

Classification with Asymmetric Label Noise: Consistency and Maximal Denoising

C. Scott¹, **G. Blanchard**², G. Handy¹

¹ U. Michigan

² U. Potsdam



PLAN

- 1 Contamination model
- 2 One contaminated class
- 3 Mutual contamination

STANDARD (GENERATIVE) SETTING FOR CLASSIFICATION

- ▶ $P_i \equiv P(X|Y = i)$: generating probability distributions for objects of class $i = 0, 1$ on space \mathcal{X} .
- ▶ Observed: samples

$$S^i = (X_1^i, \dots, X_{n_i}^i) \stackrel{i.i.d}{\sim} P_i$$

- ▶ **Goal:** estimate decision function $f : \mathcal{X} \rightarrow \{0, 1\}$
- ▶ Various performance error criteria: average classification error, min-max error, Neyman-Pearson error, ...

STANDARD CLASSIFICATION: GENERAL PRINCIPLES

- ▶ Approximate P_i by corresponding empirical distribution \widehat{P}_i
- ▶ For all error criteria, key quantities to estimate for classifiers f are

$$R_i(f) := P_i[f(X) \neq i] \rightarrow \widehat{R}_i(f) := \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{1}\{f(X_j^i) \neq i\}$$

- ▶ agnostic/distribution-free philosophy:
 - ▶ don't want a specific (parametric) model for P_i .
 - ▶ (first) theoretical goal is universal consistency
- ▶ basic strategy: uniform probabilistic control of $\left| R_i(f) - \widehat{R}_i(f) \right|$ over function/set classes \mathcal{C}_k
- ▶ use structural risk minimization to choose adapted class \mathcal{C}_k

CONTAMINATION MODEL

- ▶ Assume the observed samples are drawn according to a **contaminated** distribution:

$$\begin{cases} (X_1^0, \dots, X_{n_0}^0) \stackrel{i.i.d.}{\sim} \tilde{P}_0 & = (1 - \kappa_0)P_0 + \kappa_0 P_1, \\ (X_1^1, \dots, X_{n_1}^1) \stackrel{i.i.d.}{\sim} \tilde{P}_1 & = (1 - \kappa_1)P_1 + \kappa_1 P_0 \end{cases}$$

- ▶ Goal: find a classification function f that performs well for the **true** distributions.
- ▶ Can only access/ estimate

$$\tilde{R}_i(f) := \tilde{P}_i(f(X) \neq i)$$

via

$$\hat{\tilde{R}}_i(f) := \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{1} \{f(X_j^{(i)}) \neq i\}$$

EQUIVALENT MODEL: (ASYMMETRIC) RANDOM LABEL NOISE MODEL

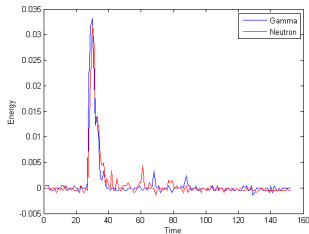
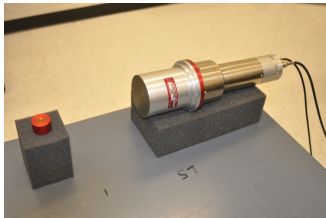
Assume

$$(X_i, Y_i) \stackrel{i.i.d.}{\sim} P;$$

- ▶ true labels Y_i unobserved, instead \tilde{Y}_i
- ▶ corrupted labels $P[\tilde{Y} = i | Y = j, X] = \zeta_{ij}$
- ▶ label corruption assumed not to depend on X
- ▶ label corruption not symmetric

MOTIVATING APPLICATION

ORGANIC SCINCILLATION DETECTOR



- ▶ Detect neutrons and gamma rays; need to classify between them
- ▶ Training using gamma ray source (e.g. Na-22) and neutron source (e.g. Cf-252)
- ▶ But: no pure neutron source – always mixed neutron/gamma ray
- ▶ Additionally, background radiation (both particles)

