

# Is Complementary-Label Learning Realistic?

Hsuan-Tien Lin

林軒田

Professor, National Taiwan University



September 25, 2024

Department of Mathematics, National University of Singapore

# About Me

## Hsuan-Tien Lin

Professor  
National Taiwan University



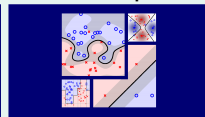
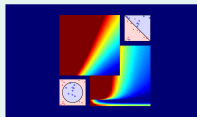
Chief Data Science Consultant  
(former Chief Data Scientist)  
Appier Inc.

The logo for Appier Inc., consisting of the word "Appier" in a stylized blue font.

Co-author  
Learning from Data

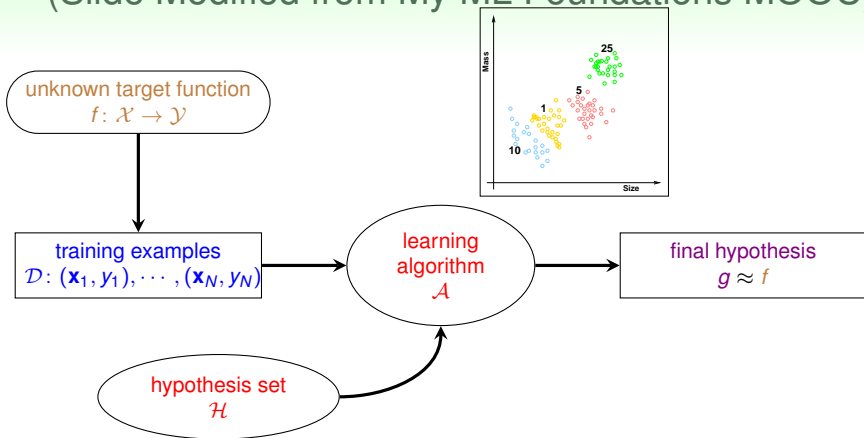


Instructor  
NTU-Coursera Mandarin MOOCs  
ML Foundations/Techniques



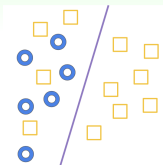
# Supervised Learning

(Slide Modified from My ML Foundations MOOC)



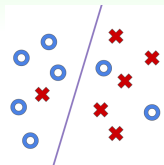
supervised learning:  
every input vector  $\mathbf{x}_n$  with  
its (possibly expensive) label  $y_n$ ,

# Weakly-supervised: Learning without True $y_n$



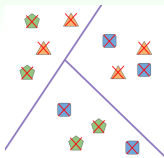
(a) Positive-unlabeled Learning [CE2008]

incomplete



(b) Learning with Noisy Labels [NN2013]

inaccurate



(c) Complementary-label Learning [TI2017]

inexact

- positive-unlabeled: **some** of true  $y_n = +1$  revealed
- noisy: **possibly incorrect** label  $y'_n$  instead of true  $y_n$
- complementary: **false label**  $\bar{y}_n$  instead of true  $y_n$

**weakly-supervised**: claimed to be a **realistic** route for reducing labeling burden

# Complementary-Label Learning

complementary label  $\bar{y}_n$  instead of true  $y_n$



Figure 1 of [XY2018]

potential to reducing labeling burden [TI2017]

- 1 ordinary label per instance
- $(K - 1)$  complementary labels per instance, **just need one of them**

complementary label: possibly **easier/cheaper**  
to obtain for some applications

# Example: Fruit Labeling Task



(left: from 2020 AICup in Taiwan; right: [publicdomainvectors.org](https://publicdomainvectors.org))

## hard: true label

- orange ?
- mango ?
- cherry
- banana

## easy: complementary label

- orange
- mango
- cherry
- banana ✗

can also help improve other ML tasks,  
like **semi-supervised learning** [QD2024]

# Comparison to Ordinary Learning

## Ordinary (Supervised) Learning

training:  $\{(\mathbf{x}_n = \text{mango}, y_n = \text{mango})\} \rightarrow \text{classifier } g(\mathbf{x})$



## Complementary-Label Learning

training:  $\{(\mathbf{x}_n = \text{mango}, \bar{y}_n = \text{banana})\} \rightarrow \text{classifier } g(\mathbf{x})$



goal during testing:  $\operatorname{argmax}_{k \in [K]} \left( g \left( \text{cherry} \right) \right)_k \rightarrow \text{cherry}$



ordinary versus complementary:  
same goal via **different training data**

# Formal Setup of Complementary-Label Learning

input	complementary label
	banana



## Given

size- $N$  data  $\mathcal{D} = \{(\text{input } \mathbf{x}_n \in \mathcal{X}, \text{complementary label } \bar{y}_n \in [K])\}_{n=1}^N$   
such that  $\bar{y}_n \neq y_n$  for some hidden ordinary label  $y_n \in [K]$

