

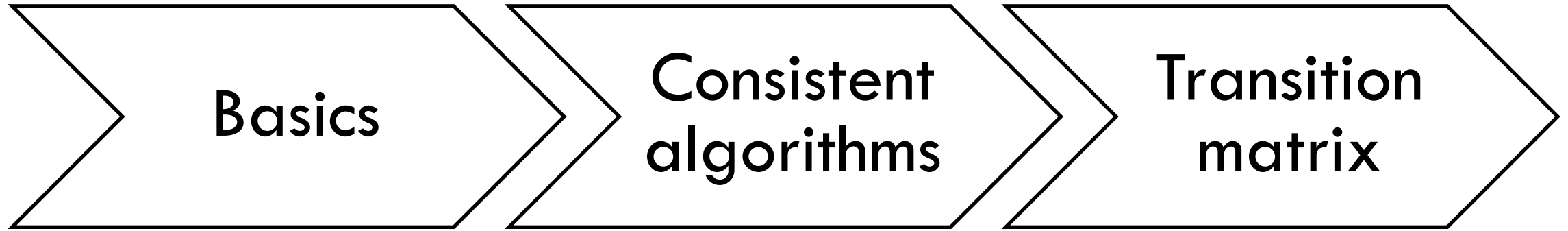
# Learning under Noisy Supervision

## Part II: Statistical Learning under Noisy Supervision

Tongliang Liu

Trustworthy Machine Learning Lab  
School of Computer Science  
University of Sydney

# Structure



# Structure



**Basics**

Consistent  
algorithms

Transition  
matrix

# Learning without noisy labels

## Problem setup:

Data:  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim D^n$ .

Aim: Learn a **classifier**  $f \in F$ , such that  $\forall (\mathbf{x}, y) \sim D$ ,  
 $f(\mathbf{x})$  is a good prediction for  $y$ .

What is the best classifier we can obtain?  
w.r.t. accuracy

To measure the **accuracy**, we define **loss function**

$$\ell(\mathbf{x}, y), f \mapsto \ell(y, f(\mathbf{x})) \in \mathbb{R}.$$

For example, 0-1 loss:  $\mathbf{1}(y \neq \text{sign}(f(\mathbf{x})))$ .

The **best classifier** should be the one that has the minimum loss on **all the possible data** from the domain.

# Expected risk, Bayes classifier

Theoretically,

The expected risk:

$$R_{D,0-1}(f) = \mathbb{E}_{(X,Y) \sim D}[\mathbf{1}(Y \neq \text{sign}(f(X)))].$$

Bayes classifier:

$$f_{\rho}(\mathbf{x}) = \operatorname{argmin}_f R_{D,0-1}(f) = \operatorname{argmax}_y P(Y = y | X = \mathbf{x}).$$

Restricted Bayes classifier:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} R_{D,0-1}(f).$$

**Empirically,**

In reality, we can only observe a sample of data

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim D^n.$$

We approximate the expected risk  $R(f)$  via **the**

**empirical risk:**  $\hat{R}_{D,\ell}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)).$

We minimize the empirical risk to find a predictor:

$$f_n = \arg \min_{f \in \mathcal{F}} \hat{R}_{D,\ell}(f).$$



Statistically consistent classifier [1,2]:

With high probability, as  $n \rightarrow \infty$ ,  
we have:  $R_{D,\ell}(f_n) \rightarrow R_{D,\ell}(f^*)$ .

[1] Mohri et al. *Foundations of machine learning*. MIT press, 2018.

[2] Devroye, et al. *A probabilistic theory of pattern recognition*. Vol. 31. Springer Science & Business Media, 2013.

Aim:

Designing algorithms whose outputs will approach  
 $f_{\rho}(\boldsymbol{x}) = \arg \max_y P(Y = y|X = \boldsymbol{x}).$

# Structure



Basics

**Consistent  
algorithms**

Transition  
matrix

# Learning with noisy labels

Noisy sample:  $\tilde{S} = \{(\mathbf{x}_1, \tilde{y}_1), \dots, (\mathbf{x}_n, \tilde{y}_n)\} \sim \tilde{D}^n$ ,  
where  $\tilde{y}$  stands for noisy labels and  $\tilde{D}$  the noisy distribution.

What is the best classifier we can learn?

Can we approach  $f_\rho(\mathbf{x}) = \arg \max_y P(Y = y | X = \mathbf{x})$ ?

# Learning with noisy labels

One category: [extracting confident examples](#).

SOTA, e.g., Co-teaching [3]; Joint Optim [4].

Another category: [label-noise learning](#) [5].

Methodology, i.e., statistically consistent algorithms.

[3] Han, Bo, et al. "Co-teaching: Robust training of deep neural networks with extremely noisy labels." *NeurIPS* 2018.

[4] Tanaka, Daiki, et al. "Joint optimization framework for learning with noisy labels." *CVPR* 2018.

[5] Xia, Xiaobo, et al. "Are anchor points really indispensable in label-noise learning?." *NeurIPS* 2019.

# Learning with noisy labels

One category: extracting confident examples.

