# Learning with Noisy Supervision, Part 5: Beyond Class-Conditional Noise

#### Gang Niu

#### Research Scientist Imperfect Information Learning Team RIKEN Center for Advanced Intelligence Project

IJCAI 2021 Tutorial August 20, 2021

Motivation	IDN	MCD	Conclusions

# Outline



- 2 Instance-dependent noise (IDN)
- 3 Mutually contaminated distributions (MCD)

## 4 Conclusions

Motivation	IDN	MCD	Conclusions
•0000			

# Class-conditional noise (CCN) model (conservation)

- All models are wrong, but some are useful (Box, "Science and Statistics", JASA 1976)
  - Following the law of total probability,  $p(\tilde{y} \mid x) = \sum_{y} p(\tilde{y} \mid x, y) p(y \mid x)$
  - Assume  $p(\tilde{y} \mid x, y) = p(\tilde{y} \mid y)$

i.e., the corruption  $y \to \tilde{y}$  is instance-independent and class-conditional

- Equivalently, using transition matrix T where  $[T]_{i,j} = p(\tilde{y} = j \mid y = i)$ 

Motivation	IDN	MCD	Conclusions
0000			

## Toy example: symmetric noise on Gaussian mixture



# Label noise is (almost) everywhere in industry

#### Active label collection Make Money **Get Results** by working on HITs Ask workers to complete HETs - Human Intelligence Tasks - and get results using Wechanical Turk. <u>Register New</u> Numan Intelligence Tasks - are individual tasks that work on, Find Hills now. As a Mechanical Turk Requester you As a Mechanical Turk Morker your ave access to a global, on-demand, 24 x 7 workfs at thousands of #ITs completed in minutes CROWPSOURCING VALUE CHAIN CROWP COMMUNITY CROWPSOURCERS (SOLVERS) (SEEKERS) MARKETPLACE (FACILITATOR)

In crowdsourcing, labels are from non-experts (Credit to Amazon Mechanical Turk and IBM Crowdsourcing and Crowdfunding)

#### Passive label collection



In search engine, labels are from users' clicks

(Credit to Google Images)

Motivation	IDN	MCD	Conclusions
00000			

#### Even test sets of most popular benchmarks!

#### Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks

Curtis G. Northcutt*	Anish Athalye
ChipBrain, MIT	MIT

Jonas Mueller Amazon

Table 1: Test set errors are prominent across common benchmark datasets. Errors are estimated using confident learning (CL) and validated by human workers on Mechanical Turk.

Deterat	Madalita	Class	Madal		Test Set Er	rors		
Dataset	Modanty	Size	Model	CL guessed	MTurk checked	validated	estimated	% error
MNIST	image	10,000	2-conv CNN	100	100 (100%)	15	-	0.15
CIFAR-10	image	10,000	VGG	275	275 (100%)	54	-	0.54
CIFAR-100	image	10,000	VGG	2235	2235 (100%)	585	-	5.85
Caltech-256	image	30,607	ResNet-152	4,643	400 (8.6%)	65	754	2.46
ImageNet	image	50,000	ResNet-50	5,440	5,440 (100%)	2,916	-	5.83
QuickDraw	image	50,426,266	VGG	6,825,383	2,500 (0.04%)	1870	5,105,386	10.12
20news	text	7,532	TFIDF + SGD	93	93 (100%)	82	-	1.11
IMDB	text	25,000	FastText	1,310	1,310 (100%)	725	-	2.9
Amazon	text	9,996,437	FastText	533,249	1,000 (0.2%)	732	390,338	3.9
AudioSet	audio	20,371	VGG	307	307 (100%)	275	-	1.35

\*Because the ImageNet test set labels are not publicly available, the ILSVRC 2012 validation set is used.

Motivation	IDN	MCD	Conclusions
0000●	00000000000	00000000	0000
Going beyond CO	CN		

# However, CCN is not enough in expressing/modeling real-world label noise!

We need to go beyond it.

