

Classification in mutual contamination models

G. Blanchard

Universität Potsdam, Institut für Mathematik

Label Noise Workshop, Nancy 30/11/2017

Joint work with C. Scott, G. Handy, J. Katz-Samuels (U. Michigan)



PLAN

- 1 Contamination models
- 2 Binary classification case
 - One contaminated class
 - Mutual contamination
- 3 Multiclass case
 - One contaminated class
 - Mutual contamination

OUTLINE

- 1 Contamination models
- 2 Binary classification case
 - One contaminated class
 - Mutual contamination
- 3 Multiclass case
 - One contaminated class
 - Mutual contamination

STANDARD (GENERATIVE) SETTING FOR CLASSIFICATION

- ▶ $P_i \equiv P(X|Y = i)$: generating probability distributions for objects of class $1 \leq i \leq L$ 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 $f : \mathcal{X} \rightarrow \{1, \dots, L\}$
- ▶ Various performance error criteria: average classification error, min-max error, Neyman-Pearson error, ...

STANDARD CLASSIFICATION: GENERAL PRINCIPLES

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

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

- ▶ “Learning”/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 $\left| R_i(f) - \hat{R}_i(f) \right|$ over function/set classes \mathcal{C}_k
- ▶ Use structural risk minimization to choose adapted class \mathcal{C}_k

CONTAMINATION MODEL

- ▶ Assume the sample S_i is drawn according to a **contaminated** distribution:

$$S_i = (X_1^i, \dots, X_{n_i}^i) \stackrel{i.i.d.}{\sim} \tilde{P}_i = \sum_{j=1}^L \pi_{ij} P_j$$

or in short form

$$\hat{P} = \Pi P$$

(Π : mixing matrix)

- ▶ **Goal:** find a classification function f that performs well for the **true** source distributions.
- ▶ **Goal:** estimate mixing weights Π and **true** sources (demixing)
- ▶ Can only access/ estimate

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

EQUIVALENT MODEL: (ASYMMETRIC) RANDOM LABEL NOISE MODEL

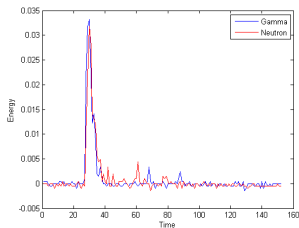
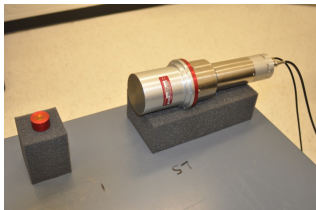
Assume

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

- ▶ True labels Y_i unobserved, instead \tilde{Y}_i
- ▶ Corrupted labels $\mathbb{P}[\tilde{Y} = i | Y = j, X] = \zeta_{ij}$
- ▶ 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)

FURTHER SETTINGS AND GOALS

- ▶ Recover source distributions *up to permutation*: **demixing** problem.
 - ▶ **Application**: Topic models (each observed document is a mixture of topics; goal is to recover “pure” topic distribution themselves)

- ▶ Recover source distributions with the additional knowledge of the **support of Π** (positions of positive entries).
 - ▶ **Application**: Partial labels models (each object comes with a subset of labels)

UNDERSTANDING LABEL NOISE

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

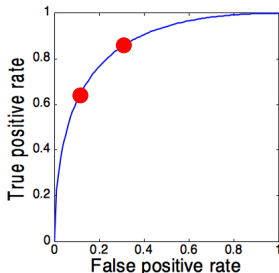
$$\begin{cases} \tilde{\rho}_0 = (1 - \kappa_0)\rho_0 + \kappa_0\rho_1 \\ \tilde{\rho}_1 = (1 - \kappa_1)\rho_1 + \kappa_1\rho_0 \end{cases}$$

Simple algebra:

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

where

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



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

RELATED WORK, PREVIOUS ASSUMPTIONS

- ▶ Previous work on related topics include:
 - ▶ Learning on positive and unlabeled data (LPUE)
 - ▶ Co-training
 - ▶ Label noise models and PAC learning
- ▶ Generally the following is assumed:
 - ▶ P_0, P_1 have non-overlapping support (\leftrightarrow deterministic target concept)
 - ▶ symmetric label noise
 - ▶ criterion is probability of error
- ▶ We do not assume the above here
- ▶ Main assumption: label noise independent of X – no adversarial noise

- ▶ 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

$$R_0(f) = \frac{(1 - \kappa_1)\tilde{R}_0(f) - \kappa_0\tilde{R}_1(f)}{1 - (\kappa_0 + \kappa_1)},$$

$$R_1(f) = \frac{(1 - \kappa_0)\tilde{R}_1(f) - \kappa_1\tilde{R}_0(f)}{1 - (\kappa_0 + \kappa_1)}$$