SOTA, e.g., Co-teaching [3]; Joint Optim [4].

Another category: **label-noise learning** [5].

Methodology, i.e., statistically consistent algorithms.

[3] Han, Bo, et al. "Co-teaching: Robust training of deep neural networks with extremely noisy labels." *NeurIPS* 2018.

[4] Tanaka, Daiki, et al. "Joint optimization framework for learning with noisy labels." *CVPR* 2018.

[5] Xia, Xiaobo, et al. "Are anchor points really indispensable in label-noise learning?." *NeurIPS* 2019.

# Model label noise

Transition matrix:

$$\begin{bmatrix} P(\tilde{Y} = 1|Y = 1, \mathbf{x}) & \cdots & P(\tilde{Y} = 1|Y = C, \mathbf{x}) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1, \mathbf{x}) & \cdots & P(\tilde{Y} = C|Y = C, \mathbf{x}) \end{bmatrix}.$$

# Transition matrix

$$\begin{bmatrix} P(\tilde{Y} = 1|\mathbf{x}) \\ \vdots \\ P(\tilde{Y} = C|\mathbf{x}) \end{bmatrix} = \underbrace{\begin{bmatrix} P(\tilde{Y} = 1|Y = 1, \mathbf{x}) & \cdots & P(\tilde{Y} = 1|Y = C, \mathbf{x}) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1, \mathbf{x}) & \cdots & P(\tilde{Y} = C|Y = C, \mathbf{x}) \end{bmatrix}}_{T^{\top}(\mathbf{x})} \begin{bmatrix} P(Y = 1|\mathbf{x}) \\ \vdots \\ P(Y = C|\mathbf{x}) \end{bmatrix}$$



# Why called “label-noise learning”?

- Label-noise learning [5]
- Learning with noisy labels [6]

[5] Xia, Xiaobo, et al. “Are anchor points really indispensable in label-noise learning?.” *NeurIPS* 2019.

[6] Natarajan, Nagarajan, et al. “Learning with noisy labels.” *NeurIPS* 2013.

# Model Label Noise

(1) Random Classification Noise (RCN) [7]:

$$\rho_{\tilde{Y},Y}(X) = P(\tilde{Y}|Y, X) = P(\tilde{Y}|Y) = \rho, \forall Y \neq \tilde{Y}.$$

(2) Class-conditional Noise (CCN) [6]:

$$\rho_{\tilde{Y},Y}(X) = P(\tilde{Y}|Y, X) = P(\tilde{Y}|Y).$$

(3) Instance-dependent Noise (IDN) [8,9]:

$$\rho_{\tilde{Y},Y}(X) = P(\tilde{Y}|Y, X).$$

[7] Angluin, Dana, and Philip Laird. "Learning from noisy examples." *Machine Learning* 2.4: 343-370, 1988.

[8] Cheng, Jiacheng, et al. "Learning with bounded instance and label-dependent label noise." *ICML* 2020.

[9] Berthon, Antonin, et al. "Confidence scores make instance-dependent label-noise learning possible." *ICML*, 2021

# Model Label Noise

(1) Random Classification Noise (RCN) [7]:

$$\rho_{\tilde{Y},Y}(X) = P(\tilde{Y}|Y, X) = P(\tilde{Y}|Y) = \rho, \forall Y \neq \tilde{Y}.$$

(2) Class-conditional Noise (CCN) [6]:

$$\rho_{\tilde{Y},Y}(X) = P(\tilde{Y}|Y, X) = P(\tilde{Y}|Y).$$

(3) Instance-dependent Noise (IDN) [8,9]:

$$\rho_{\tilde{Y},Y}(X) = P(\tilde{Y}|Y, X).$$

[7] Angluin, Dana, and Philip Laird. "Learning from noisy examples." *Machine Learning* 2.4: 343-370, 1988.

[8] Cheng, Jiacheng, et al. "Learning with bounded instance and label-dependent label noise." *ICML* 2020.

[9] Berthon, Antonin, et al. "Confidence scores make instance-dependent label-noise learning possible." *ICML*, 2021

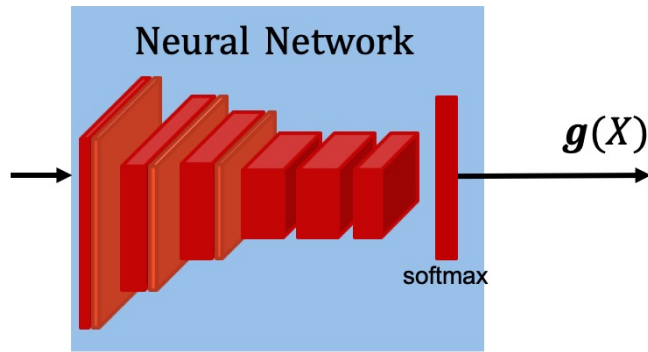
# Class-conditional Noise (CCN)