RELATED WORK, PREVIOUS ASSUMPTIONS

- ▶ Previous work on related topics include:
 - ▶ Learning on positive and unlabeled data (LPUE) (Denis et al. 05, Liu et al. 03)
 - ▶ Co-training (Blum and Mitchell 98)
 - ▶ Label noise models and noise-tolerant PAC learning (Angluin and Laird 88, Kearns 93, Aslam and Deactur 96, Cesa-Bianchi et al. 97, Bshouty et al. 98, Kalai and Servedio 03, Stempfel and Ralaivola 09, Jabbari 10)
- ▶ Generally one or several of the following is assumed:
 - ▶ P_0, P_1 have non-overlapping support (\leftrightarrow deterministic target concept)
 - ▶ symmetric label noise
 - ▶ known noise proportions
 - ▶ criterion is probability of error
- ▶ We do not assume the above here
- ▶ Main assumption: label noise independent of X – no adversarial noise

UNDERSTANDING LABEL NOISE

- ▶ Assume P_0, P_1 have densities ρ_0, ρ_1
- ▶ Then \tilde{P}_0, \tilde{P}_1 have densities

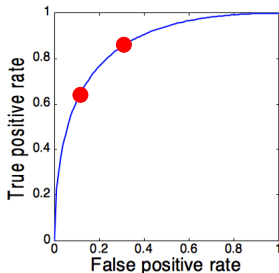
$$\begin{cases} \tilde{\rho}_0 = (1 - \kappa_0)\rho_0 + \kappa_0\rho_1 \\ \tilde{\rho}_1 = (1 - \kappa_1)\rho_1 + \kappa_1\rho_0 \end{cases}$$

Simple algebra:

$$\frac{\rho_1(x)}{\rho_0(x)} \leq \lambda \iff \frac{\tilde{\rho}_1(x)}{\tilde{\rho}_0(x)} \leq \gamma,$$

where

$$\lambda(\gamma) = \frac{\kappa_1 + \gamma(1 - \kappa_1)}{1 - \kappa_0 - \gamma\kappa_0}$$



LABEL NOISE UNDER DIFFERENT ERROR CRITERIA

- ▶ Standard classifier trained on data with noisy labels \tilde{Y} is consistent only when optimal decision is identical under P and \tilde{P} .
- ▶ Consider criterion: **misclassification probability**

$$\mathcal{E}(f) = P[f(X) \neq Y]$$

- ▶ Identical decisions only for **symmetric** label noise $\zeta_{01} = \zeta_{10}$

LABEL NOISE UNDER DIFFERENT ERROR CRITERIA

- ▶ Standard classifier trained on data with noisy labels \tilde{Y} is consistent only when optimal decision is identical under P and \tilde{P} .
- ▶ Consider criterion: **max error**

$$\mathcal{E}(f) = \max(R_0(f), R_1(f))$$

- ▶ Identical decisions if $\kappa_0 = \kappa_1$, or $P_0 = P_1$

LABEL NOISE UNDER DIFFERENT ERROR CRITERIA

- ▶ Standard classifier trained on data with noisy labels \tilde{Y} is consistent only when optimal decision is identical under P and \tilde{P} .
- ▶ Consider criterion: **balanced error**

$$\mathcal{E}(f) = R_0(f) + R_1(f),$$

$$\text{Then: } (1 - \mathcal{E}(f)) = (1 - \kappa_0 - \kappa_1)(1 - \tilde{\mathcal{E}}(f))$$

- ▶ In this case, identical decisions: OK to train on contaminated data (implicit in Blum and Mitchell, 98)

LABEL NOISE UNDER DIFFERENT ERROR CRITERIA

- ▶ Standard classifier trained on data with noisy labels \tilde{Y} is consistent only when optimal decision is identical under P and \tilde{P} .
- ▶ Consider criterion: **balanced error**

$$\mathcal{E}(f) = R_0(f) + R_1(f),$$

$$\text{Then: } (1 - \mathcal{E}(f)) = (1 - \kappa_0 - \kappa_1)(1 - \tilde{\mathcal{E}}(f))$$

- ▶ In this case, identical decisions: OK to train on contaminated data (implicit in Blum and Mitchell, 98)

Overall: training a regular classifier on contaminated data leads to **asymptotic bias and inconsistency** except in very particular circumstances.