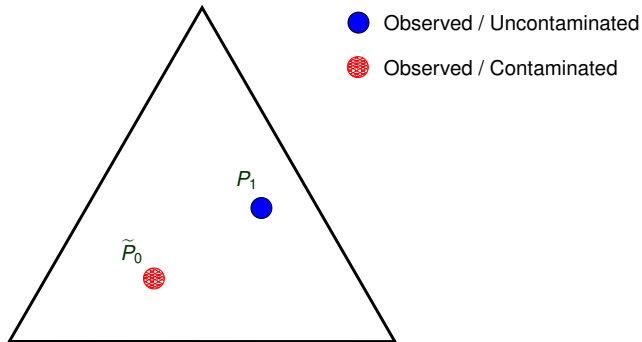
- ▶ **Key point:** estimation of contamination proportions κ_0, κ_1 .

OUTLINE

- 1 Contamination models
- 2 Binary classification case**
 - One contaminated class
 - Mutual contamination
- 3 Multiclass case
 - One contaminated class
 - Mutual contamination

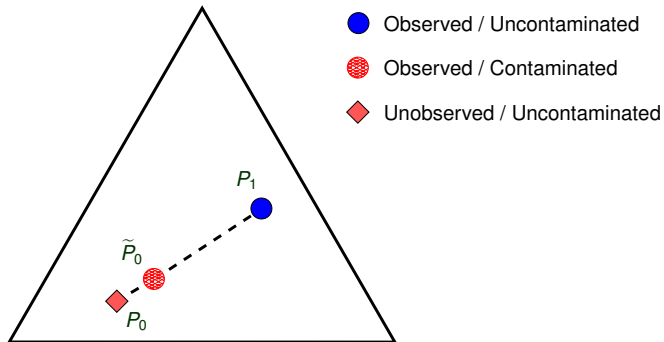
THE BINARY CASE

ONLY ONE CONTAMINATED DISTRIBUTION



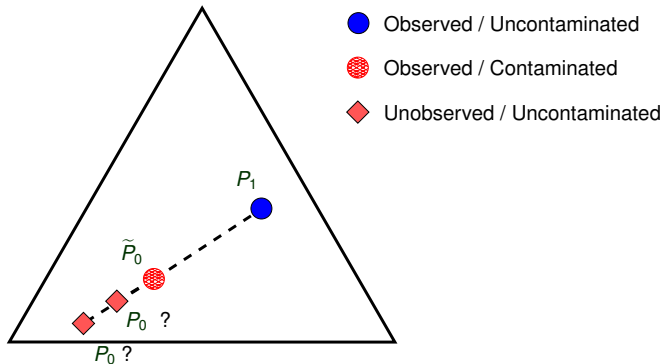
THE BINARY CASE

ONLY ONE CONTAMINATED DISTRIBUTION



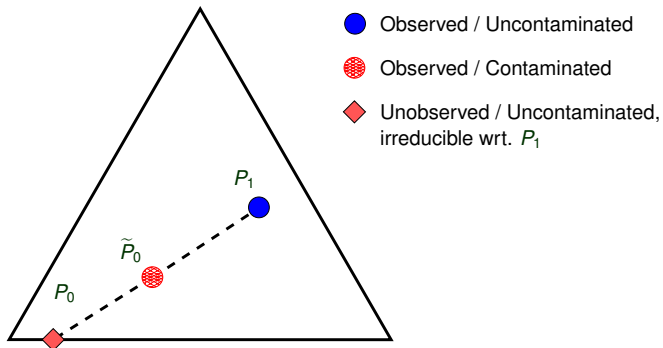
THE BINARY CASE

ONLY ONE CONTAMINATED DISTRIBUTION



THE BINARY CASE

ONLY ONE CONTAMINATED DISTRIBUTION



ONLY \tilde{P}_0 CONTAMINATED: IDENTIFIABILITY

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

- ▶ 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

- ▶ F, H distributions; Lebesgue decomposition:

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

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

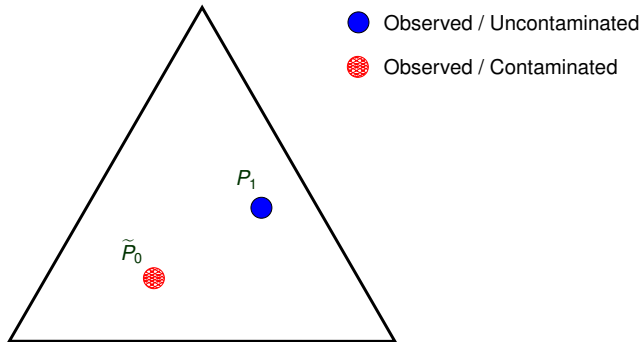
$$\kappa^*(F|H) = \text{Ess. Inf.} \frac{dF_H}{dH} = \inf_{C: H(C) > 0} \frac{F(C)}{H(C)}$$

