

Recent Advances in Classification from Noisy Labels

Masashi Sugiyama

RIKEN Center for Advanced Intelligence Project/
The University of Tokyo

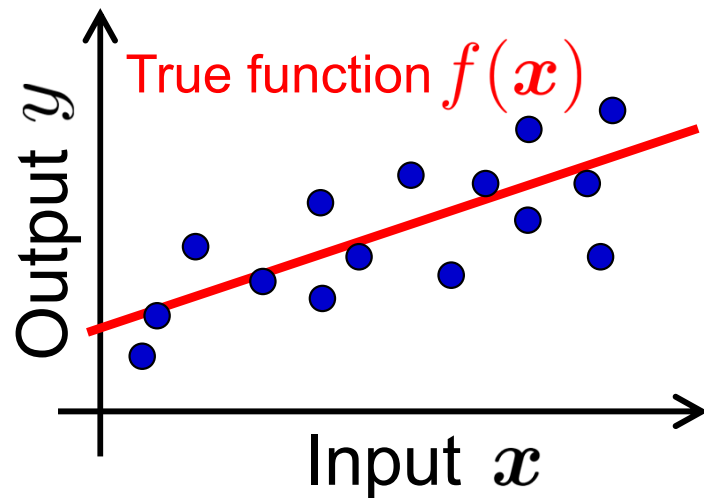


<http://www.ms.k.u-tokyo.ac.jp/sugi/>

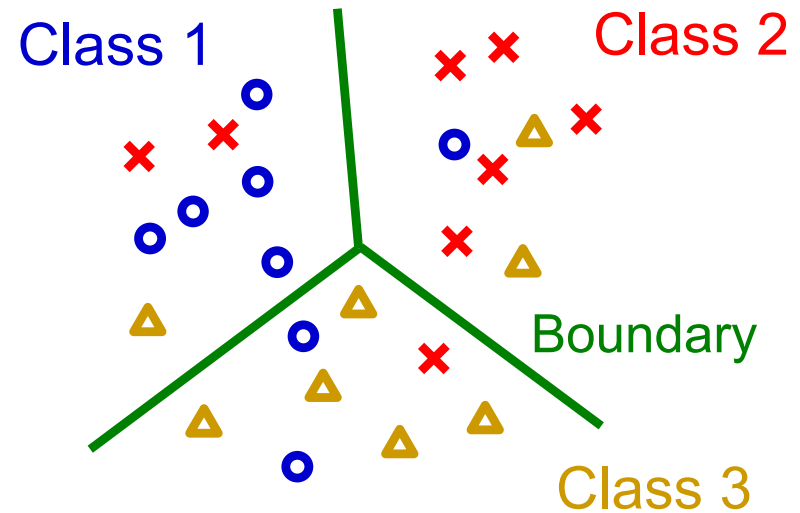


Supervised Learning from Noisy Output Data

Regression



Classification



- Classical problem. Nothing to study further?
 - For regression, just using big data is fine.
 - For classification, big data doesn't necessarily help.
Need further study to cope with label noise!



Contents

1. Formulation and existing results
2. One-step solution
3. Beyond anchor points
4. Further investigations

Formulation

■ Clean training data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$

■ Noisy training data: $\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \bar{p}(\mathbf{x}, \bar{y})$

$\mathbf{x} \in \mathbb{R}^d$: Input instance

$y \in \{1, \dots, c\}$: Clean class label

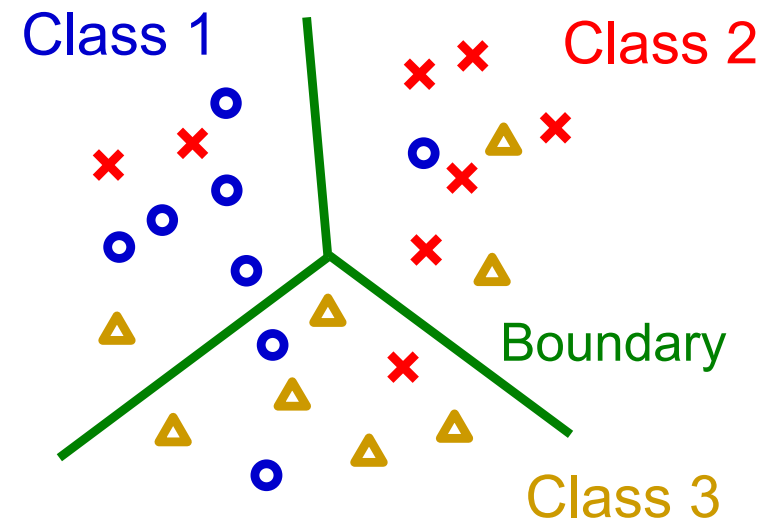
$\bar{y} \in \{1, \dots, c\}$: Noisy class label