- ▶ We can surely estimate $\tilde{R}_i(f)$ from its empirical counterpart

$$\hat{\tilde{R}}_i(f) = \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{1}\{f(X_j^i) \neq i\},$$

uniformly in f in a limited complexity classifier class \mathcal{C}_K

- ▶ Observe

$$\tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0P_1 \implies \tilde{R}_0(f) = (1 - \kappa_0)R_0(f) + \kappa_0R_1(f)$$

$$\tilde{P}_1 = (1 - \kappa_1)P_1 + \kappa_1P_0 \implies \tilde{R}_1(f) = (1 - \kappa_1)R_1(f) + \kappa_1R_0(f)$$

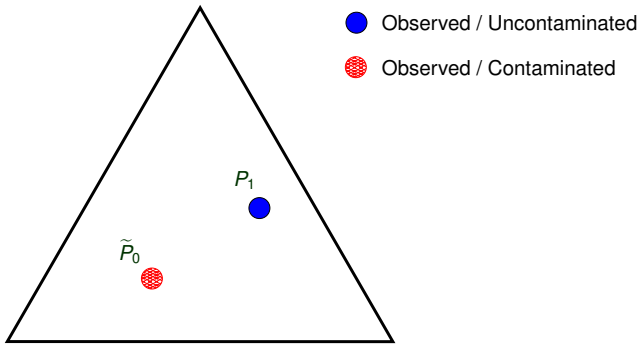
implying

$$R_0(f) = \frac{(1 - \kappa_1)R_0(f) - \kappa_0R_1(f)}{1 - (\kappa_0 + \kappa_1)},$$

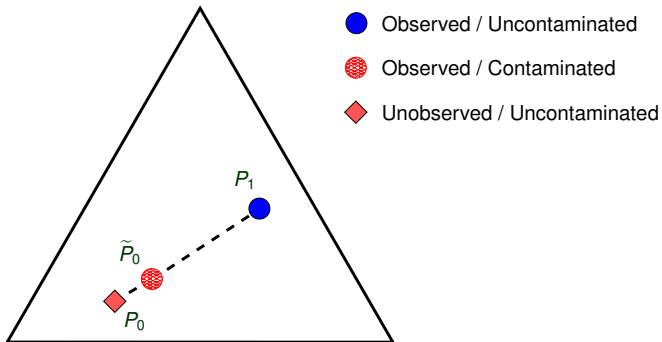
$$R_1(f) = \frac{(1 - \kappa_0)R_1(f) - \kappa_1R_0(f)}{1 - (\kappa_0 + \kappa_1)}$$

- ▶ **Key point:** estimation of contamination proportions κ_0, κ_1 .