Motivation	IDN	MCD	Conclusions
00000	•0000000000	0000000	0000

# Outline



2 Instance-dependent noise (IDN)

#### Mutually contaminated distributions (MCD)

#### Conclusions

# Mainstream approaches to DL under CCN

#### Loss correction

- Design a corrected loss function such that

minimize corrected loss on noisy data = minimize original loss on clean data

#### • Sample selection/reweighting

- Selection: select data likely with correct labels and train only on those data
- Reweighting: upweight/downweight data likely with correct/incorrect labels

#### Label correction

- Direct: correct the given labels using predicted labels
- Indirect: sample selection + semi-supervised learning

Motivation	IDN	MCD	Conclusions
00000	0000000000	0000000	0000

#### Loss correction may fail under IDN

- CCN assumes  $p(\tilde{y} \mid x, y) = p(\tilde{y} \mid y)$ 
  - It holds that  $oldsymbol{p}_{\widetilde{y}|x}=T^{ op}oldsymbol{p}_{y|x}$  where T is a matrix independent of x
  - Possible to estimate T from data since all instances share the same T
- IDN does not assume  $p(\tilde{y} \mid x, y) = p(\tilde{y} \mid y)$ 
  - It becomes  $oldsymbol{p}_{\widetilde{y}|x} = T(x)^{\!\!\top} oldsymbol{p}_{y|x}$  where T is a matrix-valued function
  - Impossible to estimate T from data since each instance x has its T(x) i.e., IDN is mathematically unidentifiable, regardless of the size of data
  - Hence, without additional assumption/information, loss correction fails

## How about sample selection?

- Sample selection may also fail (Berthon+, ICML 2021; Zhu+, CVPR 2021)
  - The memorization effect is weakened-learn mislabeled data in low-noise regions first
  - Even if sample selection is perfect, a covariate shift exists between clean distributions



Motivation	IDN	MCD	Conclusions
00000	0000000000	0000000	0000

## Wait a minute, can we approximate IDN?

- With additional assumption/information, we can obtain some approximations of IDN (list is not comprehensive)
  - Boundary consistent noise (for binary classification) (Menon+, MLJ 2018)
  - Bounded IDN (for binary classification) (Cheng+, ICML 2020)
  - Part-dependent noise (Xia+, NeurIPS 2020)
  - Difficulty-dependent noise (Wang+, AAAI 2021; Zhu+, CVPR 2021; Zhang+, arXiv 2021)
  - Confidence-scored IDN (Berthon+, ICML 2021)
- After we approximate IDN, we will perform loss correction

Motivation	IDN	MCD	Conclusions
	0000000000		

# Review of backward correction (BC) reconcision

- BC multiplies loss  $\ell$  by  $T^{-1}$  in backward pass
  - Key assumption: underlying transition matrix  ${\mathcal T}$  in CCN is invertible
- Derivation of backward correction

1. 
$$\mathbb{E}_{p(x,y)}[\ell(g(x),y)] = \mathbb{E}_{p(x)}\mathbb{E}_{p(y|x)}[\ell(g(x),y)] = \mathbb{E}_{p(x)}[\mathbf{p}_{y|x}^{\top}\ell_{y|g(x)}]$$
  
where  $\ell_{y|g(x)} = (\ell(g(x),1), \dots, \ell(g(x),c))$ 

2. Then, 
$$\boldsymbol{p}_{\boldsymbol{y}|\boldsymbol{x}}^{\top} \boldsymbol{\ell}_{\boldsymbol{y}|\boldsymbol{g}(\boldsymbol{x})} = (\boldsymbol{p}_{\boldsymbol{y}|\boldsymbol{x}}^{\top} T)(T^{-1} \boldsymbol{\ell}_{\boldsymbol{y}|\boldsymbol{g}(\boldsymbol{x})}) = \boldsymbol{p}_{\boldsymbol{y}|\boldsymbol{x}}^{\top}(T^{-1} \boldsymbol{\ell}_{\boldsymbol{y}|\boldsymbol{g}(\boldsymbol{x})})$$

3. 
$$\mathbb{E}_{p(x)}[p_{y|x}^{\top} \ell_{y|g(x)}] = \mathbb{E}_{p(x)}[p_{\tilde{y}|x}^{\top} \ell_{\tilde{y}|g(x)}^{b}]$$
 where  $\ell_{\tilde{y}|g(x)}^{b} = T^{-1} \ell_{y|g(x)}$ 

- 4. Let  $\ell^b(g(x), \tilde{y}) = [\ell^b_{\tilde{y}|g(x)}]_{\tilde{y}}$  be the (backward-)corrected loss so that  $\mathbb{E}_{p(x,y)}[\ell(g(x), y)] = \mathbb{E}_{p(x,\tilde{y})}[\ell^b(g(x), \tilde{y})]$
- BC reverses label corruption for any classifier & loss!