The loss correction method:  
Modify  $\ell$  to be  $\tilde{\ell}$  such that

$$\mathbb{E}_{(X, \tilde{Y}) \sim \tilde{D}} [\tilde{\ell}(f(X), \tilde{Y})] = \mathbb{E}_{(X, Y) \sim Y} [\ell(f(X), Y)]$$

By exploiting the model of label noise:

$$\begin{bmatrix} P(\tilde{Y} = 1 | \mathbf{x}) \\ \vdots \\ P(\tilde{Y} = C | \mathbf{x}) \end{bmatrix} = \begin{bmatrix} P(\tilde{Y} = 1 | Y = 1) & \cdots & P(\tilde{Y} = 1 | Y = C) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C | Y = 1) & \cdots & P(\tilde{Y} = C | Y = C) \end{bmatrix} \begin{bmatrix} P(Y = 1 | \mathbf{x}) \\ \vdots \\ P(Y = C | \mathbf{x}) \end{bmatrix}$$



$$\begin{bmatrix} P(\tilde{Y} = 1|\mathbf{x}) \\ \vdots \\ P(\tilde{Y} = C|\mathbf{x}) \end{bmatrix} = \begin{bmatrix} P(\tilde{Y} = 1|Y = 1) & \cdots & P(\tilde{Y} = 1|Y = C) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1) & \cdots & P(\tilde{Y} = C|Y = C) \end{bmatrix} \begin{bmatrix} P(Y = 1|\mathbf{x}) \\ \vdots \\ P(Y = C|\mathbf{x}) \end{bmatrix}$$

## Importance reweighting [11]:

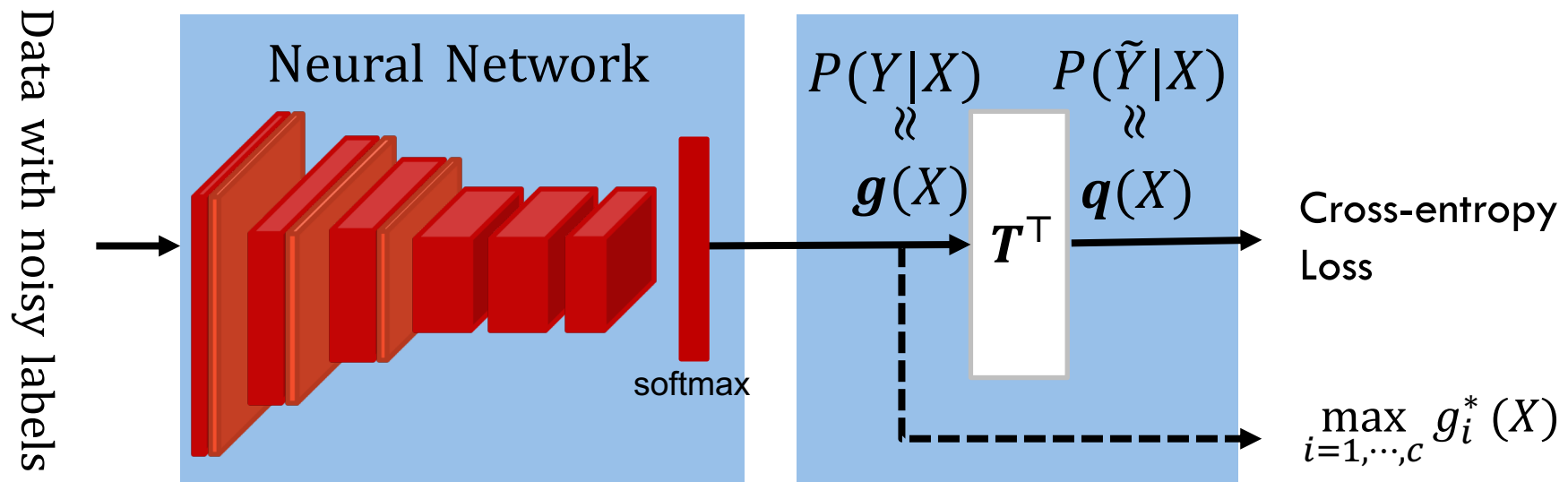
$$\tilde{\ell}_{ir}(f(\mathbf{x}), y) = \frac{P(\mathbf{x}, y)}{\tilde{P}(\mathbf{x}, y)} \ell(f(\mathbf{x}), y) = \frac{\mathbf{g}_y(\mathbf{x})}{(T^\top \mathbf{g})_y(\mathbf{x})} \ell(f(\mathbf{x}), y),$$

where  $f(\mathbf{x}) = \arg \max_{j \in \{1, \dots, C\}} \mathbf{g}_j(\mathbf{x})$ .

$$\text{Thus, } \mathbb{E}_{(X, \tilde{Y}) \sim \tilde{D}} [\tilde{\ell}_{ir}(f(X), \tilde{Y})] = \mathbb{E}_{(X, Y) \sim D} [\ell(f(X), Y)]$$

$$\begin{bmatrix} P(\tilde{Y} = 1|\mathbf{x}) \\ \vdots \\ P(\tilde{Y} = C|\mathbf{x}) \end{bmatrix} = \begin{bmatrix} P(\tilde{Y} = 1|Y = 1) & \cdots & P(\tilde{Y} = 1|Y = C) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1) & \cdots & P(\tilde{Y} = C|Y = C) \end{bmatrix} \begin{bmatrix} P(Y = 1|\mathbf{x}) \\ \vdots \\ P(Y = C|\mathbf{x}) \end{bmatrix}$$