## Goal

a multi-class classifier  $g(\mathbf{x})$  that **closely predicts the ordinary label**  $y$   
associated with some unseen inputs  $\mathbf{x}$

todo: two CLL models, **and more!**



Yu-Ting Chou, Gang Niu, Hsuan-Tien Lin, and Masashi Sugiyama. **Unbiased risk estimators can mislead: A case study of learning with complementary labels.** ICML 2020.

# Review: Risk Minimization in Ordinary Learning

- goal: minimize the 0/1 loss

$$\ell_{01}(y, g(\mathbf{x})) = \mathbb{I}[y \neq \operatorname{argmax}_{k \in [K]} (g(\mathbf{x}))_k]$$

with risk (average loss)  $R_{01} = \mathbb{E}_{(\mathbf{x}, y)} \{\ell_{01}(y, g(\mathbf{x}))\}$

- consider a surrogate loss  $\ell$  that replaces  $\ell_{01}$

$$\ell: [K] \times \mathbb{R}^K \rightarrow \mathbb{R}_+$$

with risk  $R_\ell = \mathbb{E}_{(\mathbf{x}, y)} \{\ell(y, g(\mathbf{x}))\}$

**Empirical Risk Minimization (ERM):**  
estimate  $R_\ell$  by training data and minimize it

# Unbiased Risk Estimation for CLL

## Ordinary Learning

- ERM: minimizes

$$\hat{R}_\ell = \mathbb{E}_{(\mathbf{x}_n, y_n) \in \mathcal{D}} \{ \ell(y_n, g(\mathbf{x}_n)) \},$$

the empirical version of the surrogate risk  $R_\ell = \mathbb{E}_{(\mathbf{x}, y)} \{ \ell(y, g(\mathbf{x})) \}$

## Unbiased Risk Estimator for CLL [TI2019]

- [under assumption on  $P(\bar{y}|y)$ ] rewrite  $\ell$  to **some**  $\bar{\ell}$  such that

$$\bar{R}_{\bar{\ell}} = \mathbb{E}_{(\mathbf{x}, \bar{y})} \bar{\ell}(\bar{y}, g(\mathbf{x})) = \mathbb{E}_{(\mathbf{x}, y)} \ell(y, g(\mathbf{x})) = R_\ell$$

- $\bar{R}_{\bar{\ell}}$  called **unbiased risk estimator** (URE)
- URE-CLL: minimize empirical version  $\hat{\bar{R}}_{\bar{\ell}}$  of URE

URE-CLL: **pioneer model** for CLL, with **theoretical guarantees** like consistency

# Example of URE-CLL

## cross-entropy loss

for  $g(\mathbf{x}) = \mathbf{p}(k | \mathbf{x})$ ,

- $\ell_{CE}$ : surrogate of  $\ell_{01}$  derived by maximum likelihood, with risk

$$R_{CE} = \mathbb{E}_{(\mathbf{x}, y)} \left\{ \underbrace{-\log \mathbf{p}(y | \mathbf{x})}_{\ell_{CE}} \right\}$$

## URE for cross-entropy loss [TI2019]

$$\overline{R}_{CE} = \mathbb{E}_{(\mathbf{x}, \bar{y})} \left\{ \overbrace{(K-1) \log \mathbf{p}(\bar{y} | \mathbf{x}) - \sum_{k=1}^K \log \mathbf{p}(k | \mathbf{x})}^{\bar{\ell}} \right\}$$

under uniform  $\bar{y}$  (that  $\neq y$ ) assumption

$$\text{URE-CLL: } \min_{\mathbf{p}} \hat{\overline{R}}_{CE}$$

## Issue: URE-CLL Overfits Easily

$$\ell_{CE} = -\log \mathbf{p}(y \mid \mathbf{x})$$

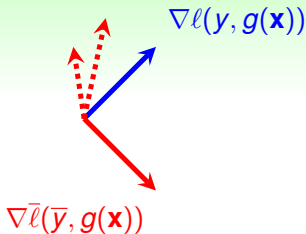
$$\bar{\ell}_{CE} = \underbrace{(K-1) \log \mathbf{p}(\bar{y} \mid \mathbf{x})}_{\text{negative}} - \sum_{k=1}^K \log \mathbf{p}(k \mid \mathbf{x})$$

ordinary risk and URE are very different

- $\ell > 0$ : ordinary risk  $R$  non-negative
- often small  $\mathbf{p}(\bar{y} \mid \mathbf{x})$ :  $\bar{\ell}$  often very negative
- empirically, negative  $\hat{\bar{R}}_{\bar{\ell}}$   
—since only some  $\bar{y}_n$  is observed
- observation: negative empirical URE  $\rightarrow$  overfitting (but why?)