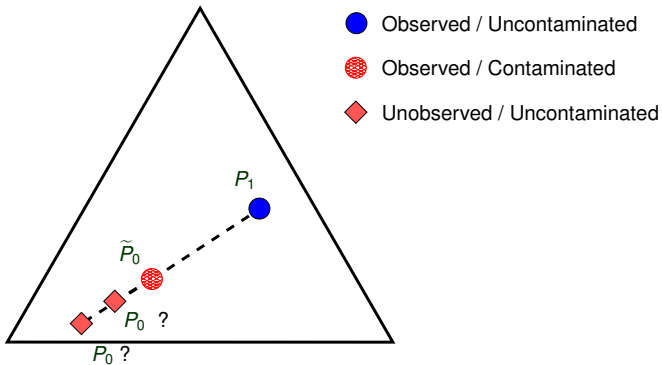
ONLY ONE CONTAMINATED DISTRIBUTION



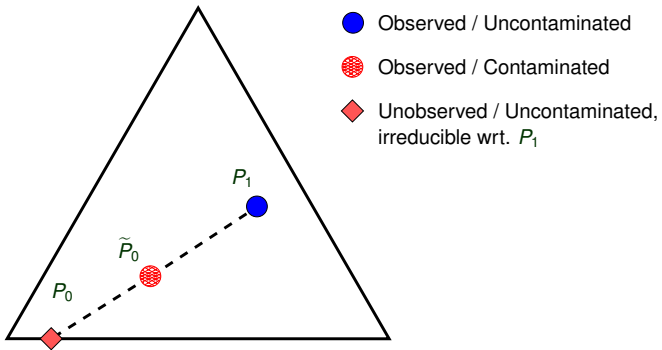
ONLY ONE CONTAMINATED DISTRIBUTION



ONLY ONE CONTAMINATED DISTRIBUTION



ONLY ONE CONTAMINATED DISTRIBUTION



ONLY \tilde{P}_0 CONTAMINATED: IDENTIFIABILITY

[BLANCHARD, SCOTT, LEE 2010]

$$\begin{cases} (X_1^0, \dots, X_{n_0}^0) \stackrel{i.i.d.}{\sim} \tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0 P_1 \\ (X_1^1, \dots, X_{n_1}^1) \stackrel{i.i.d.}{\sim} P_1 \end{cases}$$

- ▶ Define the “maximum proportion of source H in F ”

$$\kappa^*(F|H) = \max \left\{ \kappa \in [0, 1] \mid \exists \text{ a distribution } G \text{ s.t. } F = (1 - \kappa)G + \kappa H \right\} ;$$

- ▶ The following holds:

$$\kappa_0 = \kappa^*(\tilde{P}_0|P_1) \Leftrightarrow \kappa^*(P_0|P_1) = 0 \quad (P_0 \text{ is irreducible wrt. } P_1)$$

ONLY \tilde{P}_0 CONTAMINATED: ESTIMATION

[BLANCHARD, SCOTT, LEE 2010]

- ▶ F, H distributions; Lebesgue decomposition:

$$F = F_H + F_H^\perp,$$

with $F_H \ll H$ and (F_H^\perp, H) mutually singular;

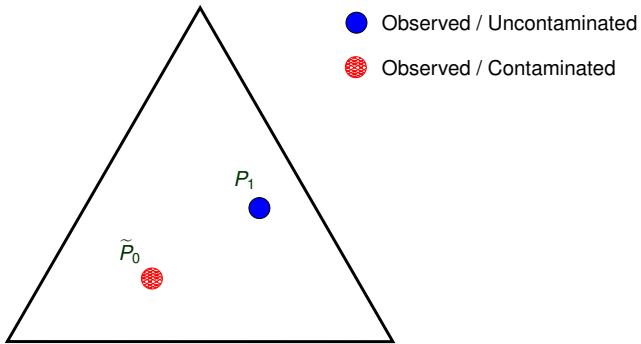
$$\kappa^*(F|H) = \text{Ess. Inf.} \frac{dF_H}{dH} = \inf_{C: H(C) > 0} \frac{F(C)}{H(C)}$$

- ▶ Suggests the estimator

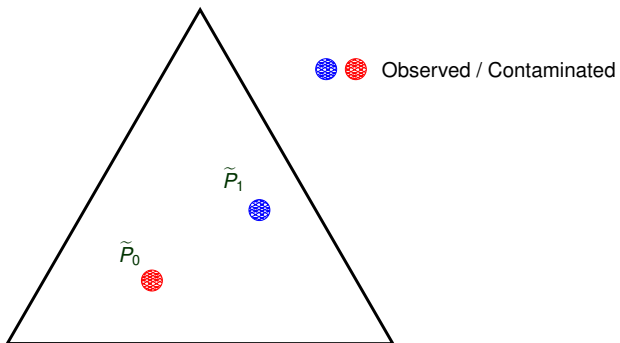
$$\hat{\kappa}(\hat{P}_0|\hat{P}_1) = \inf_{C \in \mathcal{C}_k} \frac{\hat{P}_0(C) + \varepsilon_k}{(\hat{P}_1(C) - \varepsilon_k)_+}$$

- ▶ $\hat{\kappa}(\hat{P}_0|\hat{P}_1) \geq \kappa^*(\tilde{P}_0|P_1)$ with high probability
- ▶ Appropriate choice of ε_k + take inf. over sequence of nested classes $\mathcal{C}_1 \subset \mathcal{C}_2 \subset \dots$ with universal approximation property yields universally consistent estimator