- ▶ Suggests the estimator

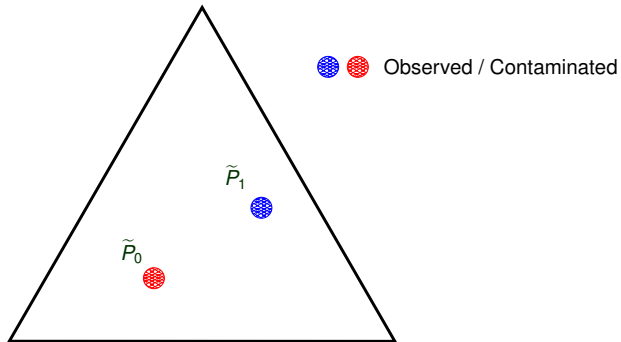
$$\hat{\kappa}(\hat{P}_0|\hat{P}_1) = \inf_{C \in \mathcal{C}_k} \frac{\hat{P}_0(C) + \varepsilon_k}{(\hat{P}_1(C) - \varepsilon_k)_+}$$

- ▶ $\hat{\kappa}(\hat{P}_0|\hat{P}_1) \geq \kappa^*(\tilde{P}_0|P_1)$ 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

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

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

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

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 $\kappa(P_0, \tilde{P}_1) = 0$

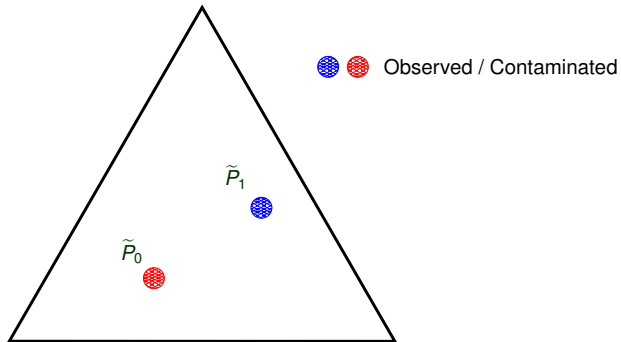
Lemma

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

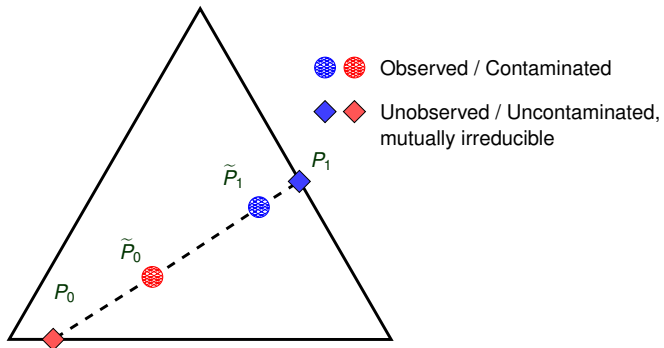
$$\mathbf{(B)} \left\{ \begin{array}{l} \kappa(P_0 | \tilde{P}_1) = 0 \\ \kappa(P_1 | \tilde{P}_0) = 0 \end{array} \right\} \iff \left\{ \begin{array}{l} \kappa(P_0 | P_1) = 0 \\ \kappa(P_1 | P_0) = 0 \end{array} \right\} \mathbf{(C)}$$

(C): P_0 and P_1 are mutually irreducible

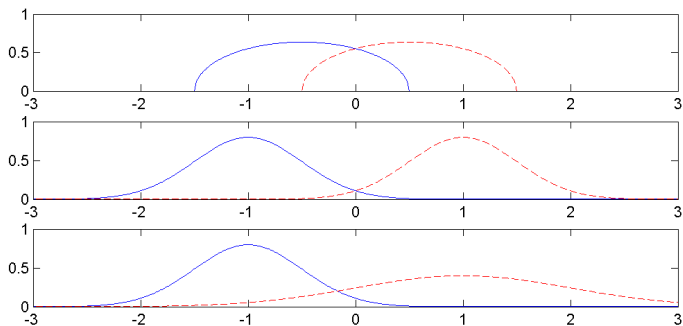
IDENTIFIABILITY



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

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

Then:

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

INTERPRETATION OF 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:

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

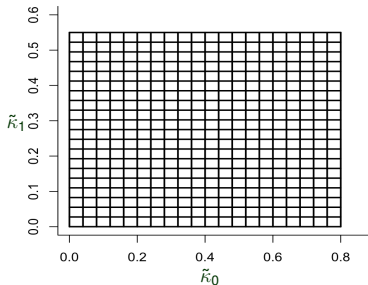
Proposition

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 maximizer over Λ of $\|P_0 - P_1\|_{TV}$.
- ▶ the unique minimizer over Λ of the Bayes error for classifying P_0 vs. P_1 with equal a priori proportions

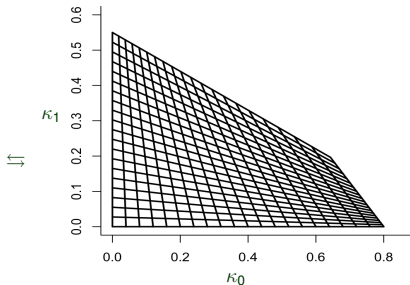
Interpretation: **maximal denoising**

THE TWO REPRESENTATIONS



Decoupled representation

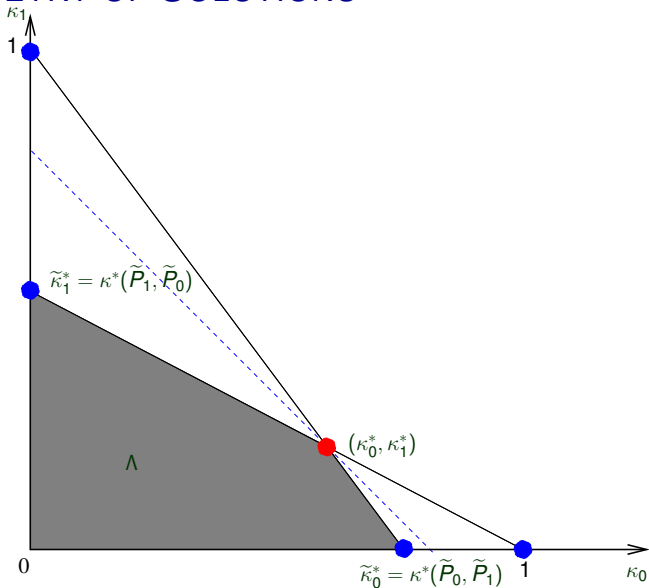
$$\begin{cases} \tilde{P}_0 = (1 - \tilde{\kappa}_0)P_0 + \tilde{\kappa}_0\tilde{P}_1, \\ \tilde{P}_1 = (1 - \tilde{\kappa}_1)P_1 + \tilde{\kappa}_1\tilde{P}_0. \end{cases}$$



Original representation

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

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:

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

- ▶ **Uniform convergence** 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

OUTLINE

- 1 Contamination models
- 2 Binary classification case
 - One contaminated class
 - Mutual contamination
- 3 Multiclass case
 - One contaminated class
 - Mutual contamination

ONLY \tilde{P}_0 CONTAMINATED

$$\left\{ \begin{array}{l} (X_1^0, \dots, X_{n_0}^0) \stackrel{i.i.d.}{\sim} \tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0 \left(\sum_{i=1}^M \mu_i P_i \right) \\ (X_1^i, \dots, X_{n_i}^i) \stackrel{i.i.d.}{\sim} P_i; \quad i = 1, \dots, M \end{array} \right. \quad \text{with } \sum_i \mu_i = 1$$

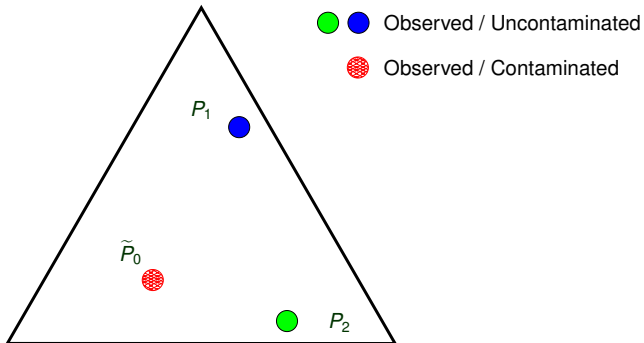
- ▶ Maximum collective proportion of H_1, \dots, H_M in F ?