■ Probabilistic classifier in simplex: $h(\mathbf{x}) \in \Delta^{c-1}$

- Each element approximates the class-posterior probability.

$$h_y(\mathbf{x}) \approx p(y|\mathbf{x})$$

■ Loss: $\ell(y, h(\mathbf{x})) \in \mathbb{R}$



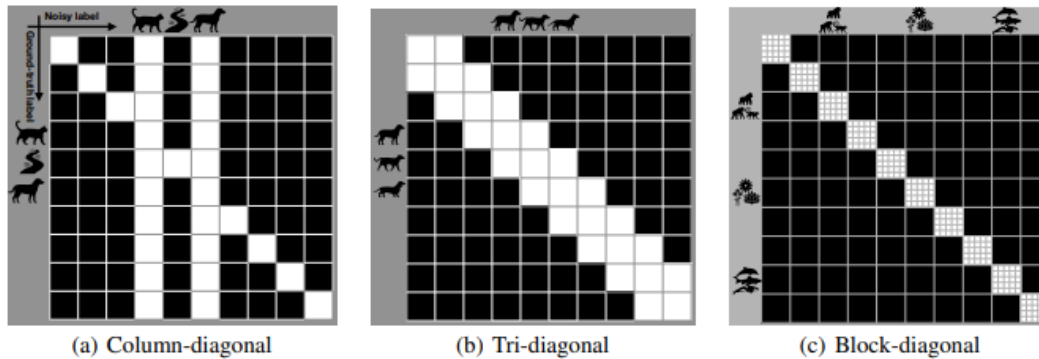
Modeling Class-Conditional Noise 5

■ **Noise transition matrix:** $T_{y, \bar{y}} = \bar{p}(\bar{y}|y)$

y	1	0	0
	0.1	0.8	0.1
	0.5	0.5	0

- Probability of flipping y to \bar{y} .

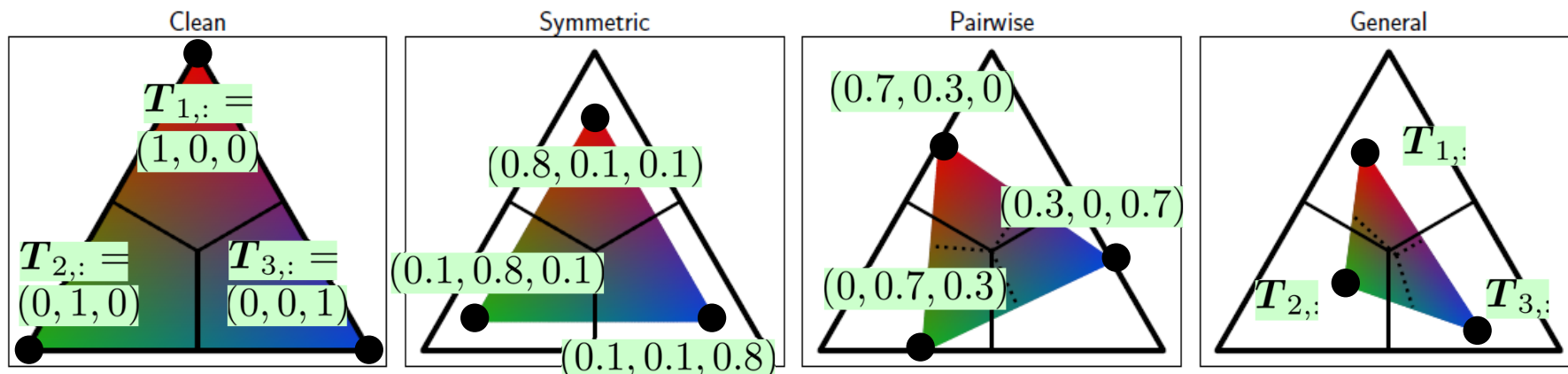
■ We may encode a human-cognitive bias: \bar{y}



