}}^\dagger; \boldsymbol{\eta})}{\partial r}\right|_{r=0} = 0$

**Supremum and infimum values coincide under the same constraint.**

**When** $\max_y \eta_y = 1 - c$

- Binary case: $\boldsymbol{\eta}$ can only be either $[1 - c, c]^\top$ or $[c, 1 - c]^\top$
- Multiclass case: $\boldsymbol{\eta}$ can be arbitrary. **Both conditions can be very different and do not hold simultaneously!**

### Case study:

$\alpha \in \mathbb{R}$   $\beta \in \mathbb{R}$   Hyperparameters
$\phi: \mathbb{R} \to \mathbb{R}$   $\psi: \mathbb{R} \to \mathbb{R}$   Convex margin losses

- **Multiplicative pairwise comparison (MPC) loss:**
$\mathcal{L}_{\text{MPC}}(r, f; \boldsymbol{x}, y) = \sum_{y' \neq y} \phi\big(\alpha(g_y(\boldsymbol{x}) - g_{y'}(\boldsymbol{x}))\big)\psi(-\alpha r(\boldsymbol{x})) + c\psi(\beta r(\boldsymbol{x}))$

- **Additive pairwise comparison (APC) loss:**
$\mathcal{L}_{\text{APC}}(r, f; \boldsymbol{x}, y) = \sum_{y' \neq y} \phi\big(\alpha(g_y(\boldsymbol{x}) - g_{y'}(\boldsymbol{x}) - r(\boldsymbol{x}))\big) + c\psi(\beta r(\boldsymbol{x}))$

Consider $\phi(z) = \psi(z) = \exp(-z)$

Condition (1) gives
$\frac{\beta}{\alpha} = (K - 2) + 2\sqrt{(K - 1)\frac{1 - c}{c}}$

Condition (2) gives
$\frac{\beta}{\alpha} = 2\sqrt{\frac{1 - c}{c}}$

Equivalent to a condition proved by (Cortes+, 2016) when considering a binary case ($K = 2$).
**In multiclass case, $(\alpha, \beta)$ that satisfies both conditions simultaneously does not exist.**
Similar results also hold when using the logistic loss $\phi(z) = \psi(z) = \log(1 + \exp(-z))$.

## References

[1] C. K. Chow. On optimum recognition error and reject tradeoff. IEEE Transaction on Information Theory, 1970
[2] C. Cortes, G. DeSalvo, and M. Mohri. Learning with rejection. ALT, 2015
[3] Cortes G. DeSalvo, and M. Mohri. Boosting with abstention. NeurIPS, 2016
[4] H.G. Ramaswamy, A. Tewari, and S. Agarwal. Consistent algorithms for multiclass classification with an abstain option. EJS, 2018.
[5] M. Yuan, M.H. Wegkamp. Classification methods with reject option based on convex risk minimization. JMLR, 2010.

## Confidence-based approach

(Bartlett+ 2008, Yuan+ 2010, Ramaswamy+ 2018)

**Rejector** depends solely on **classifier's** confidence

- **Cross-entropy (CE) loss:**
$\mathcal{L}_{\text{CE}}(f; \boldsymbol{x}, y) = -g_y(\boldsymbol{x}) + \log \sum_{y' \in \mathcal{Y}} \exp(g_{y'}(\boldsymbol{x}))$

- **One-versus-all (OVA) loss:**
$\mathcal{L}_{\text{OVA}}(f; \boldsymbol{x}, y) = \phi(g_y(\boldsymbol{x})) + \sum_{y' \neq y} \phi(-g_{y'}(\boldsymbol{x}))$

- **Rejector:**
$r_f(\boldsymbol{x}) = \underset{y \in \mathcal{Y}}{\max}\, \Psi^{-1}(\boldsymbol{g}(\boldsymbol{x})) - (1 - c)$

$\boldsymbol{g}(\boldsymbol{x}) = [g_1(\boldsymbol{x}), \dots, g_K(\boldsymbol{x})]^\top$
$\Psi^{-1}: \mathbb{R}^K \to [0, 1]^K$ Inverse link function

$\Psi^{-1}_{y, \text{OVA}}(\boldsymbol{g}) = \frac{\phi'(-g_y)}{\phi'(-g_y) + \phi'(g_y)}$   $\Psi^{-1}_{y, \text{CE}}(\boldsymbol{g}) = \frac{\exp(g_y)}{\sum_{y' \in \mathcal{Y}} \exp(g_{y'})}$
See our paper for conditions of $\phi$.   Softmax function

**Excess risk:**

$\Delta R_{0\text{-}1\text{-}c}(r_f, f) = R_{0\text{-}1\text{-}c}(r_f, f) - \underset{f':\text{measurable}}{\inf} R_{0\text{-}1\text{-}c}(r_f, f)$

$\Delta R_\ell(f) = R_\ell(f) - \underset{f':\text{measurable}}{\inf} R_\ell(f')$

If $\Delta R_{0\text{-}1\text{-}c}$ can be upper-bounded by $\Delta R_\ell$,
-> then the **minimizer of both risks are identical**.