$$\kappa^*(F|H_1, \dots, H_M) = \max_{\mu \in S_M} \kappa^*(F|H_\mu)$$

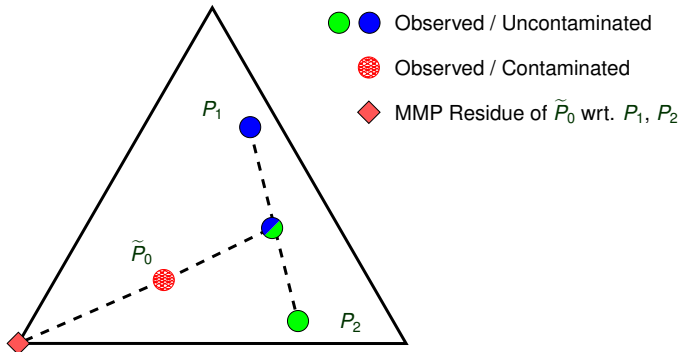
where

- ▶ S_M : $(M - 1)$ -dimensional simplex
- ▶ For $\mu \in S_M$: $H_\mu = \sum_{i=1}^M \mu_i H_i$
- ▶ Interpretation: attained for “projection” of F onto convex hull of $\{H_1, \dots, H_M\}$ for the separation distance $1 - \kappa^*(F|\bullet)$

MAXIMAL MIXTURE PROPORTION



MAXIMAL MIXTURE PROPORTION



MAXIMAL MIXTURE PROPORTION ESTIMATION (MMPE)

$$\kappa^*(F|H_1, \dots, H_M) = \max_{\mu \in S_M} \kappa^*(F|H_\mu)$$

- ▶ Estimator:

$$\hat{\kappa}(\hat{P}_0|\hat{P}_1, \dots, \hat{P}_M) = \max_{\mu \in S_M} \inf_{C \in \mathcal{C}_k} \frac{\hat{P}_0(C) + \varepsilon_{0,k}}{\left(\hat{P}_\mu(C) - \sum_i \mu_i \varepsilon_{i,k}\right)},$$

for (\mathcal{C}_k) sequence of VC-classes

- ▶ $\hat{\kappa} \geq \kappa^*$ with high probability
- ▶ Universally consistent if the VC sequence is universally approximating
- ▶ $\hat{\mu}$ attaining the max converges to the population maximum μ , whenever the latter is unique

IDENTIFIABILITY

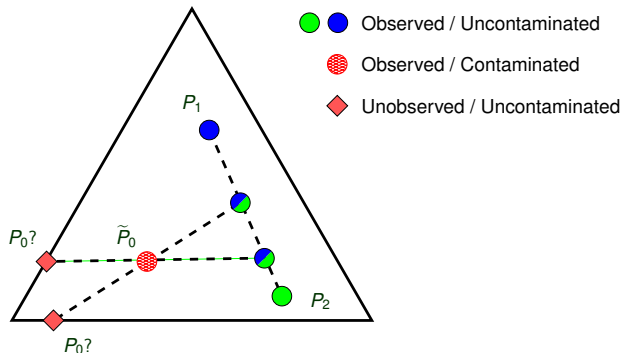
$$\tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0 \left(\sum_{i=1}^M \mu_i P_i \right) \quad \text{with } \sum_i \mu_i = 1$$

- ▶ When is it the case that $\kappa_0 = \kappa^*(\tilde{P}_0 | P_1, \dots, P_M)$?

IDENTIFIABILITY

$$\tilde{P}_0 = (1 - \kappa_0)P_0 + \kappa_0 \left(\sum_{i=1}^M \mu_i P_i \right) \quad \text{with} \quad \sum_i \mu_i = 1$$

- When is it the case that $\kappa_0 = \kappa^*(\tilde{P}_0 | P_1, \dots, P_M)$?



P_0 irreducible w.r.t. all $P_\mu, \mu \in S_M$ is **not** sufficient for identifiability

- \rightarrow joint irreducibility of $(P_i) \dots$

JOINT IRREDUCIBILITY

Call a family of distributions Q_1, \dots, Q_L **jointly irreducible** under either of the equivalent conditions:

- ▶ For any $I \subset \{1, \dots, L\}$; $1 \leq |I| \leq (L - 1)$:
any distribution in $\text{ConvHull} \{Q_i, i \in I\}$ is irreducible with respect to any distribution in $\text{ConvHull} \{Q_i, i \in I^c\}$
- ▶ If $\sum_{i=1}^L \gamma_i Q_i$ is a distribution, then $\gamma_i \geq 0$ for all i .
- ▶ If $\mathcal{M}_1(\mathcal{X})$ is the set of all probability distributions on \mathcal{X} ,

$$\mathcal{M}_1(\mathcal{X}) \cap \text{Span} \{Q_i, 1 \leq i \leq n\} = \text{ConvHull} \{Q_i, 1 \leq i \leq n\}$$

JOINT IRREDUCIBILITY – INTERPRETATION

Assume:

- ▶ (P_1, \dots, P_L) are jointly irreducible;
- ▶ $\tilde{P}_i = \pi_i^T \mathbf{P}$, with π_i (rows of the mixing matrix Π) linearly independent

Then:

$$\kappa^*(\tilde{P}_k | (\tilde{P}_i)_{i \in I}) = \kappa^*(\pi_k | (\pi_i)_{i \in I}),$$

and there is a one-to-one correspondance between the set of residues.