MUTUAL CONTAMINATION



MUTUAL CONTAMINATION



MUTUAL CONTAMINATION

$$\begin{cases} \tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0P_1, \\ \tilde{P}_1 = (1 - \kappa_1)P_1 + \kappa_1P_0 \end{cases}$$

Proposition (Decoupled Representation)

Assume $P_0 \neq P_1$ and

$$(A) \quad \kappa_1 + \kappa_2 < 1;$$

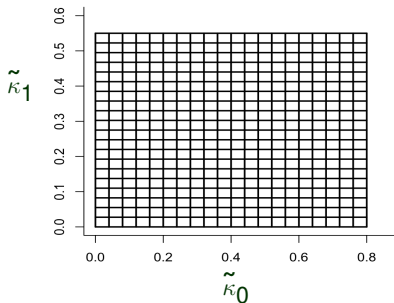
then $\tilde{P}_0 \neq \tilde{P}_1$, and there exist unique $0 \leq \tilde{\kappa}_0, \tilde{\kappa}_1 < 1$ such that

$$\begin{cases} \tilde{P}_0 = (1 - \tilde{\kappa}_0)P_0 + \tilde{\kappa}_0\tilde{P}_1, \\ \tilde{P}_1 = (1 - \tilde{\kappa}_1)P_1 + \tilde{\kappa}_1\tilde{P}_0. \end{cases}$$

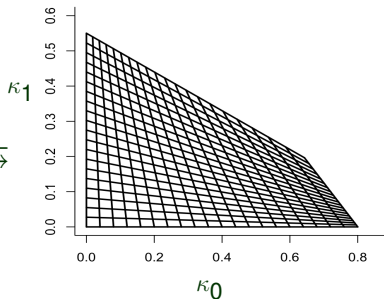
with

$$\tilde{\kappa}_0 = \frac{\kappa_0}{1 - \kappa_1} < 1; \quad \tilde{\kappa}_1 = \frac{\kappa_1}{1 - \kappa_0} < 1.$$

THE TWO REPRESENTATIONS



Decoupled representation



Original representation

$$\begin{cases} \tilde{P}_0 = (1 - \tilde{\kappa}_0)P_0 + \tilde{\kappa}_0\tilde{P}_1, \\ \tilde{P}_1 = (1 - \tilde{\kappa}_1)P_1 + \tilde{\kappa}_1\tilde{P}_0. \end{cases}$$

$$\begin{cases} \tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0P_1, \\ \tilde{P}_1 = (1 - \kappa_1)P_1 + \kappa_1P_0 \end{cases}$$

IDENTIFIABILITY

Decoupled model:

$$\begin{cases} \tilde{P}_0 = (1 - \tilde{\kappa}_0)P_0 + \tilde{\kappa}_0\tilde{P}_1, \\ \tilde{P}_1 = (1 - \tilde{\kappa}_1)P_1 + \tilde{\kappa}_1\tilde{P}_0. \end{cases}$$

From the results on mixture proportion estimation: we can estimate $\tilde{\kappa}_0$ consistently if $\kappa(P_0, \tilde{P}_1) = 0$