Motivation	IDN	MCD	Conclusions
	0000000000		

## Review of forward correction (FC) manual concerns

#### • FC multiplies classifier g by T in forward pass

- In BC, g is score &  $\ell$  is proper composite loss = inverse link + base loss
- In FC, g is score + inverse link = an estimated  $\hat{p}(y \mid x) \& \ell$  is base loss

#### Derivation of forward correction

- 1. For any target  $\boldsymbol{q}_{y|x}$  & big enough model h,  $\arg\min_{h} \mathbb{E}_{p(x)}[\boldsymbol{q}_{y|x}^{\top} \ell_{y|h(x)}] = \boldsymbol{q}_{y|x}$
- 2. Let  $\ell^f(g(x), \tilde{y}) = \ell(\mathcal{T}^{\top}g(x), \tilde{y})$  be the (forward-)corrected loss

denote by  $h(x) = T^{\top}g(x)$  as well as  $\ell^{f}_{\tilde{y}|g(x)} = \ell_{\tilde{y}|T^{\top}g(x)} = \ell_{\tilde{y}|h(x)}$ 

- 3. Then,  $\mathbf{T}^{\top} \arg\min_{g} \mathbb{E}_{p(x)}[\mathbf{p}_{\tilde{y}|x}^{\top} \boldsymbol{\ell}_{\tilde{y}|g(x)}^{f}] = \arg\min_{h} \mathbb{E}_{p(x)}[\mathbf{p}_{\tilde{y}|x}^{\top} \boldsymbol{\ell}_{\tilde{y}|h(x)}] = \mathbf{p}_{\tilde{y}|x}$
- 4. This implies  $\arg\min_{g} \mathbb{E}_{p(x)}[\boldsymbol{p}_{\tilde{y}|x}^{\top} \boldsymbol{\ell}_{\tilde{y}|g(x)}^{f}] = \boldsymbol{p}_{y|x} = \arg\min_{g} \mathbb{E}_{p(x)}[\boldsymbol{p}_{y|x}^{\top} \boldsymbol{\ell}_{y|g(x)}]$ i.e.,  $\arg\min_{g} \mathbb{E}_{p(x,\tilde{y})}[\ell^{f}(g(x),\tilde{y})] = \arg\min_{g} \mathbb{E}_{p(x,y)}[\ell(g(x),y)]$
- FC simulates label corruption for probabilistic classifier & loss!

Motivation	IDN	MCD	Conclusions
00000	0000000●0000	0000000	0000

## Part-dependent noise (PDN) (volume and volume and volum

#### • PDN is naturally motivated

- Humans perceive instances based on the parts, physiologically and psychologically
- More likely to annotate an instance based on its parts but not the whole instance
- A wrong mapping from parts to classes would cause PDN (a special case of IDN)

#### • 3 key assumptions of PDN

- Every instance can be decomposed into r parts (a convex combination of r parts)
- For every class, there are at least r anchor points
- For every x, T(x) is a convex combination of r matrices (with the same weights)







Motivation	IDN	MCD	Conclusions
	00000000000		

# Effective learning under PDN (communication)

- 1. Learn the parts and the combination weights
  - 1.1. Estimate  $p(\tilde{y} \mid x)$  from noisy data, and extract latent representations of instances 1.2. Learn the parts and the combination weights by non-negative matrix factorization
- Estimate the rows of T(x) for anchor points
  When x is an anchor point for class i, we obtain that ∀j, [T(x)]<sub>i,j</sub> = p(ỹ = j | x)
- Recover T(x) for all training data (including non-anchor points)
  3.1. Estimate P<sup>1</sup>,..., P<sup>r</sup> given the weights and those rows of T(x) for anchor points
  3.2. Recover T(x) for every training instance x based on the weights and P<sup>1</sup>,..., P<sup>r</sup>





CSIDN ≥ boundary consistent noise + difficulty-dependent noise

- Binary boundary consistent noise: noise gets higher if p(y = 1 | x) is closer to 0.5

- Difficulty-dependent noise: x influences the noise magnitude but not its dynamics

 $p(\tilde{y} \mid \tilde{y} \neq y, x, y) = p(\tilde{y} \mid \tilde{y} \neq y, y) \iff [T(x)]_{i,j|j\neq i} = (1 - [T(x)]_{i,i})[E]_{i,j}$ where  $1 - [T(x)]_{i,i} = p(\tilde{y} \neq y \mid y = i, x)$  controls the magnitude of the noise and  $[E]_{i,j} = p(\tilde{y} = j \mid \tilde{y} \neq y, y = i)$  is CCN and controls the dynamics of the noise