RECOVERABILITY

Recall the general contamination model:

$$\tilde{P}_i = \sum_{j=1}^L \pi_{ij} P_j \quad \iff \quad \tilde{P} = \Pi P$$

Call mixing weight matrix Π *recoverable* under either of the equivalent conditions:

- ▶ Π^{-1} has strictly positive diagonal entries and nonpositive off-diagonal entries
- ▶ For all ℓ , $\kappa^*(\pi_\ell | \{\pi_j, j \neq \ell\}) = \kappa_\ell$ is uniquely attained for decomposition

$$\pi_\ell = (1 - \kappa_\ell) \mathbf{e}_\ell + \kappa_\ell \pi'_\ell, \quad (*)$$

where π_ℓ is ℓ -th row of Π and $\mathbf{e} = \ell$ -th canonical basis vector, $\mathbf{e}_\ell = (0, \dots, 0, 1, 0, \dots, 0)$

DECONTAMINATION UNDER THE RECOVERABILITY ASSUMPTION

- ▶ Recoverability implies $\pi_\ell = (1 - \kappa_\ell)\mathbf{e}_\ell + \kappa_\ell\pi'_\ell$, unique maximal decomposition
- ▶ Irreducibility implies one-to-one correspondance, therefore

$$\tilde{P}_\ell = (1 - \kappa_\ell)P_\ell + \sum_{j \neq \ell} \nu_{\ell j} \tilde{P}_j;$$

unique maximal decomposition.

DECONTAMINATION UNDER THE RECOVERABILITY ASSUMPTION

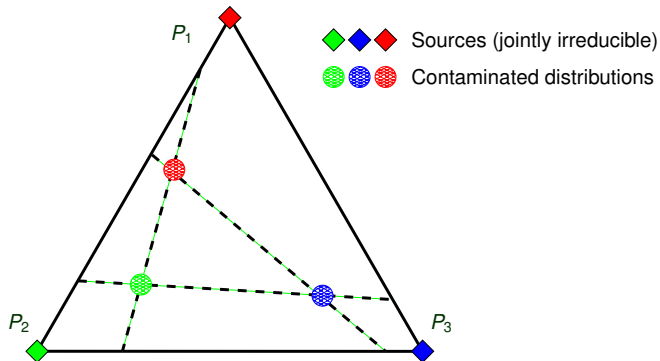
- ▶ Recoverability implies $\pi_\ell = (1 - \kappa_\ell)\mathbf{e}_\ell + \kappa_\ell\pi'_\ell$, unique maximal decomposition
- ▶ Irreducibility implies one-to-one correspondance, therefore

$$\tilde{P}_\ell = (1 - \kappa_\ell)P_\ell + \sum_{j \neq \ell} \nu_{\ell j} \tilde{P}_j;$$

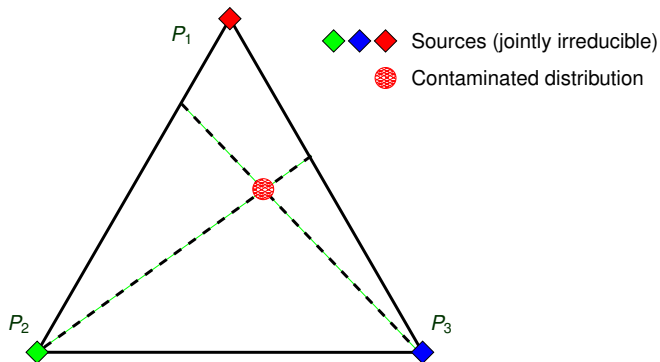
unique maximal decomposition.

- ▶ **Conclusion:** κ_ℓ can be estimated consistently by MMPE estimators $\hat{\kappa}(\tilde{P}_\ell | \{\tilde{P}_j, j \neq \ell\})$
- ▶ We estimate also consistently the sources P_ℓ (residues), and further $\nu_{\ell j}$, Π^{-1} and finally Π

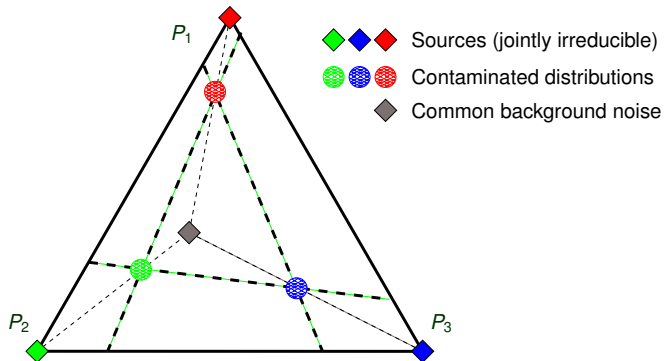
WHEN DOES RECOVERABILITY HOLD?