Han, Yao, Niu, Zhou, Tsang, Zhang & Sugiyama (NeurIPS2018)

■ **Visualization as a simplex:**

Zhang, Niu & Sugiyama (ICML2021)



■ Forward correction: Add noise by \mathbf{T}^\top

- $\ell^{\rightarrow}(\mathbf{h}(\mathbf{x})) = \ell(\mathbf{T}^\top \mathbf{h}(\mathbf{x}))$ $\ell_y(\mathbf{h}(\mathbf{x})) = \ell(y, \mathbf{h}(\mathbf{x}))$

Classifier-consistency

$$\operatorname{argmin}_h \mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} [\ell^{\rightarrow}(y, \mathbf{h}(\mathbf{x}))] = \operatorname{argmin}_h \mathbb{E}_{p(\mathbf{x}, y)} [\ell(y, \mathbf{h}(\mathbf{x}))]$$

■ Backward correction: Remove noise by \mathbf{T}^{-1}

- $\ell^{\leftarrow}(\mathbf{h}(\mathbf{x})) = \mathbf{T}^{-1} \ell(\mathbf{h}(\mathbf{x}))$

Classifier-consistency

$$\operatorname{argmin}_h \mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} [\ell^{\leftarrow}(y, \mathbf{h}(\mathbf{x}))] = \operatorname{argmin}_h \mathbb{E}_{p(\mathbf{x}, y)} [\ell(y, \mathbf{h}(\mathbf{x}))]$$

Risk-consistency

$$\forall \mathbf{x}, \mathbb{E}_{\bar{p}(\bar{y}|\mathbf{x})} [\ell^{\leftarrow}(y, \mathbf{h}(\mathbf{x}))] = \mathbb{E}_{p(y|\mathbf{x})} [\ell(y, \mathbf{h}(\mathbf{x}))]$$

■ If \mathbf{T} is given, consistency can be guaranteed!

Identifiability of Noise Transition

7

- In practice, we need to estimate T from noisy training data $\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n$.
- However, T is **non-identifiable** in general:
 - T can be decomposed as $T = UV$, where U, V are some transition matrices.
 - Then $\bar{\mathbf{p}}_x = T^\top \mathbf{p}_x = V^\top (U^\top \mathbf{p}_x)$

$T_{y, \bar{y}} = \bar{p}(\bar{y} y)$
$[\bar{\mathbf{p}}_x]_{\bar{y}} = \bar{p}(\bar{y} \mathbf{x})$
$[\mathbf{p}_x]_y = p(y \mathbf{x})$
- Let's use **anchor points** (100%-certain samples) for each class: $\{\mathbf{x}^y \mid p(y|\mathbf{x}^y) = 1\}_{y=1}^c$

Estimation of Noise Transition with Anchor Points

- Given anchor points $\{\mathbf{x}^y \mid p(y|\mathbf{x}^y) = 1\}_{y=1}^c$, $T_{y,\bar{y}} = \bar{p}(\bar{y}|y)$ can be naively estimated as

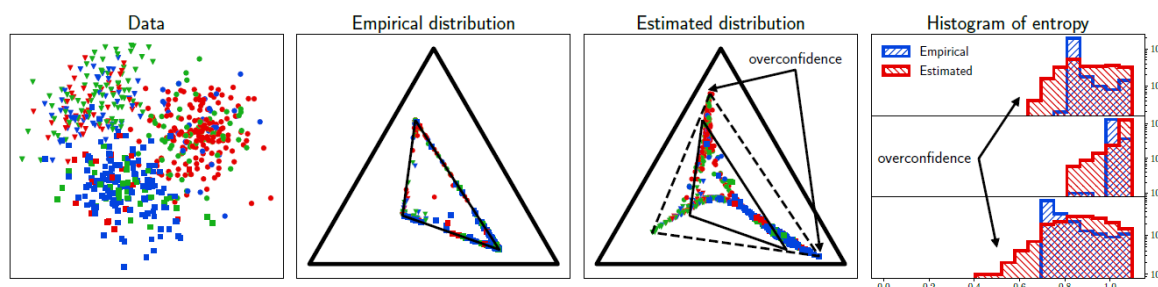
$$T_{y,\bar{y}} = \sum_{y'=1}^c p(\bar{y}|y')p(y'|\mathbf{x}^y) = \bar{p}(\bar{y}|\mathbf{x}^y) \approx \bar{h}_{\bar{y}}(\mathbf{x}^y)$$