practical remedy NN-URE [TI2019]:  
constrain empirical URE to be non-negative

# Our Contributions



(to be discussed)

an analytical and algorithmic study of URE-CLL, which ...

- constructs a **novel loss-design framework**
- results in **promising empirical performance**
- leads to **novel insights** light on why negative empirical URE causes overfitting

will first describe **key idea**  
behind our proposed framework

# Key Idea: URE on 0/1 instead of $\ell$

## Minimize Complementary 0/1

- goal: minimize  $R_{01}$ , not surrogate  $R_\ell$
- URE of  $R_{01}$ : need

$$\bar{R}_{01} = \mathbb{E}_{(\mathbf{x}, \bar{y})} \bar{\ell}_{01}(\bar{y}, g(\mathbf{x})) = \mathbb{E}_{(\mathbf{x}, y)} \underbrace{\ell_{01}(y, g(\mathbf{x}))}_{\mathbb{I}[y \neq \arg\max_k (g(\mathbf{x}))_k]}$$

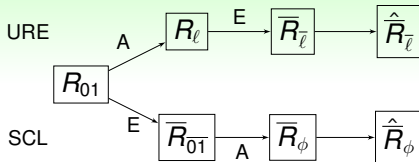
- simple solution:

$$\bar{\ell}_{01}(\bar{y}, g(\mathbf{x})) = \mathbb{I}[\bar{y} = \arg\max_k (g(\mathbf{x}))_k]$$

- intuition: all we need is to discourage  $g(\mathbf{x})$  from predicting  $\bar{y}$   
—minimum likelihood “principle”

Surrogate Complementary Loss (SCL):  
minimize (empirical) surrogate risk of  $\bar{\ell}_{01}$

# Illustrative Difference between URE and SCL



## URE: ripple effect of error

- theoretical motivation [TI2017]
- **estimation step (E)** amplifies **approximation error (A)** in  $\bar{\ell}$

## SCL: “directly” minimize complementary likelihood

- **non-negative surrogate loss  $\phi$**  for  $\bar{\ell}_{01}$  to be minimized
- potentially preventing ripple effect
- **unify previous studies** as different  $\phi$  [XY2018, YK2019]

SCL: swapping **(E)** and **(A)** for loss design



# Example of Avoiding Negative Risk

## Unbiased Risk Estimator (URE)

URE loss  $\bar{\ell}_{CE}$  [TI2019] from  $\ell_{CE}$ ,

$$\bar{\ell}_{CE}(\bar{y}, g(\mathbf{x})) = \underbrace{(K-1) \log \mathbf{p}(\bar{y} | \mathbf{x})}_{\text{negative}} - \sum_{j=1}^K \log \mathbf{p}(j | \mathbf{x})$$

## Surrogate Complementary Loss (SCL)

[YK2019]

$$\phi_{NL}(\bar{y}, g(\mathbf{x})) = -\log(1 - \mathbf{p}(\bar{y} | \mathbf{x}))$$

—a non-negative surrogate of  $\bar{\ell}_{01}$

SCL opens new possibilities  
on studying different  $\phi$

# Experimental Results

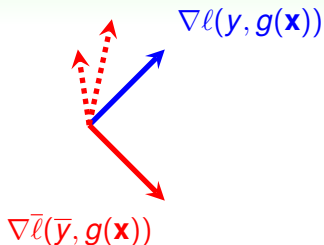
## Models

- 1 Unbiased Risk Estimator (URE) with  $\bar{\ell}_{CE}$  [TI2017]
- 2 Non-Negative Correction of URE (NN-URE) with  $\bar{\ell}_{CE}$  [TI2019]
- 3 Surrogate Complementary Loss (SCL) with exponential  $\phi$  (ours)

Dataset + Model	URE	NN-URE	SCL
MNIST + Linear	0.850	0.818	0.902
MNIST + MLP	0.801	0.867	0.925
CIFAR10 + ResNet	0.109	0.308	0.492
CIFAR10 + DenseNet	0.291	0.338	0.544

