

Label Shift Quantification

with Robustness Guarantees via Distribution Feature Matching

Dussap Bastien

Laboratoire de mathématiques d'Orsay
Université Paris-Saclay, Inria

May 30, 2023



Flow Cytometry

Principe de fonctionnement d'un analyseur-trieur

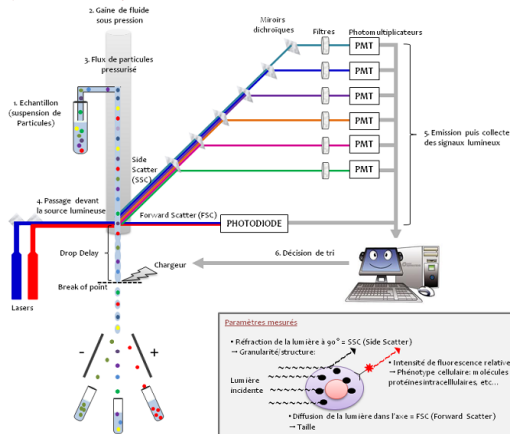


Figure: <https://bfa.u-paris.fr/cytometrie-en-flux-et-tri-cellulaire/>

Flow Cytometry

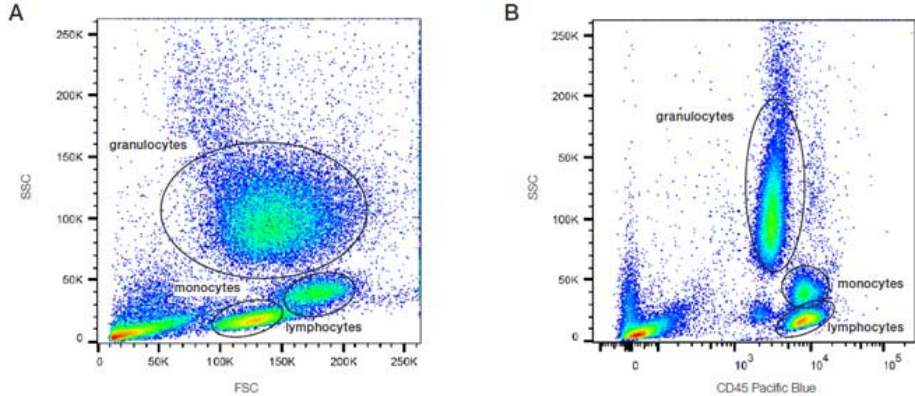


Figure: www.bio-rad-antibodies.com

Metaflow



Label Shift

Model

- \mathcal{X} : the data space.
- \mathcal{Y} : the label space, $\{1, \dots, c\}$.
- $\mathbb{P}_1, \dots, \mathbb{P}_c$: A list of c distributions, one for each population.
- $\mathbb{P} = \sum_{i=1}^c \beta_i \mathbb{P}_i$: A "source" distribution.
- $\mathbb{Q} = \sum_{i=1}^c \alpha_i^* \mathbb{P}_i$: A "target" distribution.

Datasets

Datasets

- $\{(x_j, y_j)\}_{j \in [n]} \in (\mathcal{X} \times \mathcal{Y})^n$, a labelled dataset : "source".
- $\{x_{n+j}\}_{j \in [m]} \in \mathcal{X}^m$, an unlabelled dataset : "target".
- $\hat{\mathbb{P}}_i := \frac{1}{n_i} \sum_{j \in [n]: y_j = i} \delta_{x_j}(\cdot)$.
- $\hat{\mathbb{Q}} := \frac{1}{m} \sum_{j=1}^m \delta_{x_{n+j}}(\cdot)$.
- $\tilde{\beta}_i := \frac{n_i}{n}$.
- $\tilde{\alpha}_i := \frac{m_i}{m} \neq \alpha_i^*$.

Quantification

Quantification

- Using $\hat{\mathbb{P}}_i$ and $\hat{\mathbb{Q}}$, estimate α^* .
- Using $\hat{\mathbb{P}}_i$ and $\hat{\mathbb{Q}}$, estimate $\tilde{\alpha}$.