- $\bar{h}(\mathbf{x})$ is a probabilistic classifier learned from noisy training data $\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n$.
- Even if anchor points are unknown, as long as they exist in noisy training data, we may find them as $\mathbf{x}^y \leftarrow \mathbf{x}_i$ s.t. $\bar{h}_y(\mathbf{x}_i) \approx 1$.

Further Improvements

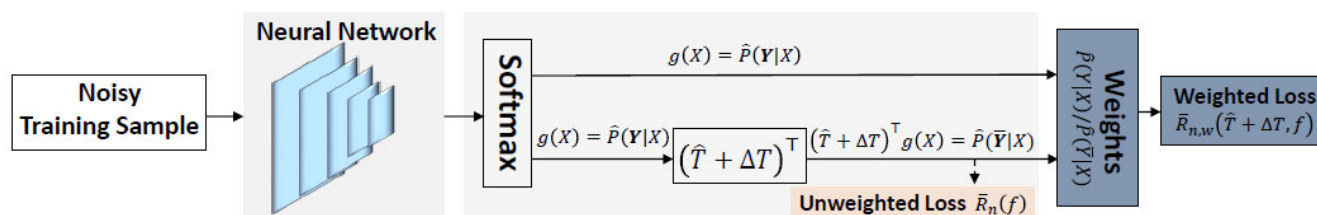
$$\mathbf{x}^y \leftarrow \mathbf{x}_i \text{ s.t. } \bar{h}_y(\mathbf{x}_i) \approx 1$$

- We typically use deep learning to obtain $\bar{h}(x)$:
 - Then it is often **over-confident** and unreliable.



Zhang, Niu & Sugiyama
(ICML2021)

- Estimated T is **revised** during classifier training:



Xia, Liu, Wang,
Han, Gong, Niu
& Sugiyama
(NeurIPS2019)

- Instead of explicitly finding anchor points,
latent labels are utilized: $y'_i = \operatorname{argmax}_{y'} \bar{h}_{y'}(\mathbf{x}_i)$

Yao, Liu, Han, Gong, Deng, Niu, Sugiyama & Tao (NeurIPS2020)



Contents

1. Formulation and existing results
2. **One-step solution**
3. Beyond anchor points
4. Further investigations

- Current approaches are in **two-step**:
 1. Estimate transition matrix T .
 2. Use estimated T to train a classifier $h(x)$.

- Step 1 is done without regard to Step 2:
 - Estimation error of T in Step 1 can be magnified in Step 2.

- We want to estimate T and $h(x)$ **simultaneously in one-step**.

Naïve Solution

- Naively, we may learn the noise transition and classifier at the same time as

$$\min_{\mathbf{U}, \mathbf{h}} \mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} [\ell(\bar{y}, \mathbf{U}^\top \mathbf{h}(\mathbf{x}))] \rightarrow (\mathbf{T}, \mathbf{p}_x)$$

$$T_{y, \bar{y}} = \bar{p}(\bar{y} | y)$$

$$[\mathbf{p}_x]_y = p(y | \mathbf{x})$$

- However, **the solution is not unique**:

- With any invertible transition matrix \mathbf{Q} , any $(\hat{\mathbf{U}}, \hat{\mathbf{h}}) = (\mathbf{Q}^{-1} \mathbf{T}, \mathbf{Q}^\top \mathbf{p}_x)$ are solutions.