## Forward correction [12]:



# A summary of consistent algorithms

- Many methods for dealing with noisy labels  
**Loss correction**, Sample selection, label correction, ...
- Model label noise  
Random Classification Noise (RCN)  
**Class-conditional Noise (CCN)**  
Instance-dependent Noise (IDN)
- Loss correction methods  
Importance reweighting (risk-consistent), forward correction (classifier-consistent)

# Structure



Basics

Consistent  
algorithms

**Transition  
matrix**



# How to estimate the transition matrix

Given the noisy data

$$\tilde{S} = \{(\mathbf{x}_1, \tilde{y}_1), \dots, (\mathbf{x}_n, \tilde{y}_n)\} \sim \tilde{D}.$$

How to estimate the transition matrix  $T$ ?

# Definition

If  $P(Y = i | \mathbf{x}^i) = 1$ , then  $\mathbf{x}^i$  is called the **anchor point** for the  $i$ -th class.

# Anchor point assumption

$$\begin{bmatrix} P(\tilde{Y} = 1|X) \\ \vdots \\ P(\tilde{Y} = C|X) \end{bmatrix} = \underbrace{\begin{bmatrix} P(\tilde{Y} = 1|Y = 1) & \cdots & P(\tilde{Y} = 1|Y = C) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1) & \cdots & P(\tilde{Y} = C|Y = C) \end{bmatrix}}_T \begin{bmatrix} P(Y = 1|X) \\ \vdots \\ P(Y = C|X) \end{bmatrix}$$

$$\begin{bmatrix} P(\tilde{Y} = 1|X = \mathbf{x}^1) \\ \vdots \\ P(\tilde{Y} = C|X = \mathbf{x}^1) \end{bmatrix} = \begin{bmatrix} P(\tilde{Y} = 1|Y = 1) & \cdots & P(\tilde{Y} = 1|Y = C) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1) & \cdots & P(\tilde{Y} = C|Y = C) \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} P(\tilde{Y} = 1|X = \mathbf{x}^i) \\ \vdots \\ P(\tilde{Y} = C|X = \mathbf{x}^i) \end{bmatrix} = \begin{bmatrix} P(\tilde{Y} = 1|Y = i) \\ \vdots \\ P(\tilde{Y} = C|Y = i) \end{bmatrix}$$

# How to find anchor points

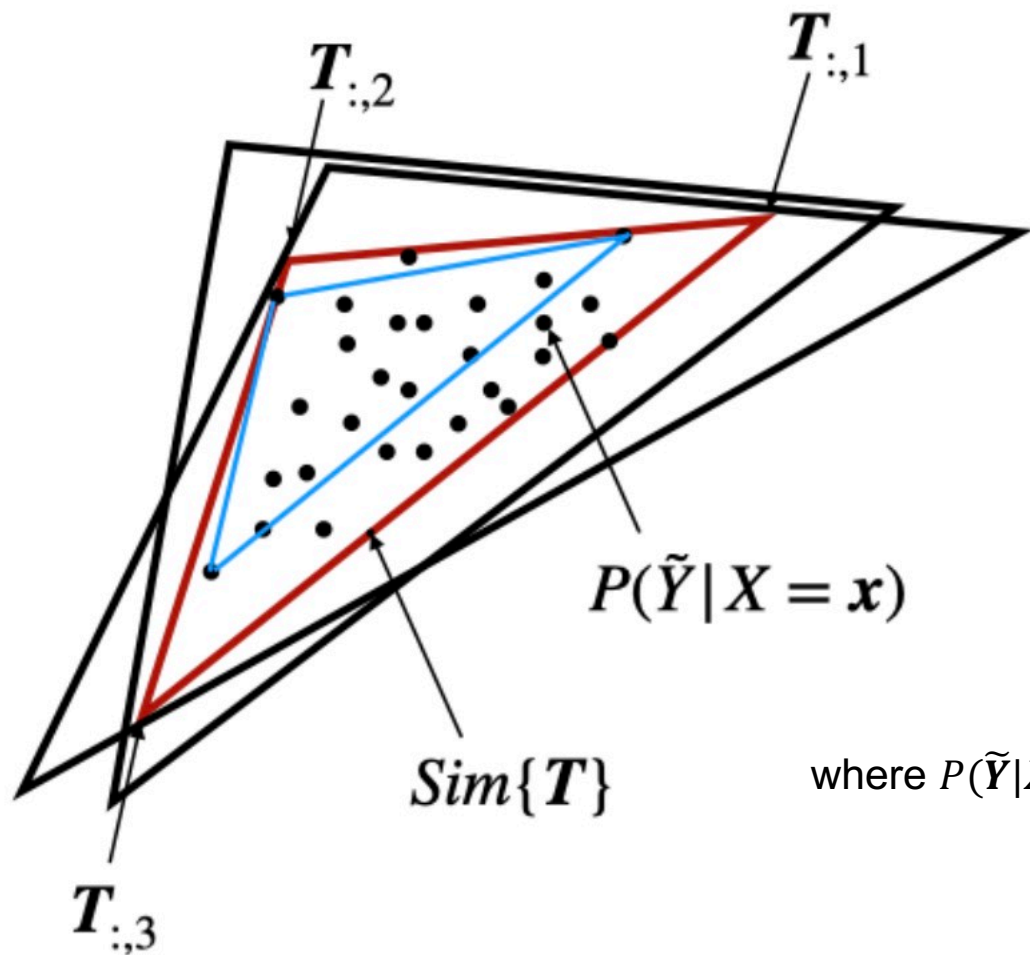
Binary classification, find the anchor points:

$$\mathbf{x}^y = \operatorname{argmax}_{\mathbf{x} \in X} P(\tilde{Y} = y | \mathbf{x}).$$

Multi-classification, approximate the anchor points for multi-class learning:

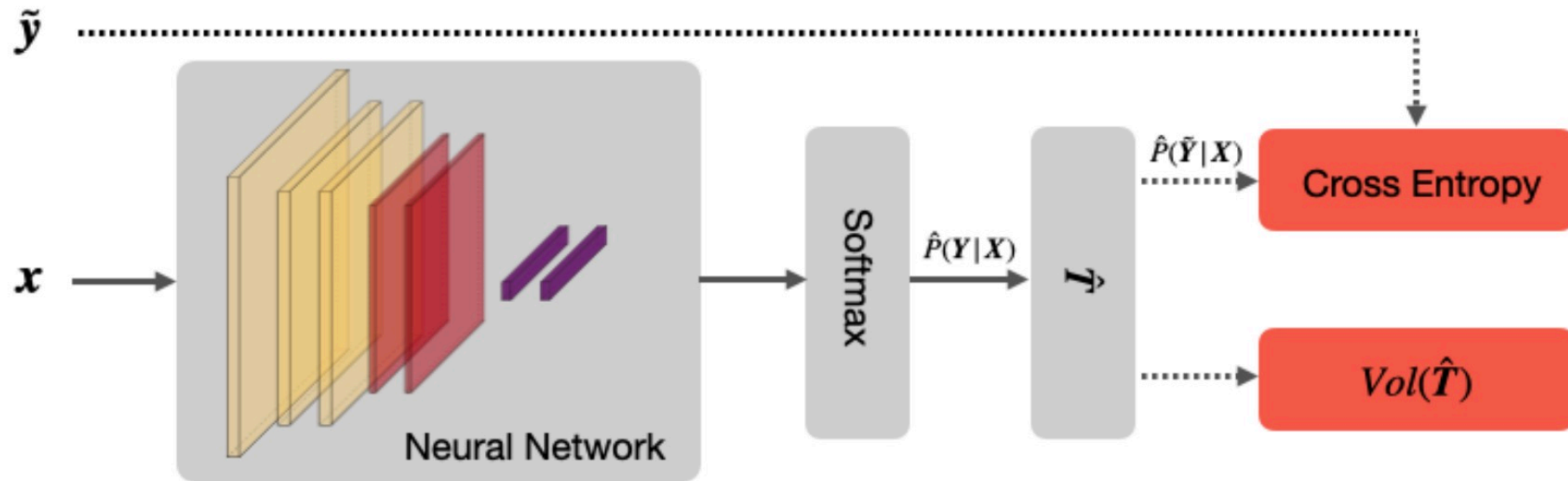
$$\mathbf{x}^y \approx \operatorname{argmax}_{\mathbf{x} \in X} P(\tilde{Y} = y | \mathbf{x}).$$