- CSIDN assumes that the confidence information  $r_{x_i} = p(y = \tilde{y}_i | x_i, \tilde{y}_i)$  is available which can indicate both of the boundary information and the difficulty information

Motivation	IDN	MCD	Conclusions
	00000000000		

## Instance-level forward correction (ILFC) (masses musical)

- ILFC minimizes  $\ell(T(x_i)^{\top}g(x_i), \tilde{y}_i)$  for each  $(x_i, \tilde{y}_i, r_i)$ 
  - Without loss of generality, assume that  $\ell$  is the cross-entropy loss
  - Hence, we need the  $\tilde{y}_i$ -th column of  $T(x_i)$  for computing the loss
- How to effectively estimate  $[T(x_i)]_{:,\tilde{y}_i}$ ?
  - 1. The matrix *E* is CCN and thus can be estimated from anchor points and  $\hat{p}(\tilde{y} \mid x)$
  - 2.  $[T(x_i)]_{\tilde{y}_i,\tilde{y}_i}$  can be estimated as  $r_i \hat{\rho}(\tilde{y} = \tilde{y}_i \mid x_i) / \hat{\rho}(y = \tilde{y}_i \mid x_i)$  in an iterative way
  - 3. Note that for  $j \neq \tilde{y}_i$ ,  $r_i = p(y = \tilde{y}_i | x_i, \tilde{y}_i)$  is uninformative to estimate  $[T(x_i)]_{j,j}$ We heuristically set  $[\hat{T}(x_i)]_{i,j}$  as the empirical average of  $[\hat{T}(x_k)]_{j,j}$  where  $\tilde{y}_k = j$
  - 4. Finally,  $[T(x_i)]_{j,\tilde{y}_i|j\neq\tilde{y}_i}$  can be estimated as  $(1 [\hat{T}(x_i)]_{j,j})[\hat{E}]_{j,\tilde{y}_i}$



000000000000000000000000000000000000000	Motivation	IDN	MCD	Conclusions
	00000	0000000000	0000000	0000

# A summary of IDN settings and methods

- IDN strictly generalizes CCN
  - Transition matrix  $T \Longrightarrow$  Matrix-valued function T(x)
- IDN is notably more challenging than CCN
  - The memorization effect acts differently in regions with different T(x)
  - T(x) is not identifiable unless we (roughly or nicely) approximate IDN
  - Rely on part-dependent noise if we can decompose our data into parts
  - Rely on confidence-scored IDN if we collected or can assign the scores
  - Otherwise, try boundary consistent noise or difficulty-dependent noise

Motivation	IDN	MCD	Conclusions
00000	0000000000	•000000	0000

# Outline



Instance-dependent noise (IDN)

#### 3 Mutually contaminated distributions (MCD)

#### 4 Conclusions



# When $p(x \mid y)$ instead of $p(y \mid x)$ is corrupted sourcements





Clean N component

Noisy N mixture (0.4P+0.6N)

Is it still a problem of noisy supervision? Yes! Does it belong to CCN or IDN? No...

Motivation	IDN	MCD	Conclusions
00000	00000000000	00●00000	0000

# MCD also (strictly) generalizes CCN measurements

• In common:  $\{(x_1, \tilde{y}_1), \dots, (x_n, \tilde{y}_n)\}$  drawn from  $p(x, \tilde{y})$ 

• CCN corrupts class-posterior probability:  $\boldsymbol{p}_{\tilde{y}|x} = \mathcal{T}^{\top} \boldsymbol{p}_{y|x}$ 

- T is a label transition matrix such that  $[T]_{i,j} = p(\tilde{y} = j \mid y = i)$
- It is a label-noise model for the corruption of the labeling process
- p(x) remains the same so that the memorization effect is reliable
- $p(\tilde{y})$  is determined once  $p(\tilde{y} \mid x)$  or T is fixed
- MCD corrupts class-conditional density:  $\boldsymbol{p}_{x|\tilde{y}} = S \boldsymbol{p}_{x|y}$ 
  - S is a mixture proportion matrix such that  $[S]_{i,j} = p(y = j \mid \tilde{y} = i)$
  - It is a "label-noise" model for the corruption of the sampling process It is often not viewed as label noise, since instances are also "wrong"
  - $p(\tilde{y})$  is totally free after  $p(x \mid \tilde{y})$  or S is fixed
  - Depending on  $p(\tilde{y})$ , p(x) may notably change (with probability one) The only chance of the same p(x) is when MCD is reduced to CCN Thus, just the memorization effect can be practically very unreliable