- We need a certain **constraint** to obtain the right solution: $(\hat{\mathbf{U}}, \hat{\mathbf{h}}) = (\mathbf{T}, \mathbf{p}_x)$

Total Variation Regularization

13

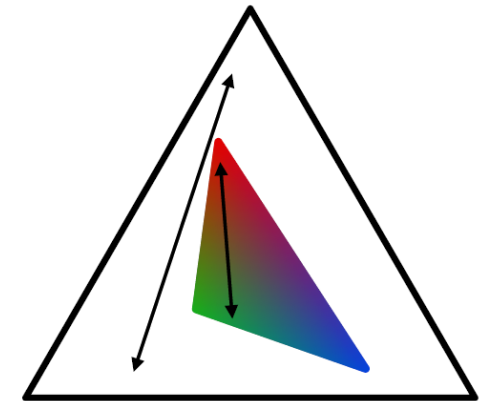
Zhang, Niu & Sugiyama (ICML2021)

- Noise transition $p_x \rightarrow U^\top p_x$ is **contraction** in total variation distance:

$$[p_x]_y = p(y|x)$$

$$\|U^\top p_x - U^\top p_{x'}\|_1 \leq \|p_x - p_{x'}\|_1$$

- **Cleaner class-posteriors have a larger total variation distance!**



- Let's use this knowledge as a regularizer:

$$\min_{U, h} \left[\mathbb{E}_{\bar{p}(x, \bar{y})} [\ell(\bar{y}, U^\top h(x))] - \lambda \mathbb{E}_{p(x), p(x')} \|h(x) - h(x')\|_1 \right]$$

$$\lambda > 0$$

- Under the anchor point assumption, the empirical solution has **statistical consistency**.



Contents

1. Formulation and existing results
2. One-step solution
3. **Beyond anchor points**
4. Further investigations

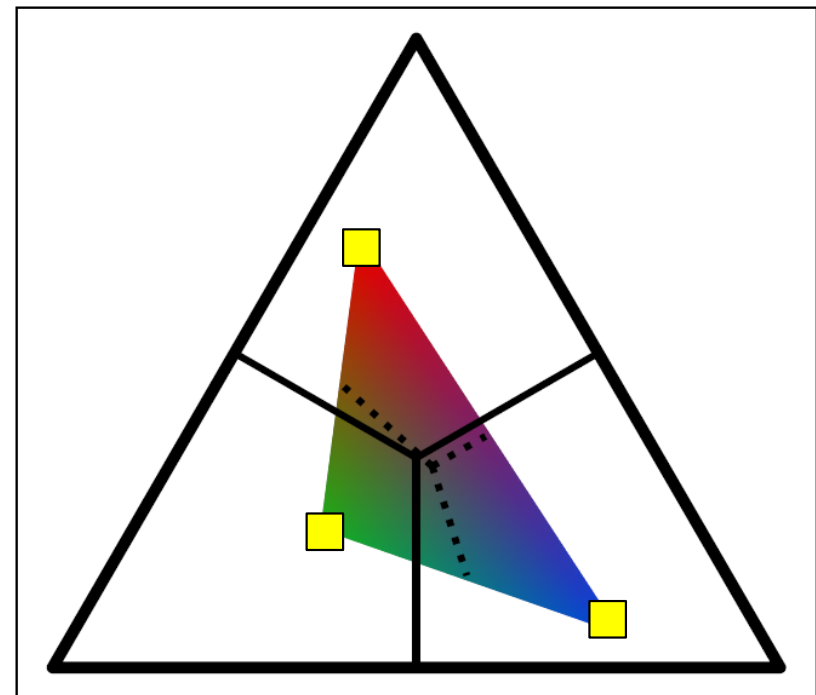
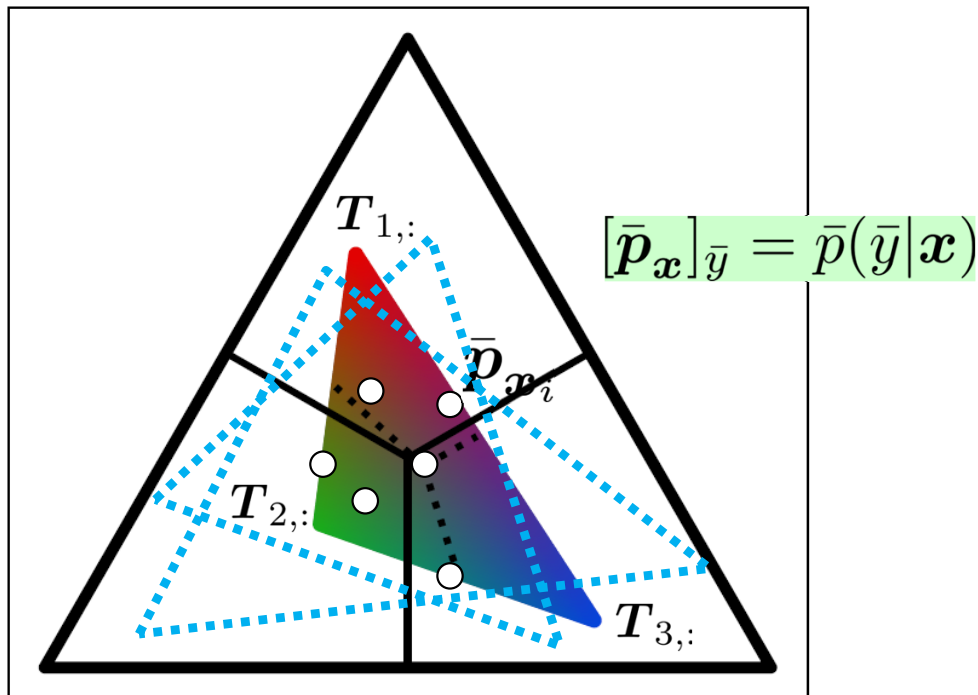
Beyond Anchor Points