**Excess risk bound of OVA loss:**
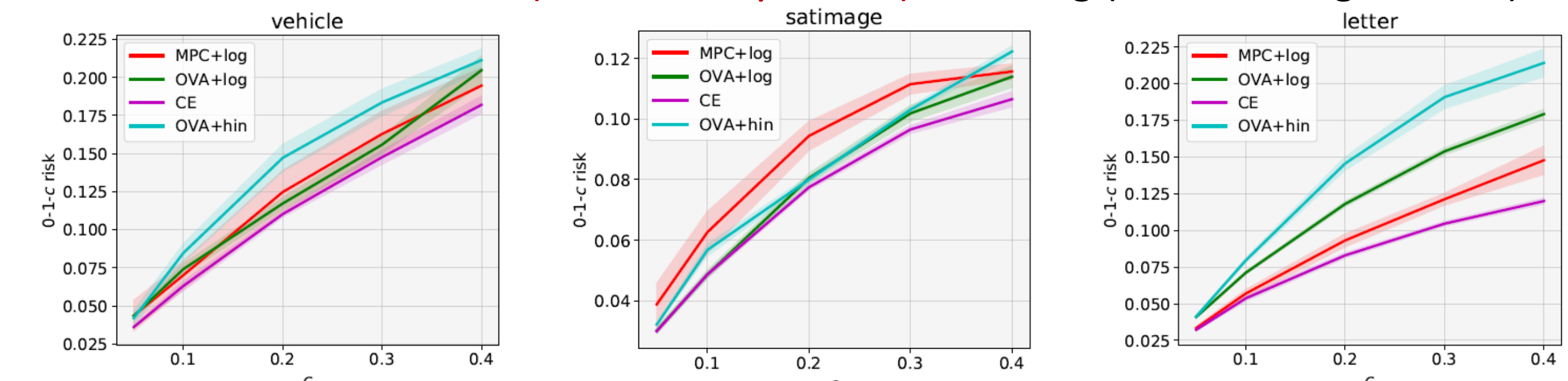$(2C)^{-s}\Delta R_{0\text{-}1\text{-}c}(r_f, f)^s \leq \Delta R_{\text{OVA}}(f)$

**Excess risk bound of CE loss:**
$\frac{1}{2}\Delta R_{0\text{-}1\text{-}c}(r_f, f)^2 \leq \Delta R_{\text{CE}}(f)$

| Loss Name | $\phi(z)$ | $C$ | $s$ |
|---|---|---|---|
| Logistic | $\log(1 + \exp(-z))$ | $\frac{1}{2}$ | 2 |
| Exponential | $\exp(-z)$ | $\frac{1}{\sqrt{2}}$ | 2 |
| Squared | $(1 - z)^2$ | $\frac{1}{2}$ | 2 |
| Squared Hinge | $(1 - z)_+^2$ | $\frac{1}{2}$ | 2 |

See our paper for more results, e.g., estimation error bound using Rademacher complexity.

## Experiment

**Classifier-rejector:** MPC+log (MPC with logistic loss), APC+log (APC with logistic loss)
**Confidence-based:** OVA+hin (Ramaswamy+ 2018), OVA+log (OVA with logistic loss), CE

vehicle   satimage   letter

**Accuracy of non-rejected data:** "- (-)" indicates all data were rejected.

| dataset | $c$ | APC+log | MPC+log | OVA+log | CE |
|---|---|---|---|---|---|
| vehicle | 0.05 | - (-) | 96.6 (2.3) | 100 (0.0) | **100 (0.0)** |
|  | 0.2 | 98.4 (1.9) | 92.4 (3.0) | 97.9 (0.7) | **97.4 (0.1)** |
|  | 0.4 | **89.1 (2.9)** | 85.3 (4.2) | 90.2 (1.6) | **91.7 (0.9)** |
| satimage | 0.05 | **99.1 (0.2)** | 97.2 (1.4) | **98.7 (0.1)** | **98.3 (0.1)** |
|  | 0.2 | 95.0 (1.0) | 92.6 (1.2) | 96.2 (0.2) | **95.7 (0.1)** |
|  | 0.4 | 91.5 (0.7) | 89.0 (1.1) | 92.2 (0.3) | **91.8 (0.2)** |
| yeast | 0.05 | - (-) | - (-) | - (-) | - (-) |
|  | 0.2 | - (-) | - (-) | - (-) | **80.6 (6.2)** |
|  | 0.4 | - (-) | - (-) | 75.0 (3.9) | **76.6 (1.7)** |

| dataset | $c$ | APC+log | MPC+log | OVA+log | CE |
|---|---|---|---|---|---|
| covtype | 0.05 | **79.5 (2.1)** | 79.8 (1.7) | **82.1 (2.7)** | **82.0 (3.2)** |
|  | 0.2 | 74.0 (1.8) | 73.8 (1.0) | 74.9 (1.4) | **77.1 (0.3)** |
|  | 0.4 | **69.8 (1.3)** | 64.9 (3.4) | 68.7 (1.1) | **69.4 (1.8)** |
| letter | 0.05 | **99.8 (0.1)** | 98.6 (0.2) | **99.6 (0.2)** | **99 8 (0.0)** |
|  | 0.2 | 97.9 (0.3) | 96.9 (0.5) | **98.3 (0.2)** | **98.4 (0.1)** |
|  | 0.4 | **95.2 (0.5)** | 94.6 (3.8) | 94.6 (0.2) | **94.9 (0.3)** |