Classify and Count

→ Use a classifier f .

Classifier

Count the number of times your classifier outputs each class.

Classify and Count (CC)

$$\hat{\alpha}_{cc} = \left(\frac{1}{m} \sum_{j=1}^m 1_{f(x_{n+j})=i} \right)_i$$

$$\alpha_{cc} = q(f(x) = i)_i$$

Black-Box Shift Estimator

→ Use the confusion matrix of a classifier.

$$\begin{aligned}q(f(x) = i) &= \sum_{j=1}^c q(f(x) = i | y = j) q(y = j) \\&= \sum_{j=1}^c p(f(x) = i | y = j) q(y = j) \\ \alpha_{cc} &= M_f \times \alpha^*\end{aligned}$$

Black-Box Shift Estimator

$$\hat{\alpha} = \hat{M}_f^{-1} \hat{\alpha}_{cc}$$

Black-Box Shift Estimator

Black-Box Shift Estimator

$$\hat{\alpha} = \hat{M}_f^{-1} \hat{\alpha}_{cc}$$

Alternative version: BBSE+

$$\hat{\alpha} = \arg \min_{\alpha \in \Delta^c} \|\alpha \hat{M}_f - \hat{\alpha}_{cc}\|_2$$

where $\Delta^c := \{x \in \mathbb{R}^c : x \geq 0, \sum x_i = 1\}$.

Kernel Mean Embedding

Kernel Methods

For a kernel k there exists a mapping $\Phi: \mathcal{X} \mapsto \mathcal{H}$ such as:

$$k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}$$

Kernel Mean Embedding

$$\begin{aligned}\Phi: \mathcal{M}_1^+(\mathcal{X}) &\rightarrow \mathcal{H} \\ \mathbb{P} &\mapsto \mathbb{E}_{X \sim \mathbb{P}}[\Phi(X)] = \Phi(\mathbb{P})\end{aligned}$$

Maximum Mean Discrepancy

Maximum Mean Discrepancy

$$\begin{aligned}\text{MMD}(\mathbb{P}, \mathbb{Q}) &= \|\Phi(\mathbb{P}) - \Phi(\mathbb{Q})\|_{\mathcal{H}} \\ &= \mathbb{E}_{\mathbb{P}, \mathbb{P}}[k(X, X)] + \mathbb{E}_{\mathbb{Q}, \mathbb{Q}}[k(Y, Y)] - 2\mathbb{E}_{\mathbb{P}, \mathbb{Q}}[k(X, Y)]\end{aligned}$$

Kernel Mean Matching

$$\arg \min_{\alpha \in \Delta^c} \text{MMD} \left(\sum_{i=1}^c \alpha_i \mathbb{P}_i, \mathbb{Q} \right) = \arg \min_{\alpha \in \Delta^c} \text{MMD} \left(\sum_{i=1}^c \alpha_i \mathbb{P}_i, \sum_{i=1}^c \alpha_i^* \mathbb{P}_i \right) \\ = \alpha^*$$

Kernel Mean Matching

$$\hat{\alpha} = \arg \min_{\alpha \in \Delta^c} \text{MMD} \left(\sum_{i=1}^c \alpha_i \hat{\mathbb{P}}_i, \hat{\mathbb{Q}} \right) \\ = \arg \min_{\alpha \in \Delta^c} \left\| \sum_{i=1}^c \alpha_i \Phi(\hat{\mathbb{P}}_i) - \Phi(\hat{\mathbb{Q}}) \right\|_{\mathcal{H}}$$

Kernel

Classical kernel

- $k(x, y) = x^T y$, linear.
- $k(x, y) = (\gamma x^T y + c_0)^d$, polynomial.
- $k(x, y) = \tanh(\gamma x^T y + c_0)$, sigmoid.
- $k(x, y) = \exp(\gamma \|x - y\|^2)$, gaussian.
- $k(x, y) = \exp(-\gamma \|x - y\|_1)$, laplacian.
- $k(x, y) = \|x\| + \|y\| - \|x - y\|$, energy.
- $k(x, y) = \left(1 + \frac{\|x - y\|^2}{\sigma^2}\right)^{-1}$, cauchy.