$$\{x^y \mid p(y|x^y) = 1\}_{y=1}^c$$

- To overcome the non-identifiability of T :
 - Initial works used **given anchor points** explicitly.
- Later, it was relaxed to only assuming
 - **Existence of anchor points** in training data.
- Further, it was relaxed to assuming
 - Only **existence of anchor regions** (no noise regions) in the true distribution.
- Can we further relax this assumption?
 - Anchor regions rarely exist in reality.

Non-identifiability of T

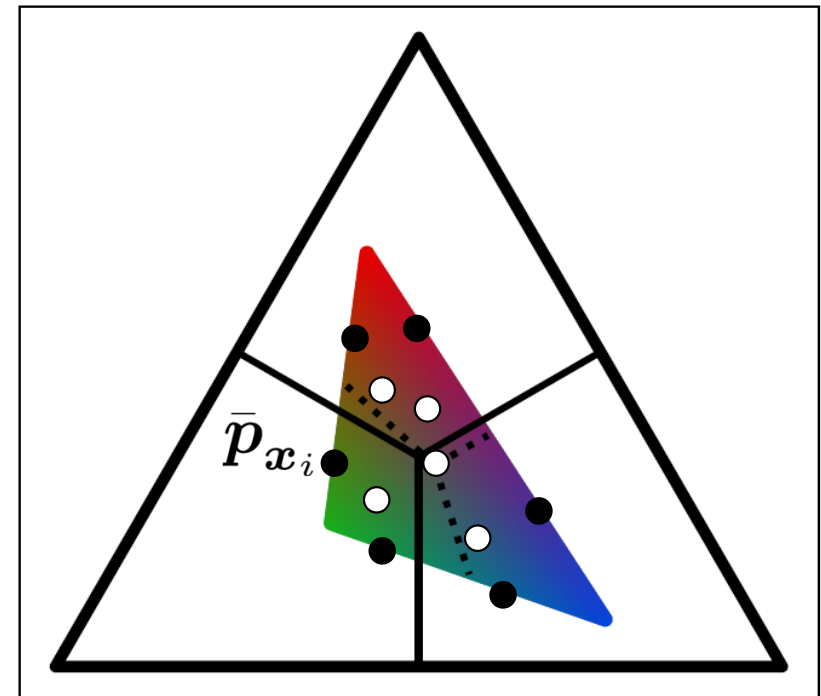
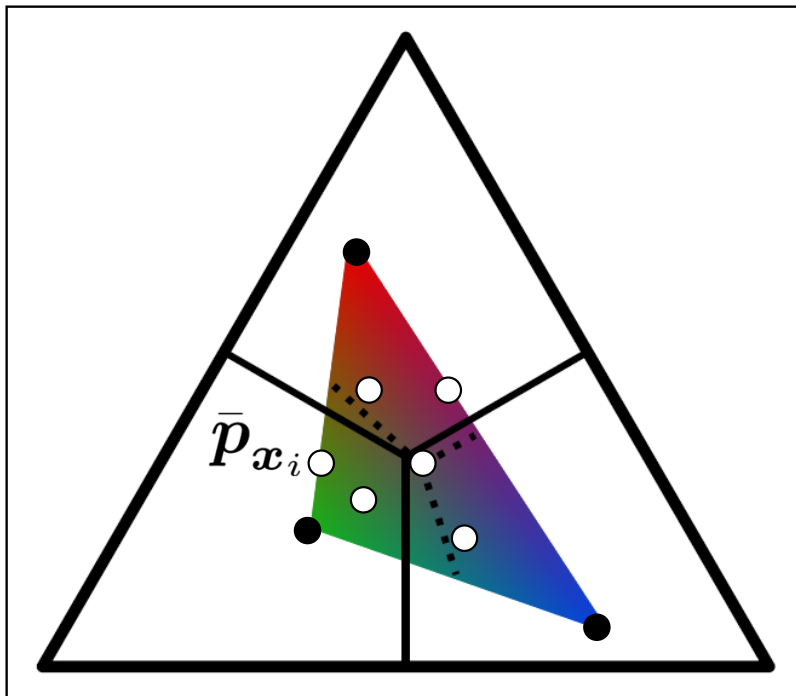
- T can be visualized as a **simplex (triangle)**, containing all training data.
- Generally, such a simplex is not unique.
- **Anchor points are vertices of the true simplex:**
 - Explicitly using anchor points naively recovers T .