# Sufficiently scattered assumption vs anchor point assumption

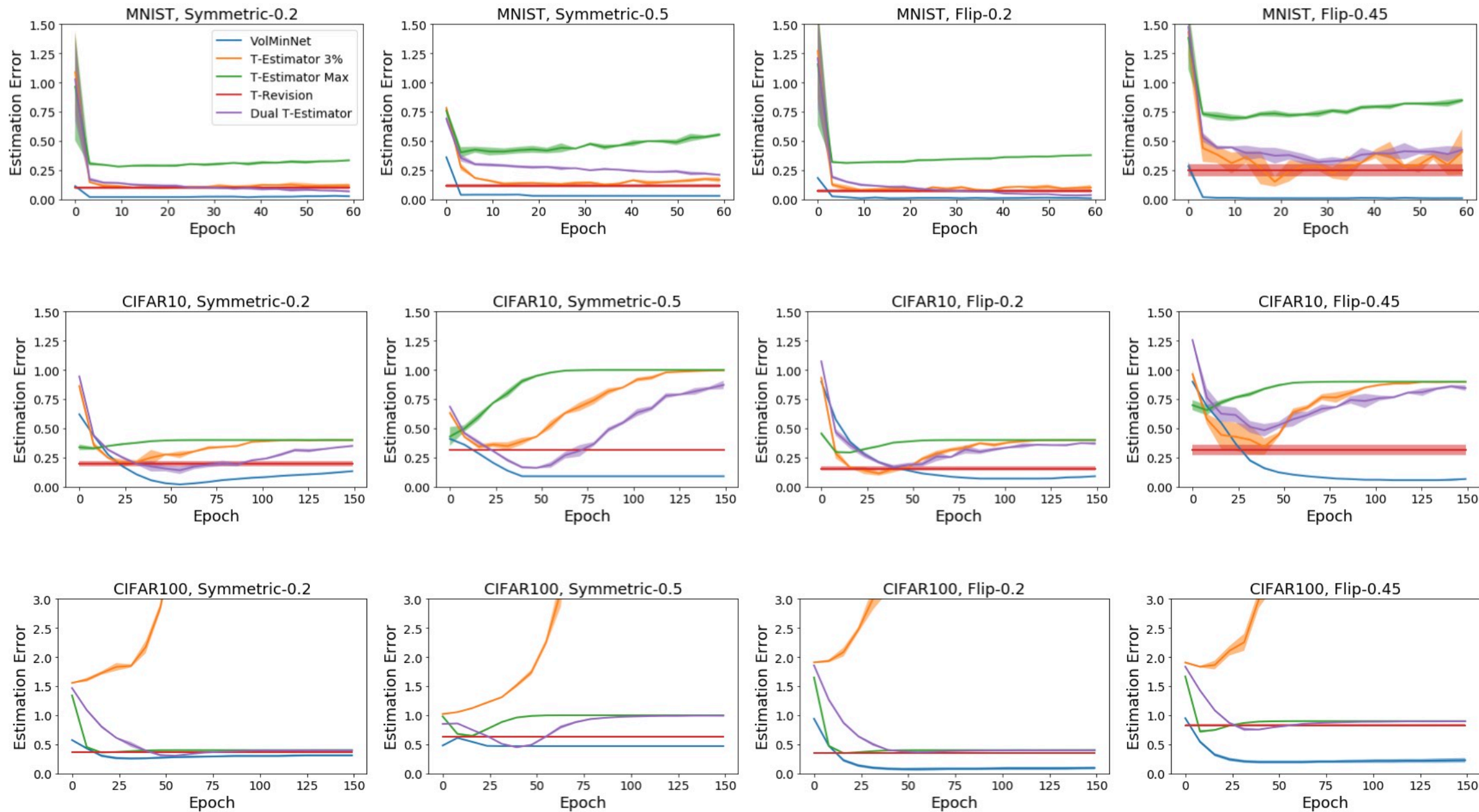


$$\text{where } P(\tilde{Y}|X) = \begin{bmatrix} P(\tilde{Y} = 1|X) \\ \vdots \\ P(\tilde{Y} = C|X) \end{bmatrix} = \underbrace{\begin{bmatrix} P(\tilde{Y} = 1|Y = 1) & \cdots & P(\tilde{Y} = 1|Y = C) \\ \vdots & \ddots & \vdots \\ P(\tilde{Y} = C|Y = 1) & \cdots & P(\tilde{Y} = C|Y = C) \end{bmatrix}}_T \begin{bmatrix} P(Y = 1|X) \\ \vdots \\ P(Y = C|X) \end{bmatrix}$$

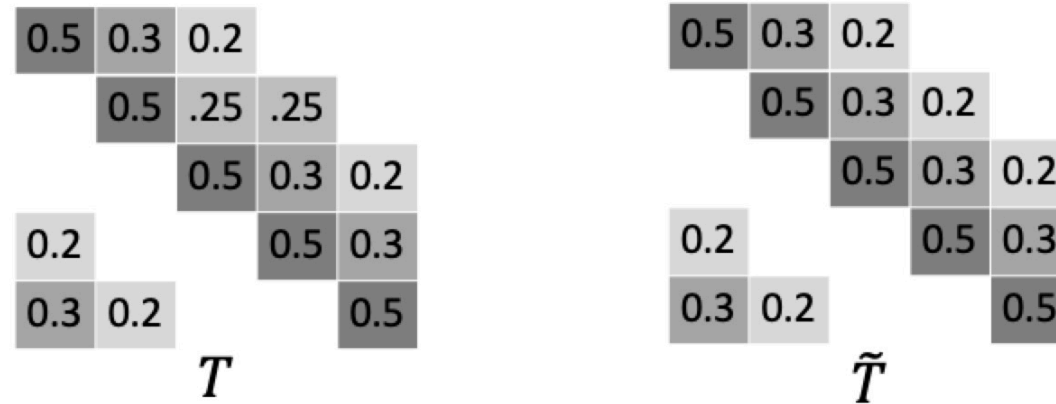
# VolMinNet [14]



$$\begin{aligned} & \min_{\hat{T} \in \mathbb{T}} \text{vol}(\hat{T}) \\ & \text{s. t. } \hat{T} h_{\theta} = P(\tilde{Y}|X) \end{aligned}$$