Distribution Feature Matching

Label Shift Quantification with Robustness Guarantees via Distribution Feature Matching

Bastien Dussap¹, Gilles Blanchard¹², and Badr-Eddine Chérif-Abdellatif³

¹ Université Paris-Saclay, Inria, Laboratoire de mathématiques d'Orsay

² Université Paris-Saclay, CNRS, Inria, Laboratoire de mathématiques d'Orsay

³ CNRS, LPSM, Sorbonne Université, Université Paris Cité

A general Approach

→ Let $\phi : \mathcal{X} \rightarrow \mathcal{F}$ be a fixed feature mapping from \mathcal{X} into a Hilbert space \mathcal{F} (possibly $\mathcal{F} = \mathbb{R}^D$).

Embedding

$$\begin{aligned}\phi : \mathcal{M}_1^+(\mathcal{X}) &\rightarrow \mathcal{F} \\ \mathbb{P} &\mapsto \mathbb{E}_{X \sim \mathbb{P}}[\phi(X)] = \phi(\mathbb{P})\end{aligned}$$

Distribution Feature Matching

Distribution Feature Matching (DFM)

$$\hat{\alpha} = \arg \min_{\alpha \in \Delta^c} \left\| \sum_{i=1}^c \alpha_i \phi(\hat{\mathbb{P}}_i) - \phi(\hat{\mathbb{Q}}) \right\|_{\mathcal{F}}$$

Black-Box Shift Estimator

$$\rightarrow \phi(x) = (1\{f(x) = i\})_{i=1,\dots,c} \in \mathbb{R}^c$$

Black-Box Shift Estimator

$$\rightarrow \phi(x) = (1\{f(x) = i\})_{i=1,\dots,c} \in \mathbb{R}^c$$

BBSE+ as DFM

- $\phi(\mathbb{Q})_i = q(f(x) = i) = \alpha_{cc}$
- $\phi(\mathbb{P}_j)_i = p(f(x) = i | y = j) = (M_f)_{i,j}$

Black-Box Shift Estimator

$$\rightarrow \phi(x) = (1\{f(x) = i\})_{i=1,\dots,c} \in \mathbb{R}^c$$

BBSE+ as DFM

- $\phi(\mathbb{Q})_i = q(f(x) = i) = \alpha_{cc}$
- $\phi(\mathbb{P}_j)_i = p(f(x) = i | y = j) = (M_f)_{i,j}$
- $\phi(\hat{\mathbb{P}}_j)_i = (\hat{M}_f)_i$
- $\phi(\hat{\mathbb{Q}})_i = \hat{\alpha}_{cc}$

Theoretical guarantees

Definition

For every set of vectors $\{b_1, \dots, b_c\}$ in a Hilbert space, denote $\Delta_{\min}(b_1, \dots, b_c)$, the second smallest eigenvalue of the centered Gram matrix:
$$M_{i,j}^c = \langle b_i - \frac{1}{c} \sum b_k, b_j - \frac{1}{c} \sum b_k \rangle.$$

Theoretical guarantees

Assumption

$$\sum_{i=1}^c \beta_i \phi(\mathbb{P}_i) = 0 \iff \beta = 0 \quad (\mathcal{A}_1)$$

and

$$\exists C > 0 : \|\phi(x)\| \leq C \text{ for all } x. \quad (\mathcal{A}_2)$$

Main Theorem

Theorem

For any $\delta \in (0, 1)$, with probability greater than $1 - \delta$:

$$\begin{aligned}\|\hat{\alpha} - \alpha^*\|_2 &\leq \frac{2CR_{(\delta+\log c)}}{\sqrt{\Delta_{\min}}} \left(\frac{\|w\|_2}{\sqrt{n}} + \frac{1}{\sqrt{m}} \right) \\ &\leq \frac{2CR_{(\delta+\log c)}}{\sqrt{\Delta_{\min}}} \left(\frac{1}{\sqrt{\min_i n_i}} + \frac{1}{\sqrt{m}} \right)\end{aligned}$$

where $R_x = 2 + \sqrt{2 \log(2/x)}$, $w_i = \frac{\alpha_i^*}{\tilde{\beta}_i}$, and $\Delta_{\min} := \Delta_{\min}(\phi(\hat{\mathbb{P}}_1), \dots, \phi(\hat{\mathbb{P}}_c))$.