SCL is **significantly better**  
than URE and NN-URE

# Analysis Using Gradients



## Gradient Direction of URE

- **very diverse directions** on each  $\bar{y}$  to maintain unbiasedness
- **low correlation** to the target  $\ell_{01}$

## Gradient Direction of SCL

- targets towards **minimum likelihood** objective
- **higher correlation** to the target  $\bar{\ell}_{01}$

will quantify this  
with **bias-variance decomposition**

# Gradient Estimation Error

## Gradient Estimation

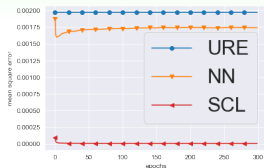
- 1 ordinary gradient  $\mathbf{f} = \nabla \ell(y, g(\mathbf{x}))$
- 2 complementary gradient  $\mathbf{c} = \nabla \bar{\ell}(\bar{y}, g(\mathbf{x}))$
- 3 expected complementary gradient  $\mathbf{h} = \text{average of } \mathbf{c} \text{ over } \bar{y}$

## Bias-Variance Decomposition

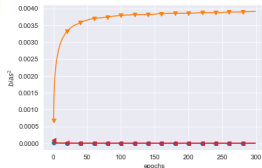
$$\begin{aligned}\text{MSE} &= \mathbb{E}[(\mathbf{f} - \mathbf{c})^2] \\ &= \underbrace{\mathbb{E}[(\mathbf{f} - \mathbf{h})^2]}_{\text{Bias}^2} + \underbrace{\mathbb{E}[(\mathbf{h} - \mathbf{c})^2]}_{\text{Variance}}\end{aligned}$$

is **unbiased** risk/gradient estimator good?

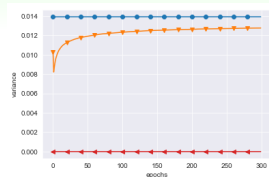
# Bias-Variance Tradeoff on Gradient



(a) MSE



(b) Bias<sup>2</sup>



(c) Variance

	Bias	Variance	MSE
URE	0	Big	Big
NN-URE	Big	Smaller	Big
SCL	Small	Smallest	Small

SCL reduces variance of URE while introducing small bias

## Some Issues for Mathematicians

minimize  $\bar{\ell}_{01}$ —hypothesis that **least matches** complementary data:

is this **minimum likelihood** principle well-justified? **Not yet.**

bias-variance decomposition of gradient based on **empirical findings**:

is there a theoretical guarantee to play with the trade-off? **Not yet.**

current results mostly based on **uniform** complementary labels:

do we understand the assumptions to make CLL ‘learnable’? **Not yet.**

some (but not all) answered in the **next paper**

# Mini-Summary

## Explain Overfitting of URE

- URE only **expected** to do well
- fixed CLs cause **high variance** (hence overfitting)

## Surrogate Complementary Loss (SCL)

- **avoids negative risk** issue by design
- **minimum likelihood** principle

## Experiment Results

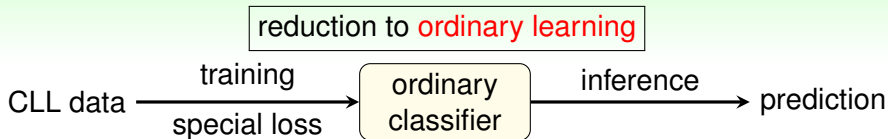
- SCL **significantly outperforms** others
- trade small gradient bias for **lower variance**

“traditional” statistics tools  
can be useful for **modern problem**

Wei-I Lin and Hsuan-Tien Lin.  
Reduction from complementary-label learning  
to probability estimates. PAKDD 2023  
Best Paper Runner-up Award.



# Reflection on CLL Model Design



## Inference: Easy

simply  $\operatorname{argmax}_k (g(\mathbf{x}))_k$

## Training: Challenging

- indirect estimation from CLs
- prone to overfitting
- mostly only tested on deep models

can we make training easier?

# Our Contributions

$$R_{01}(\text{dec}(\bar{g}, L_1)) \leq \frac{4\sqrt{2}}{\gamma} \sqrt{R(\bar{g}, \ell_{KL})}$$

(to be discussed)

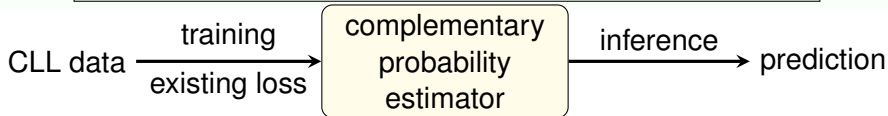
a principled study of CLL Model Design, which ...