WHEN DOES RECOVERABILITY HOLD?



WHEN DOES RECOVERABILITY HOLD?



CONSISTENT ESTIMATION OF RISK

From

$$\tilde{P}_\ell = (1 - \kappa_\ell)P_\ell + \sum_{j \neq \ell} \nu_{\ell j} \tilde{P}_j$$

we get, denoting $\tilde{R}_{ij}(f) = \tilde{\mathbb{P}}_i(f(X) \neq j)$

$$R_\ell(f) = \frac{\tilde{R}_{\ell\ell}(f) - \sum_{j \neq \ell} \nu_{\ell j} \tilde{R}_{\ell j}}{1 - \kappa_\ell} \quad \longrightarrow \quad \hat{R}_\ell(f) = \frac{\hat{\tilde{R}}_{\ell\ell}(f) - \sum_{j \neq \ell} \hat{\nu}_{\ell j} \hat{\tilde{R}}_{\ell j}}{1 - \hat{\kappa}_\ell}$$

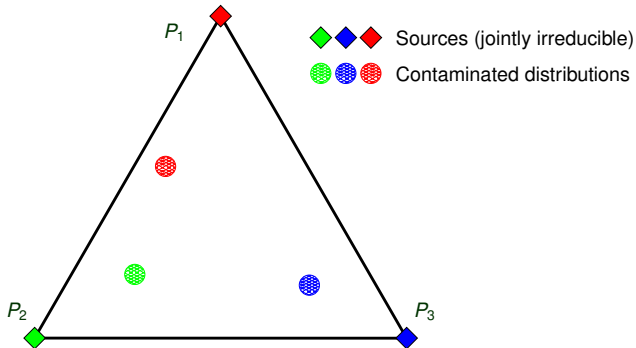
Then it holds:

$$\sup_{f \in \mathcal{F}_{k(n)}} \left| \hat{R}_\ell(f) - R_\ell(f) \right| \rightarrow 0 \text{ in probability,}$$

as $\mathbf{n} = \min(n_1, \dots, n_L) \rightarrow \infty$, for VC-classes \mathcal{F}_k of dimension V_k ,
provided $\frac{V_{k(n)} \log \mathbf{n}}{\mathbf{n}} \rightarrow 0$

DEMIXING WITHOUT RECOVERABILITY

- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



DEMIXING WITHOUT RECOVERABILITY

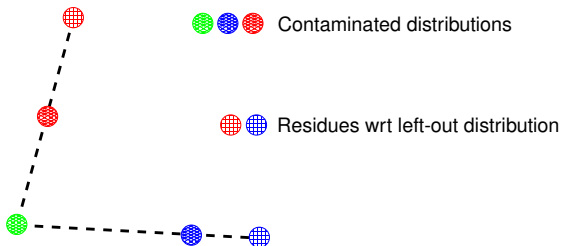
- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:

 Contaminated distributions



DEMIXING WITHOUT RECOVERABILITY

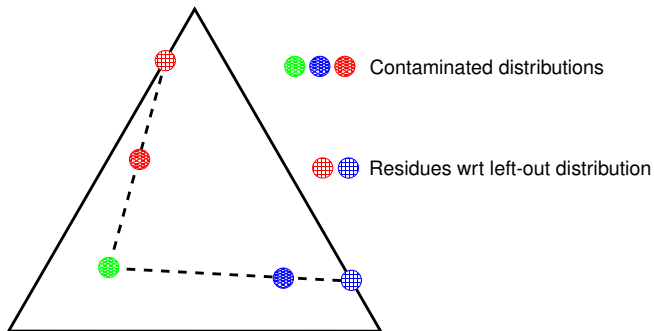
- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



- ▶ Residues **always** belong to the boundary

DEMIXING WITHOUT RECOVERABILITY

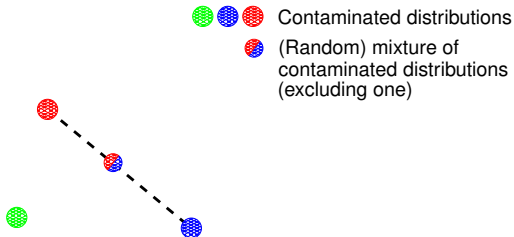
- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