The same result holds when replacing α^* by the (unobserved) vector of empirical proportions $\tilde{\alpha}$ in the target sample, both on the left-hand side and in the definition of w .

Properties of Δ_{\min}

Let (b_1, \dots, b_c) be a c -uple of vectors of \mathbb{R}^D assumed to be linearly independent. We denote M the Gram matrix of those vectors, i.e. $M_{ij} = \langle b_i, b_j \rangle$. We also write M^c the centered Gram matrix of the vectors : $(M^c)_{i,j} = \langle b_i - \bar{b}, b_j - \bar{b} \rangle$.

Theorem

For any number of classes c , Δ_{\min} is equal to $\min_{\substack{\|u\|=1 \\ \mathbf{1}^T u = 0}} u^T M u$.

→ $\Delta_{\min}(b_1, \dots, b_c)$ is always greater than the smallest eigenvalue of the Gram matrix.

Properties of Δ_{\min}

Theorem

In particular for two classes, $\Delta_{\min}(b_1, b_2) = \frac{1}{2}\|b_1 - b_2\|^2$.

Properties of Δ_{\min}

Theorem

In particular for two classes, $\Delta_{\min}(b_1, b_2) = \frac{1}{2}\|b_1 - b_2\|^2$.

Theorem (informal)

$\Delta_{\min}(b_1, \dots, b_c) \propto \text{vol}(\text{ConvHull}\{b_i, i \in [c]\})^2$.

Robustness to contamination

Robustness to contamination

Contaminated Label Shift

$$\mathbb{Q} = \sum_{i=1}^c \alpha_i^* \mathbb{P}_i + \alpha_0^* \mathbb{Q}_0 \quad (\mathcal{CLS})$$

Robustness to contamination

Contaminated Label Shift

$$\mathbb{Q} = \sum_{i=1}^c \alpha_i^* \mathbb{P}_i + \alpha_0^* \mathbb{Q}_0 \quad (\mathcal{CLS})$$

Soft Distribution Feature Matching (soft-DFM)

$$\hat{\alpha} = \arg \min_{\alpha \in \text{int}(\Delta^c)} \left\| \sum_{i=1}^c \alpha_i \phi(\hat{\mathbb{P}}_i) - \phi(\hat{\mathbb{Q}}) \right\|_{\mathcal{F}}^2,$$

where $\text{int}(\Delta^c) := \{x \in \mathbb{R}^c : x \geq 0, \sum x_i \leq 1\}$.

Main Theorem for soft-DFM

Theorem

For any $\delta \in (0, 1)$, with probability greater than $1 - \delta$:

$$\begin{aligned}\|\hat{\alpha} - \alpha^*\|_2 &\leq \frac{2CR_{(\delta+\log c)}}{\sqrt{\lambda_{\min}}} \left(\frac{\|w\|_2}{\sqrt{n}} + \frac{1}{\sqrt{m}} \right) \\ &\leq \frac{2CR_{(\delta+\log c)}}{\sqrt{\lambda_{\min}}} \left(\frac{1}{\sqrt{\min_i n_i}} + \frac{1}{\sqrt{m}} \right)\end{aligned}$$

where $R_x = 2 + \sqrt{2 \log(2/x)}$, $w_i = \frac{\alpha_i^*}{\beta_i}$, and $\lambda_{\min} := \lambda_{\min}(M)$.

The same result holds when replacing α^* by the (unobserved) vector of empirical proportions $\tilde{\alpha}$ in the target sample, both on the left-hand side and in the definition of w .