- promotes a **novel reduction framework**
- leads to **sound explanations** on several existing models
- results in **promising empirical performance** in some scenarios

again, will first describe **key idea**  
behind our proposed framework

# Key Idea: Complementary Probability Estimation

reduction to **complementary probability estimation (CPE)**



## Training: Easy

learn complementary probability estimates  $\bar{g}(\mathbf{x})$  with CLs

- **direct learning** from CLs
- many **existing deep/non-deep models**
- **easy to validate** too

inference: **how (under what assumption)?**

# Assumption: How are CLs Generated?

## uniform assumption

$$P(\bar{y} | y) = \frac{1}{K-1} \mathbb{I}[\bar{y} \neq y]$$

## conditional generation assumption

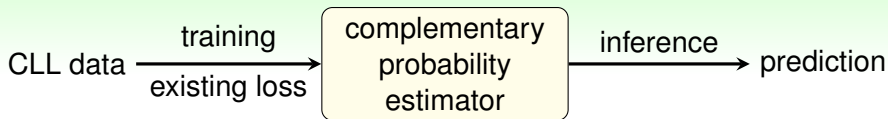
$$P(\bar{y} | \mathbf{x}, y) = P(\bar{y} | y) = T_{y, \bar{y}}$$

e.g. transition matrix

$$T = \begin{bmatrix} 0 & 0.3 & 0.3 & 0.4 \\ 0.4 & 0 & 0.3 & 0.3 \\ 0.3 & 0.4 & 0 & 0.3 \\ 0.4 & 0.3 & 0.3 & 0 \end{bmatrix}$$

how to do inference **with known  $T$**  after CPE?

# Nearest Transition Vector Decoder



$$T = \begin{bmatrix} 0 & 0.3 & 0.3 & 0.4 \\ 0.4 & 0 & 0.3 & 0.3 \\ 0.3 & 0.4 & 0 & 0.3 \\ 0.4 & 0.3 & 0.3 & 0 \end{bmatrix}$$

looks like  $y = 1$  if  $\bar{g}(\mathbf{x}) = [0.03, 0.27, 0.25, 0.45]$

proposed nearest-transition-vector decoder  
for inference:

$$\text{dec}(\bar{g}, d): \mathbf{x} \rightarrow \underset{y \in [K]}{\text{argmin}} d(\bar{g}(\mathbf{x}), T_y)$$

# Theoretical Guarantee of CPE

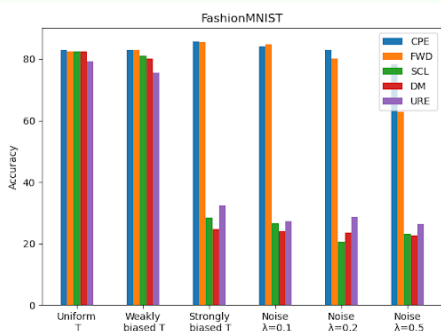
When using  $d = L_1$  distance,

$$R_{01}(\text{dec}(\bar{g}, L_1)) \leq \frac{4\sqrt{2}}{\gamma} \sqrt{R_{KL}(\bar{g})}$$

- $\gamma$ : **minimum  $L_1$  distance** between rows of transition vectors
- **smaller CPE error** (KL divergence)  $\rightarrow$  **smaller  $R_{01}$**
- explains **SCL as special case** of  $L_1$  decoding under uniform assumption
- can be used to **validate with CLs only**

**other distance measures possible**  
(but we did not study much)

# Experimental Results



## Models

- 1 Unbiased Risk Estimator (URE) [TI2017]
- 2 Discriminative model (DM\*) [YG2021]
- 3 Surrogate Complementary Loss (SCL\*, our previous work)
- 4 Forward (FWD\*) [XY2018]
- 5 Complementary Probability Estimator (CPE, ours)

CPE better than others and special cases (\*),  
especially with noisy  $T$

## Some Issues for Mathematicians Revisited

minimize  $\bar{\ell}_{01}$ —hypothesis that **least matches** complementary data:

is **minimum likelihood** principle well-justified? **Yes, special case of CPE.**

bias-variance decomposition of gradient based on **empirical findings**:

is there a theoretical guarantee to play with the trade-off? **Not yet.**

current results mostly based on **uniform** complementary labels:

the assumptions to make CLL ‘learnable’? **any known  $T$  with  $\gamma > 0$ .**

some answered in the **this paper**



# Mini-Summary

## Explain SCL (and Others)

- via a **different reduction** route

## Complementary Probability Estimation (CPE)

- **estimate complementary probabilities** during training (easy)
- **nearest transition vector decoding** (theoretical guarantees)