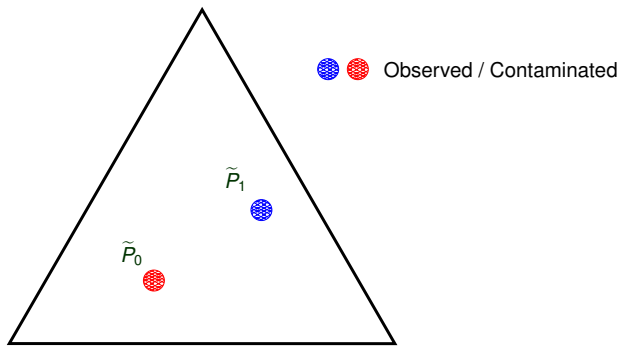
Lemma

Under assumption **(A)**: $\kappa_0 + \kappa_1 < 1$, it holds

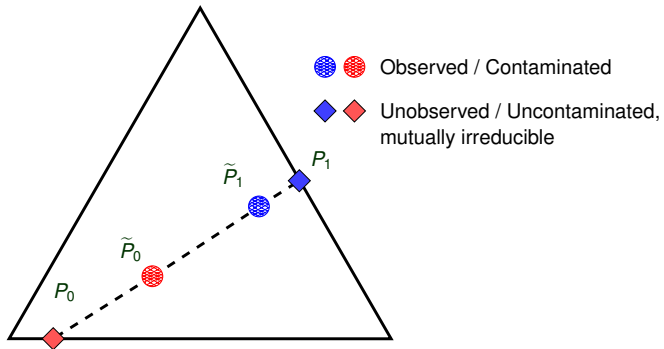
$$\mathbf{(B)} \left\{ \begin{array}{l} \kappa(P_0|\tilde{P}_1) = 0 \\ \kappa(P_1|\tilde{P}_0) = 0 \end{array} \right\} \iff \left\{ \begin{array}{l} \kappa(P_0|P_1) = 0 \\ \kappa(P_1|P_0) = 0 \end{array} \right\} \mathbf{(C)}$$

(C): P_0 and P_1 are mutually irreducible

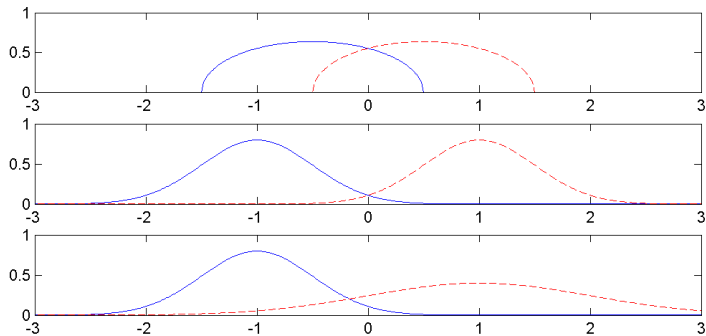
IDENTIFIABILITY



IDENTIFIABILITY



MUTUAL IRREDUCIBILITY



- ▶ Top: mutually irreducible
- ▶ Middle: mutually irreducible
- ▶ Bottom: P_1 irreducible wrt P_0 , but P_0 not irreducible wrt P_0 .

MUTUAL IRREDUCIBILITY

Under joint distribution model

$$(X, Y) \sim \mathbb{P}_{XY}, \quad \eta(x) = P_{XY}[Y = 1|X = x]$$

Then:

$$\left. \begin{array}{l} \kappa(P_0|P_1) = 0 \\ \kappa(P_1|P_0) = 0 \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{l} \text{Ess.Sup.}_x \eta(x) = 1, \\ \text{Ess.Inf.}_x \eta(x) = 0, \end{array} \right.$$

CHARACTERIZING THE IRREDUCIBLE SOLUTION

For given observed contaminated $\tilde{P}_0 \neq \tilde{P}_1$, let Δ be the convex set of quadruples $(\kappa_0, \kappa_1, P_0, P_1)$ satisfying **(A)** and solution of:

$$\begin{cases} \tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0 P_1, \\ \tilde{P}_1 = (1 - \kappa_1)P_1 + \kappa_1 P_0 \end{cases} \quad (1)$$