Theorem for \mathcal{CLS}

Introduce $\bar{V} := \text{Span}\{\Phi(\mathbb{P}_i), i \in [c]\}$ and let $\Pi_{\bar{V}}$ be the orthogonal projection on \bar{V} .

Corollary

With probability greater than $1 - \delta$:

$$\|\hat{\alpha} - \alpha^*\|_2 \leq \frac{1}{\sqrt{\Delta_{\min}}} \left(3\epsilon_n + \epsilon_m + \sqrt{2\alpha_0^* \epsilon_n \|\phi(\mathbb{Q}_0)\|} + \|\Pi_{\bar{V}}(\phi(\mathbb{Q}_0))\|_{\mathcal{F}} \right),$$

$$\epsilon_n = C \frac{R_{\delta + \log c}}{\sqrt{\min_i n_i}}; \quad \epsilon_m = C \frac{R_{\delta}}{\sqrt{m}};$$

Mixture of Gaussians

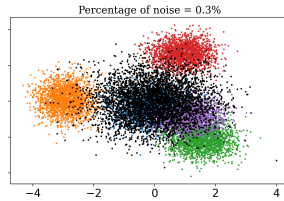
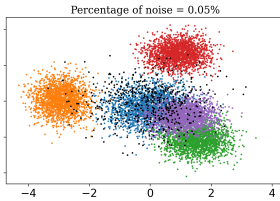
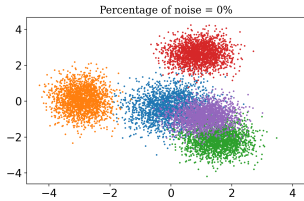
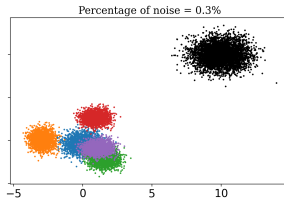
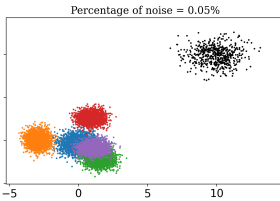
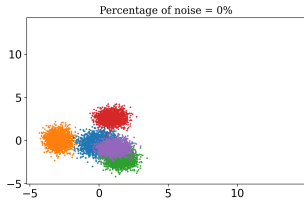
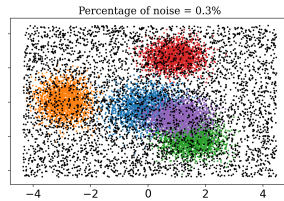
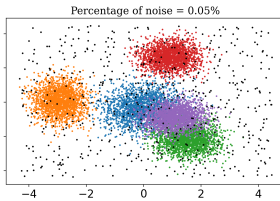
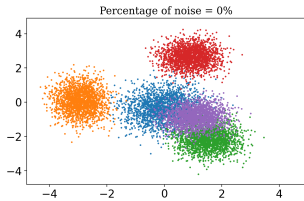
Contaminated Label Shift

$$\mathbb{Q} = \sum_{i=1}^c \alpha_i^* \mathbb{P}_i + \alpha_0^* \mathbb{Q}_0$$

The source is a list of c Gaussian distributions: $\mathbb{P}_1, \dots, \mathbb{P}_c$.

Contamination

1. Background uniform noise: \mathbb{Q}_0 is uniformly distributed over the data range.
2. New class **far** from the other distributions. In that case \mathbb{Q}_0 is Gaussian with a mean distant from the other means.
3. New class **close** to the other distributions. In that case \mathbb{Q}_0 is Gaussian with a similar mean to the others.