Motivation	IDN	MCD	Conclusions
00000	0000000000	0000000	0000

## Backward correction for MCD: an overview

- We are going to rewrite the risk  $R(g) = \mathbb{E}_{p(x,y)}[\ell(g(x), y)]$
- Specifically, R(g) could be decomposed into c partial risks
- We create a loss  $\ell^b$ , such that  $\mathbb{E}_{p(x,\tilde{y})}[\ell^b(g(x),\tilde{y})] = R(g)$
- It could be achieved by solving a set of  $c^2$  linear equations
- The solution is simple:  $\ell^b(\cdot,j) = \sum_{k=1}^{c} \frac{[S^{-1}]_{k,j}p(y=k)}{p(\tilde{y}=j)}\ell(\cdot,k)$

Motivation	IDN	MCD	Conclusions				
00000	00000000000	0000●000	0000				

KISK decomposition (on Regener, MLR 2018, Los, ICLR 2018)

• Key idea: Work directly on class-wise risks

- For CCN, we make use of  $R(g) = \mathbb{E}_{p(x,y)}[\ell(g(x),y)] = \mathbb{E}_{p(x)}\mathbb{E}_{p(y|x)}[\ell(g(x),y)]$ 

- For MCD,  $p(x \mid y)$  gets corrupted, based on which we should write down R(g)

• For 
$$j = 1, \ldots, c$$
, denote by

-  $\pi_j = p(y = j)$  the clean class-prior probability of the *j*-th class

- $\tilde{\pi}_j = p(\tilde{y} = j)$  the noisy class-prior probability of the j-th class
- $p_j(x) = p(x \mid y = j)$  the clean class-conditional density of the *j*-th class
- $\tilde{p}_j(x) = p(x \mid \tilde{y} = j)$  the noisy class-conditional density of the *j*-th class
- Then, R(g) can be decomposed into a class-wise manner
  - $R(g) = \mathbb{E}_{p(y)}\mathbb{E}_{p(x|y)}[\ell(g(x), y)] = \sum_{j=1}^{c} \pi_j \mathbb{E}_{p_j(x)}[\ell(g(x), j)] = \sum_{j=1}^{c} \pi_j R_j(g)$ where  $R_j(g) = \mathbb{E}_{p_j(x)}[\ell(g(x), j)]$  is the (partial) risk of the *j*-th class

Motivation	IDN	MCD	Conclusions
00000	00000000000	00000000	0000

#### Risk alignment and rewrite (paragraph data and detailed and)

- Consider a pseudo risk  $\widetilde{R}(g)$  on the noisy distribution  $p(x, \tilde{y})$ -  $\widetilde{R}(g) = \mathbb{E}_{p(\tilde{y})} \mathbb{E}_{p(x|\tilde{y})}[\ell^{b}(g(x), \tilde{y})] = \sum_{j=1}^{c} \tilde{\pi}_{j} \mathbb{E}_{\tilde{p}_{j}(x)}[\ell^{b}(g(x), j)] = \sum_{j=1}^{c} \tilde{\pi}_{j} \widetilde{R}_{j}(g)$ where  $\ell^{b}$  is the corrected loss and  $\widetilde{R}_{j}(g) = \mathbb{E}_{\tilde{p}_{j}(x)}[\ell^{b}(g(x), j)]$
- We would like to align the pseudo risk  $\widetilde{R}(g)$  to the risk R(g)
  - Theoretically,  $\ell^b(\cdot, j)$  as a linear combination of  $\{\ell(\cdot, 1), \dots, \ell(\cdot, c)\}$  suffices Let U be the coefficient matrix for  $\ell^b$  such that  $\ell^b(\cdot, j) = \sum_{k=1}^{c} [U]_{k,j}\ell(\cdot, k)$
  - $\begin{array}{l} \widetilde{R}_{j}(g) = \mathbb{E}_{\sum_{l}[S]_{j,l} \mathcal{P}_{l}(x)} \left[ \sum_{k} [U]_{k,j} \ell(g(x),k) \right] = \sum_{k,l} [U]_{k,j} [S]_{j,l} \mathbb{E}_{\mathcal{P}_{l}(x)} [\ell(g(x),k)] \\ \text{and thus the coefficient of } \mathbb{E}_{\mathcal{P}_{l}(x)} [\ell(g(x),k)] \text{ in } \widetilde{R}(g) \text{ is } \sum_{j} \widetilde{\pi}_{j} [U]_{k,j} [S]_{j,l} \end{array}$
  - By matching the two coefficients of  $\mathbb{E}_{p_l(x)}[\ell(g(x), k)]$  in  $\widetilde{R}(g)$  and R(g)we can see that  $\sum_i \tilde{\pi}_i[U]_{k,i}[S]_{j,l}$  should be  $\pi_k$  if l = k and 0 otherwise
  - Solving the set of linear equations, we can derive  $[U]_{k,j} = [S^{-1}]_{k,j} \pi_k / \tilde{\pi}_j$
- By risk alignment, we successfully rewrite R(g) into R(g)