Non-identifiability of T (cont.)

17

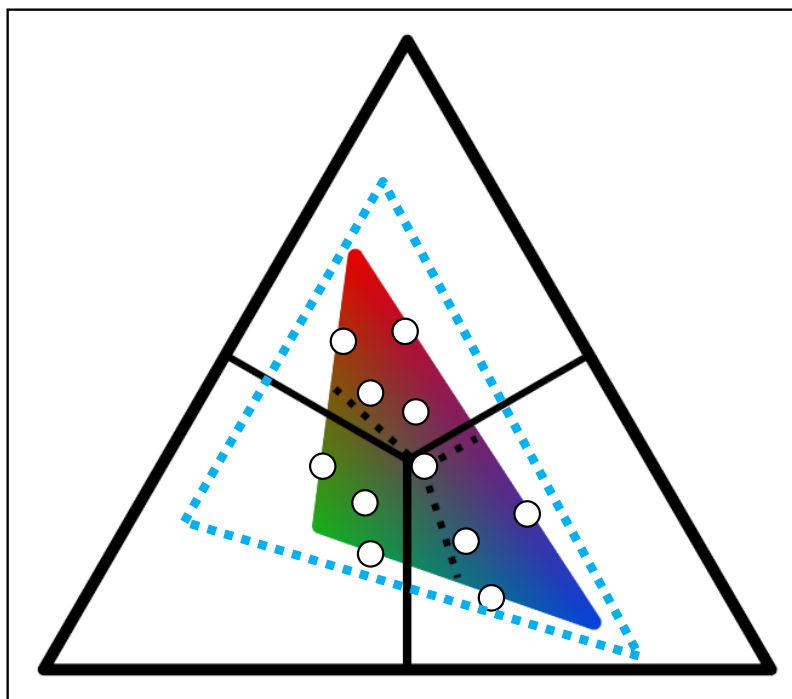
- Only the existence of anchor points still guarantees the identifiability of T .
- Even without anchor points, “sufficiently scattered” training data can guarantee the identifiability.