- ▶ Residues **always** belong to the boundary
- ▶ Need a test of whether the residues belong to the same “face”

DEMIXING WITHOUT RECOVERABILITY

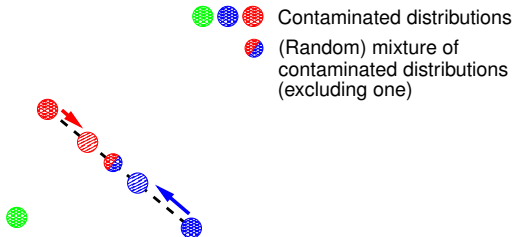
- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



- ▶ Residues **always** belong to the boundary
- ▶ Need a test of whether the residues belong to the same “face”

DEMIXING WITHOUT RECOVERABILITY

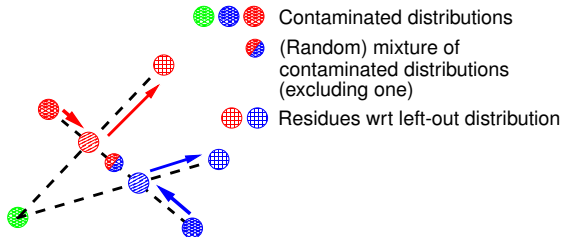
- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



- ▶ Residues **always** belong to the boundary
- ▶ Need a test of whether the residues belong to the same “face”

DEMIXING WITHOUT RECOVERABILITY

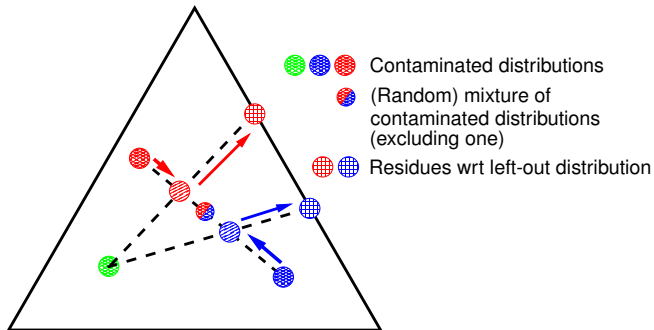
- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



- ▶ Residues **always** belong to the boundary
- ▶ Need a test of whether the residues belong to the same “face”

DEMIXING WITHOUT RECOVERABILITY

- ▶ **Goal:** estimating sources up to permutation (demixing problem)
- ▶ Try to “reduce dimension”:



- ▶ Residues **always** belong to the boundary
- ▶ Need a test of whether the residues belong to the same “face”
- ▶ If test does not reject, apply algorithm recursively

DEMIXING WITHOUT RECOVERABILITY

- ▶ **Advantage:** only need estimator $\hat{\kappa}$ for **two distributions** (much simpler to implement)
- ▶ **Advantage:** only need full column rank (weaker than recoverability) to establish population consistency
- ▶ **Disadvantage:** need more iterations/retries, theoretical consistency of estimation only established under the stronger assumption of

$$\forall i \quad \text{Supp}(\mathbf{P}_i) \not\subseteq \bigcup_{j \neq i} \text{Supp}(\mathbf{P}_j)$$

- ▶ **Extension:** If support \mathbf{S} of Π is known, and all columns of \mathbf{S} are unique, can recover the specific sources by support matching.


CONCLUSIONS

Contributions:

- ▶ Nonparametric/distribution-free point of view
- ▶ 2-class case: characterization of irreducible solution and consistent estimation
- ▶ Multiclass case:
 - ▶ Consistent maximal mixture proportions estimation
 - ▶ Consistent de-contamination under irreducibility + recoverability
 - ▶ Consistent de-mixing (up to permutation) under support irreducibility + full column rank
 - ▶ Consistent de-contamination under the same conditions as the previous point, if support of mixing weights known

THANK YOU FOR YOUR ATTENTION!

REFERENCES

-  G. Blanchard, G. Lee, C. Scott.
Semi-Supervised Novelty Detection.
Journal of Machine Learning Research 11: 2973-3009, 2010.
-  C. Scott, G. Blanchard, G. Handy.
Classification with Asymmetric Label Noise: Consistency and Maximal Denoising.
COLT, 2013
-  G. Blanchard, C. Scott.
Decontamination of mutually contaminated models.
AISTATS, 2014
-  C. Scott.
A Rate of Convergence for Mixture Proportion Estimation, with Application to Learning from Noisy Labels.
AISTATS, 2015.
-  J. Katz-Samuels, G. Blanchard, C. Scott.
Decontamination of Mutual Contamination Models
ArXiv: 1710.01167