Motivation	IDN	MCD	Conclusions
		00000000	

## Consistent risk correction (see many and us astrong and

- However, BC for MCD tends to overfit the training data
  - $\ell^b(\cdot, j)$  is a linear combination but not convex combination of  $\{\ell(\cdot, k)\}$
  - We may suffer from that  $[U]_{k,j} = [S^{-1}]_{k,j} \pi_k / \tilde{\pi}_j < 0$  for some j and k



- Aggressive ideas: enforce  $[U]_{k,j} \ge 0$  or  $\ell^b(g(x_i), \tilde{y}_i) \ge 0$
- Least aggressive idea: just enforce  $\widehat{\mathbb{E}}_{p_i(x)}[\ell(g(x), j)] \ge 0$

# Connection to learning from unlabeled data

- Binary classification (based on empirical risk minimization)
  - Classifier training is impossible given a single set of U data  $({\tt Lu+, \ ICLR \ 2019})$
  - This becomes possible given two sets of U data with different class priors by assuming/forcing  $p(y = +1) = \frac{1}{2}$  (du Plessis+, TAAI 2013; Menon+, ICML 2015)
  - p(y) becomes free (Lu+, ICLR 2019), and practical solution (Lu+, AISTATS 2020)
  - Able to train from  $\geq 3$  different-class-prior U datasets (Lu+, ICML 2021)
- Multi-class classification (based on empirical risk minimization)
  - Should be possible if the number of U datasets = the number of classes
  - However, mapping U datasets to right corrupted classes is combinatorial

Motivation	IDN	MCD	Conclusions
00000	0000000000	0000000	<b>●</b> 000

# Outline



Instance-dependent noise (IDN)

#### Mutually contaminated distributions (MCD)

#### 4 Conclusions

## Two ways to go beyond CCN

#### • Instance-dependent noise (IDN)

- $\boldsymbol{p}_{\tilde{y}|x} = T(x)^{\top} \boldsymbol{p}_{y|x}$ , the best model for the labeling-process corruption
- When we confirm/believe p(x) does not change, apply IDN methods
- Very hard to estimate T(x):

Rely on part-dependent noise if we can decompose our data into parts Rely on confidence-scored IDN if we collected or can assign the scores

- Mutually contaminated distributions (MCD)
  - $p_{x|\tilde{y}} = Sp_{x|y}$ , the best model for the sampling-process corruption
  - When we confirm/believe p(x) may change, apply MCD methods
  - Very hard to estimate S: Best to (re)label a small subset of data
  - Don't forget learning rate decay and/or consistent risk correction

Motivation	IDN	MCD	Conclusions
			0000

## Future directions

- IDN and MCD are huge future directions of noisy supervisions
  - How to adjust/modify the sample selection/label correction methods for them
- Within IDN
  - What assumptions, besides part-dependent noise, can make T(x) identifiable
  - What information, besides confidence scores, can also help to estimate T(x)
- Within MCD
  - How to better mitigate the overfitting of its backward corrections
  - How to accurately estimate S, i.e., the mixture proportion matrix
- Even beyond IDN and MCD
  - A partial label for  $x_i$  is a set  $Y_i$  of candidate labels, including the true label  $y_i$
  - It belongs to inexact supervision rather than inaccurate/noisy supervision but the key ideas here can be applied (Lv+, ICML 2020; Feng+, ICML 2020 & NeurIPS 2020)

Motivation	IDN	MCD	Conclusions
00000	00000000000	00000000	

# Thanks

## Q & A