Li, Liu, Han, Niu & Sugiyama (ICML2021)

- Under the “sufficiently scattered” assumption, **minimizing the volume** of the transition matrix guarantees consistency!

$$\min_{\mathbf{U}, \mathbf{h}} \left[\mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} [\ell(\bar{y}, \mathbf{U}^\top \mathbf{h}(\mathbf{x}))] + \lambda \log \det(\mathbf{U}) \right] \quad \lambda > 0$$

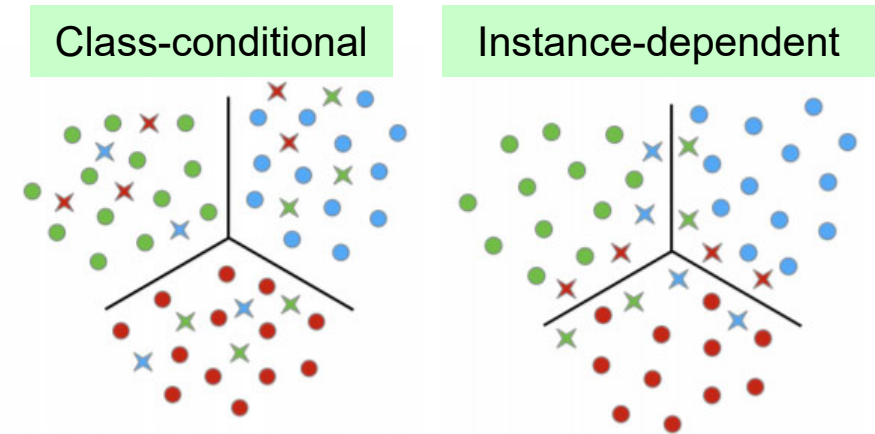




Contents

1. Formulation and existing results
2. One-step solution
3. Beyond anchor points
4. Further investigations

- Instance-independence in class-conditional noise is restrictive.



- Instance-dependent noise: $T_{y, \bar{y}}(\mathbf{x}) = \bar{p}(\bar{y} | y, \mathbf{x})$

- Extremely challenging problem!

- Various new solutions emerge:

- Parts-based estimation
- Use of additional confidence scores
- Manifold regularization

Xia, Liu, Han, Wang,
Gong, Liu, Niu, Tao
& Sugiyama (NeurIPS2020)

Berthon, Han, Niu, Liu
& Sugiyama (ICML2021)

Cheng, Liu, Ning, Wang, Han, Niu,
Gao & Sugiyama (CVPR2022)

■ Memorization of neural nets:

Arpit et al. (ICML2017)
Zhang et al. (ICLR2017)

- Stochastic gradient descent fits clean data faster.
- However, naïve early stopping does not work well.

■ “Co-teaching” between two neural nets:

- Teach small-loss data each other.

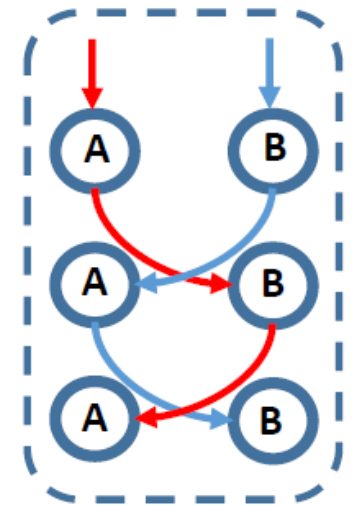
Han, Yao, Yu, Niu, Xu, Hu, Tsang & Sugiyama (NeurIPS2018)

- Teach only disagreed data.

Yu, Han, Yao, Niu, Tsang & Sugiyama (ICML2019)

- Gradient ascent for large-loss data.

Han, Niu, Yu, Yao, Xu, Tsang & Sugiyama (ICML2020)



■ No theory but very robust in experiments:

- Works well even if 50% labels are randomly flipped.

Machine Learning, Neuroscience, and Society 22

- So far, various neuroscientific findings were brought to machine learning with great success:
 - Learning rule, model architecture, adversarial attack, etc.
 - How humans handle noisy observations?
 - Beyond performance improvement, next-generation AI should take into account various **constraints in society**:
 - Culture, common sense, ethics, curiosity, friendliness, etc.
 - Combining
 - neuroscientific findings (**internal learning mechanism**)
 - social demands (**external constraints**)
- would be a promising direction.