# $T$ revision [15]



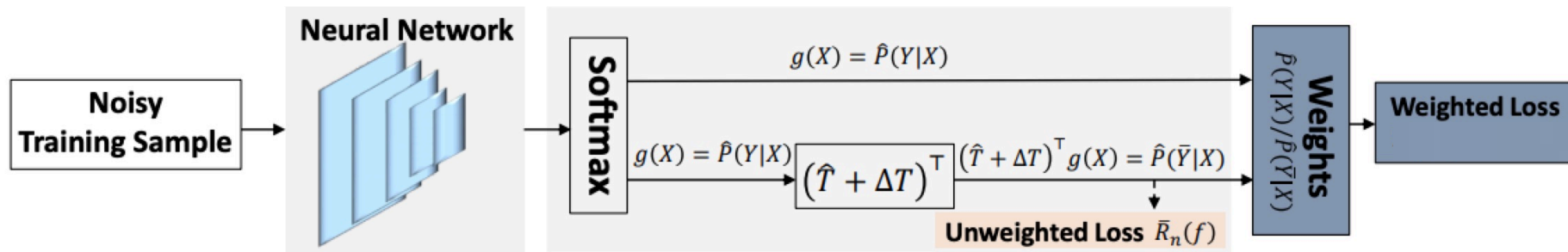
If  $P(\tilde{Y}|X = \mathbf{x}) = [0.141; 0.189; 0.239; 0.281; 0.15]$ ,

then,  $P(Y|X = \mathbf{x}) = (T^\top)^{-1}P(\tilde{Y}|X = \mathbf{x}) =$   
 $[0.15; 0.28; 0.25; 0.3; 0.02]$ .

$P(Y|X = \mathbf{x}) = (\tilde{T}^\top)^{-1}P(\tilde{Y}|X = \mathbf{x})$   
 $= [0.1587; 0.2697; 0.2796; 0.2593; 0.0325]$ .