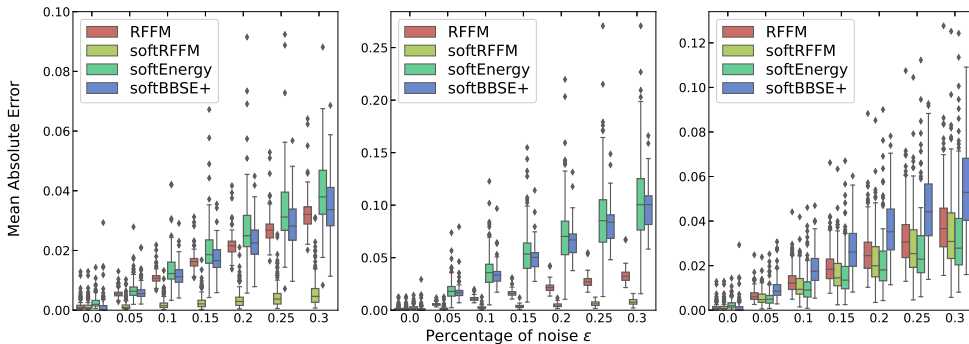
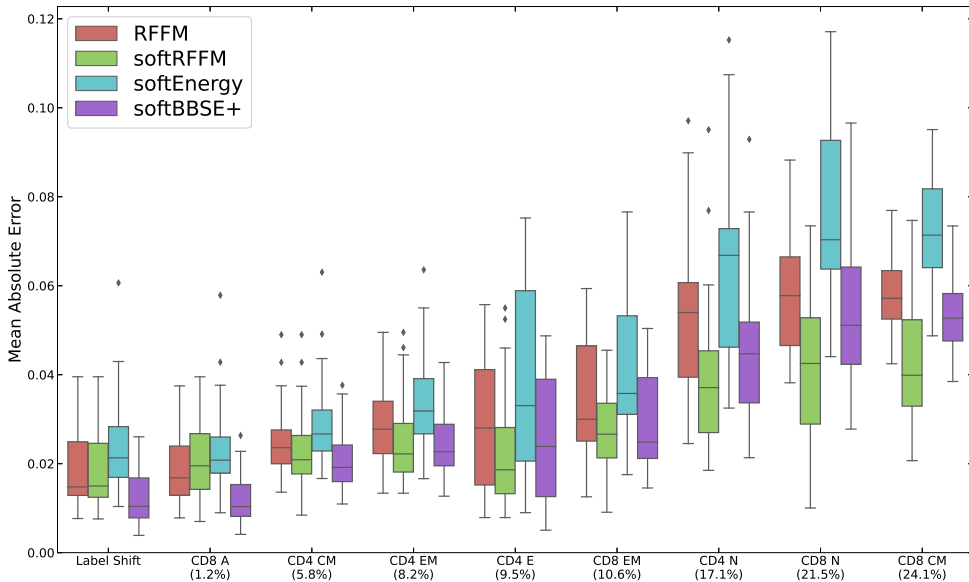


Figure: RFFM is Kernel Mean Matching with Random Fourier Features to speed up the computation.



Conclusion

- We introduced a general approach for Label Shift Quantification.
- We proposed to use Random Fourier Features to speed up the computation.
- We provided a general theoretical analysis of DFM, improving over previously known bounds derived for specific instantiations only.
- We analysed theoretically the behavior of DFM under a new setting and put into light a robustness against certain types of perturbations, depending on the feature mapping used.

Thank you for your attention.

References



González, Pablo et al. (2017). "A review on quantification learning". In: *ACM Computing Surveys (CSUR)* 50.5, pp. 1–40.



Gretton, Arthur et al. (2006). "A kernel method for the two-sample-problem". In: *Advances in neural information processing systems* 19.



Iyer, Arun, Saketha Nath, and Sunita Sarawagi (2014). "Maximum mean discrepancy for class ratio estimation: Convergence bounds and kernel selection". In: *International Conference on Machine Learning*. PMLR, pp. 530–538.



Kawakubo, Hideko, Marthinus Christoffel Du Plessis, and Masashi Sugiyama (2016). "Computationally efficient class-prior estimation under class balance change using energy distance". In: *IEICE Transactions on Information and Systems* 99.1, pp. 176–186.



Lipton, Zachary, Yu-Xiang Wang, and Alexander Smola (2018). "Detecting and correcting for label shift with black box predictors". In: *International conference on machine learning*. PMLR, pp. 3122–3130.



Rahimi, Ali and Benjamin Recht (2007). "Random features for large-scale kernel machines". In: *Advances in neural information processing systems* 20.