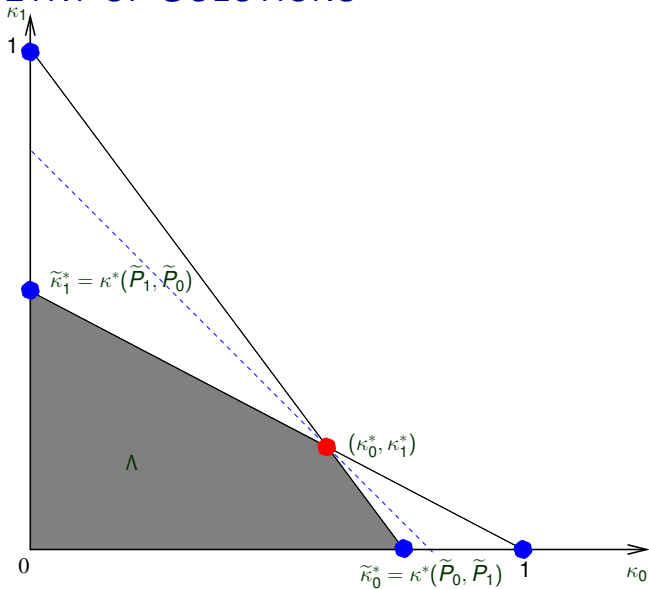
Proposition

The solution $(\kappa_0^*, \kappa_1^*, P_0^*, P_1^*)$ is characterized as either of:

- ▶ the unique quadruple for which (P_0, P_1) are mutually irreducible;
- ▶ the unique nontrivial $(\kappa_0 \neq 0, \kappa_1 \neq 0)$ extremal point of Λ ;
- ▶ the unique maximizer of $(\kappa_0 + \kappa_1)$ over Λ ;
- ▶ the unique maximizer over Λ of $\|P_0 - P_1\|_{TV}$.
- ▶ the unique minimizer over Λ of optimal balanced error for classifying P_0 vs. P_1 .

Interpretation: maximal denoising / source separation

GEOMETRY OF SOLUTIONS



CONSISTENT ESTIMATION OF CONTAMINATION PROPORTIONS

Decoupled representation:

$$\begin{cases} \tilde{P}_0 = (1 - \tilde{\kappa}_0)P_0 + \tilde{\kappa}_0\tilde{P}_1, \\ \tilde{P}_1 = (1 - \tilde{\kappa}_1)P_1 + \tilde{\kappa}_1\tilde{P}_0. \end{cases}$$

- ▶ (P_0, P_1) mutually irreducible $\Rightarrow P_0$ irreducible wrt \tilde{P}_1 , and P_1 irreducible wrt. \tilde{P}_0
- ▶ leverage case of only one contaminated distribution (twice):

$$\hat{\kappa}_0 = \hat{\kappa}(\hat{\tilde{P}}_0 | \hat{\tilde{P}}_1); \quad \hat{\kappa}_1 = \hat{\kappa}(\hat{\tilde{P}}_1 | \hat{\tilde{P}}_0)$$

- ▶ Then

$$\hat{\kappa}_0 = \frac{\hat{\kappa}_0(1 - \tilde{\kappa}_1)}{1 - \tilde{\kappa}_0\tilde{\kappa}_1}; \quad \hat{\kappa}_1 = \frac{\hat{\kappa}_1(1 - \tilde{\kappa}_0)}{1 - \tilde{\kappa}_0\tilde{\kappa}_1}$$

are universally consistent estimators of κ_0, κ_1 under **(A)**, **(C)**.

CONSISTENT ESTIMATION OF RISK

- ▶ Construction of estimator for type II error:

$$\begin{aligned}\tilde{P}_0 &= (1 - \tilde{\kappa}_0)P_0 + \tilde{\kappa}_0\tilde{P}_1 \Rightarrow R_0(f) = \frac{\tilde{R}_0(f) - \tilde{\kappa}_0(1 - \tilde{R}_1(f))}{1 - \tilde{\kappa}_0} \\ &\rightarrow \hat{R}_0(f) = \frac{\hat{\tilde{R}}_0(f) - \hat{\tilde{\kappa}}_0(1 - \hat{\tilde{R}}_1(f))}{1 - \hat{\tilde{\kappa}}_0}\end{aligned}$$

- ▶ Convergence uniform over e.g. VC-Classes of classifiers f
- ▶ Can apply SRM principle to choose appropriate model
- ▶ Can construct universally consistent estimators for various error measures

CONCLUSION

Contributions:

- ▶ nonparametric/distribution-free point of view of the (asymmetric) contamination problem
- ▶ existence and unicity of irreducible solution ; characteristic properties
- ▶ consistent estimation of contamination weights under irreducibility
- ▶ consistent estimation of optimal decisions under different error criteria

Further work:

- ▶ rates
- ▶ multiclass case → more challenging

THANK YOU