## Experiment Results

- CPE **outperforms (?)** others
- potential for **noisy CLL and CL-only validation**

now, is CLL **realistic?**

Hsiu-Hsuan Wang, Tan-Ha Mai,  
Nai-Xuan Ye, Wei-I Lin, Hsuan-Tien Lin.

CLImage: Human-Annotated Datasets for  
Complementary-Label Learning. arXiv:2305.08295

## Recall: Assumptions in CLL Model Design

noise-free assumption

$$P(\bar{y} = y \mid y) = 0$$

uniform assumption

$$P(\bar{y} \mid y) = \frac{1}{K-1} \mathbb{I}[\bar{y} \neq y]$$

conditional generation assumption

$$P(\bar{y} \mid \mathbf{x}, y) = P(\bar{y} \mid y) = T_{y, \bar{y}}$$

do they **hold in reality**?

# CLImage: Protocol for Collecting CL from Annotators

air-  
plane

auto-  
mobile

bird

cat

deer

dog

frog

horse

ship

truck

Randomly pick four classes

air-  
plane

auto-  
mobile

bird

cat

deer

dog

frog

horse

ship

truck

Ask the annotators to select any incorrect label



auto-  
mobile

ship

bird

frog

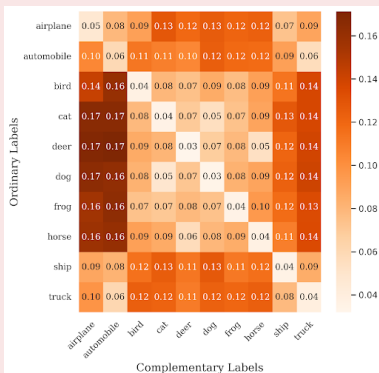
(courtesy of Wei-I Lin)

play here: [https://github.com/ntucllab/CLImage\\_Dataset/](https://github.com/ntucllab/CLImage_Dataset/)

# Analysis of Collected Data

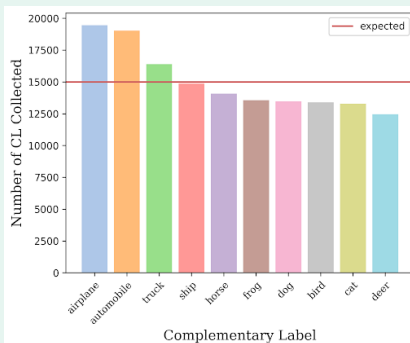
is it noise-free?

no (not surprisingly), and it affects performance significantly



is it uniform?

no (not surprisingly), and it affects performance a bit



more studies on noisy CLL is needed

# An Insider Secret

## CLImage

- CLCIFAR10
- CLCIFAR20 (20 meta-classes)
- CLMicroImageNet10 (10 random classes)
- CLMicroImageNet20 (20 random classes)

—why **only data of 10 or 20 classes**?

## Truth

tried CIFAR100 **but failed**

- **higher accuracy than random guess**
- much lower than ordinary classification, **even after noise cleaning**

pure CLL **overly weak** and may not be realistic

# Summary (Finally)

## Surrogate Complementary Loss

run URE **before doing surrogate** instead

## Complementary Probability Estimation

consider **probability estimation on CLs** instead

## CLImage

attempt to benchmark how realistic CLL is, with a **dataset collection** and a library in its beta version

<https://github.com/ntucllab/libcll>

**Thank you and all my students/collaborators!**

# References

- [CE2008] Learning classifiers from only positive and unlabeled data, KDD 2008
- [NN2013] Learning with noisy labels, NeurIPS 2013
- [TI2017] Learning from complementary labels, NeurIPS 2017
- [XY2018] Learning with biased complementary labels, ECCV 2018
- [TI2019] Complementary-Label Learning for Arbitrary Losses and Models, ICML 2019
- [YK2019] NLNL: Negative learning for noisy labels, ICCV 2019
- [YG2021] Discriminative complementary-label learning with weighted loss, ICML 2021
- [QD2024] Boosting Semi-Supervised Learning with Contrastive Complementary Labeling, Neural Networks 2024