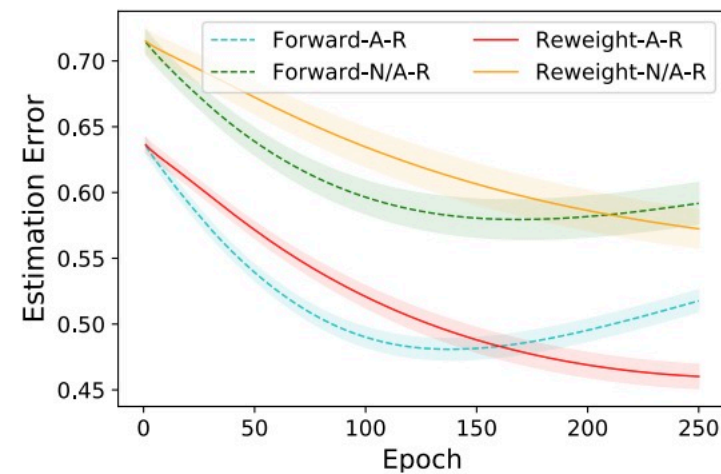
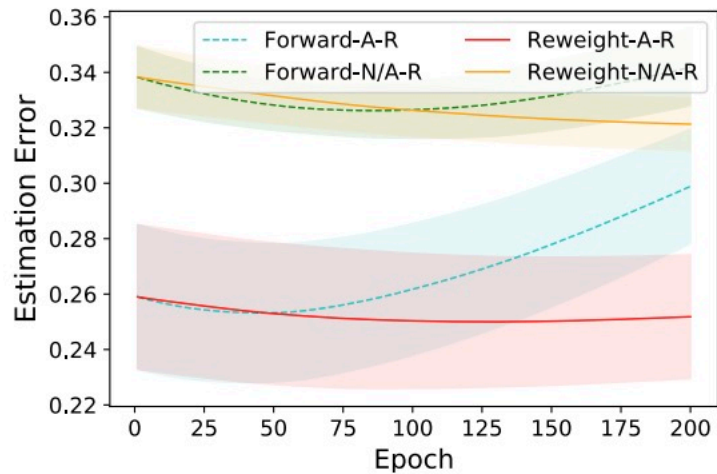
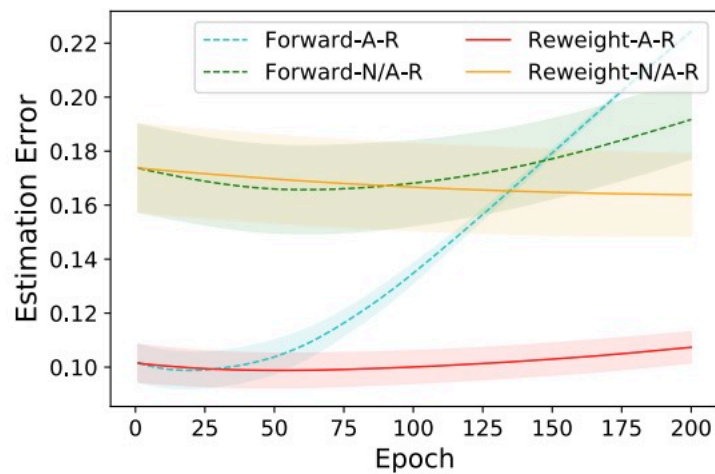
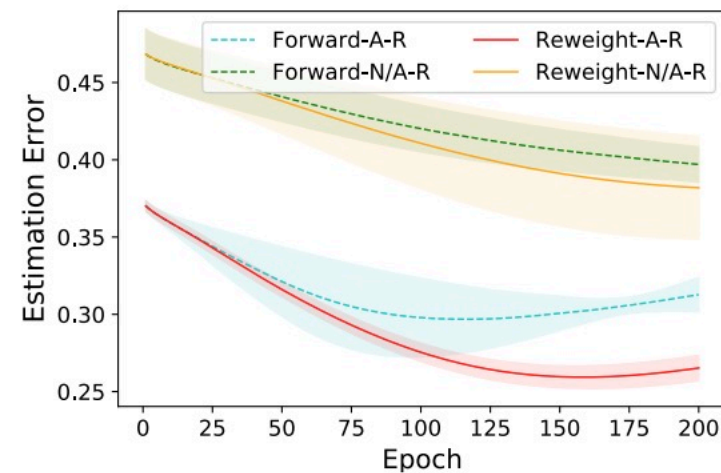
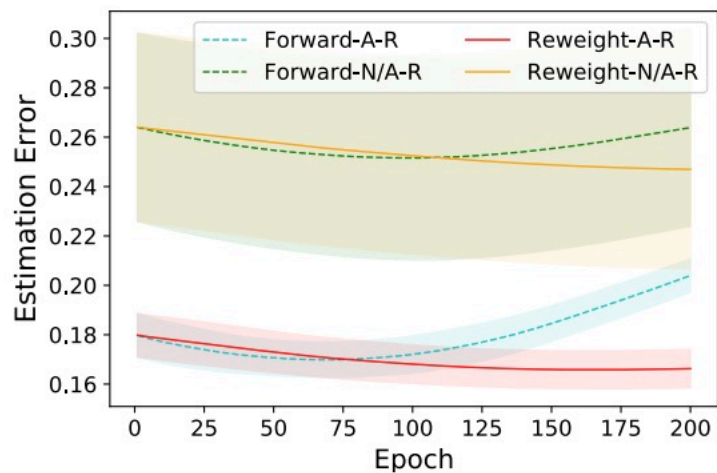
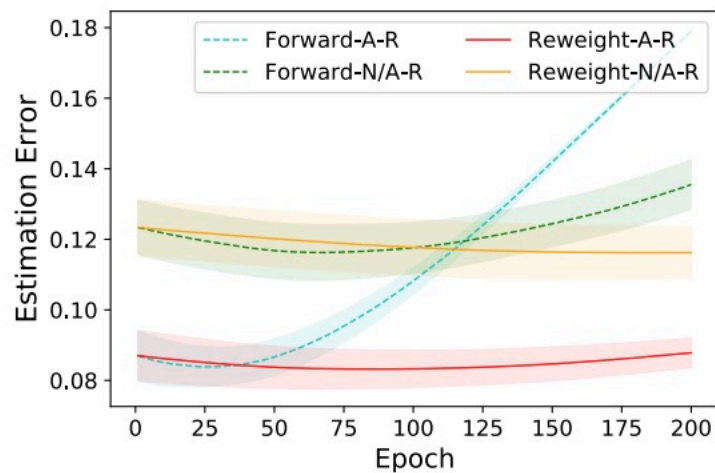


# $T$ revision [15]



$$\text{Weighted loss} = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{g}_{\tilde{y}_i}(\mathbf{x}_i)}{(\mathbf{T}^\top \mathbf{g})_{\tilde{y}_i}(\mathbf{x}_i)} L(f(\mathbf{x}_i), \tilde{y}_i),$$

where  $f(\mathbf{x}) = \operatorname{argmax}_{i \in \{1, \dots, C\}} \mathbf{g}_i(\mathbf{x})$ .



(a) *MNIST*

(b) *CIFAR-10*

(c) *CIFAR-100*

# A summary of estimating transition matrix

- How to estimate the transition matrix given only noisy data?  
Method:  $T$  estimator (by exploiting anchor points)
- How about if there is no anchor points?  
Method: VolMinNet (using the sufficiently scattered assumption)
- How to deal with poorly estimated transition matrix  
Method:  $T$  revision (revising the matrix by using a slack variable)

# Conclusion and future directions

## ➤ Conclusion

- Statistically consistent algorithms: the classifier learned by using noisy data will converge to the optimal one defined by using clean data
- Statistically consistent algorithms are robust to the data distribution and label noise type
- Modelling the label noise and estimating the transition matrix are cores in label-noise learning

## ➤ Future directions

- Design effectively loss correction methods for deep learning
- How to address the finite/small sample problem
- How to use a small set of clean data to better estimate the transition matrix
- How to model and estimate the instance-dependent label noise (IDN)