



Ensuring Fairness under Prior Probability Shifts

Arpita Biswas

Harvard University

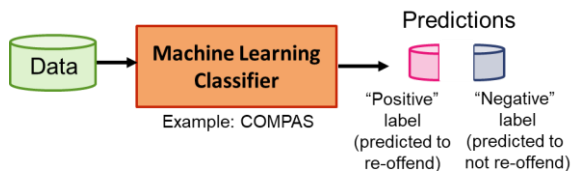
arpitabiswas@seas.harvard.edu

Suvam Mukherjee

Microsoft Corporation

sumukherjee@microsoft.com

Problem Statement



Problem: **COMPAS discriminates!** [Angwin et al. 2016]
More likely to wrongly deny bail to an African-American.

Prior Probability Shifts

The phenomenon where the prior probability $P(Y)$ changes between the training and test datasets, but the class conditional probability $P(X|Y)$ remains unaltered. [Moreno-Torres et al., 2012]

		(2013) Training Data True Prevalence	(2014) Test Data True Prevalence
	Z	ρ_d^z	ρ_d^z
Caucasian	0	0.327	0.636
African-American	1	0.486	0.706

Proportional Equality [Biswas and Mukherjee 2019]

The predictions on the test dataset d is defined to be fair if, for all z, z'

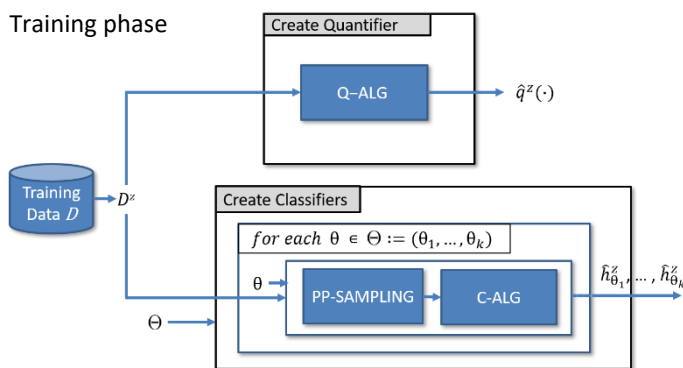
$$\left| \frac{\hat{\rho}_d^z}{\hat{\rho}_d^{z'}} - \frac{\rho_d^z}{\rho_d^{z'}} \right| \leq \epsilon$$

Goal:

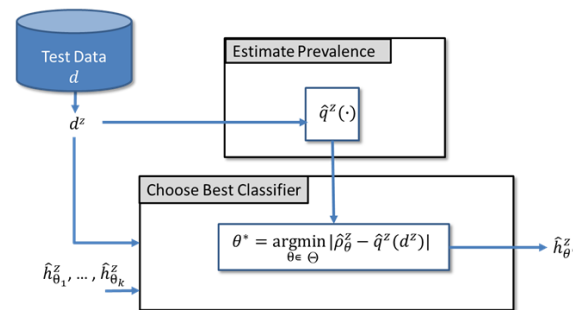
Maximize Accurate + fair + robust to prior probability shifts

Proposed Framework: CAPE

Training phase



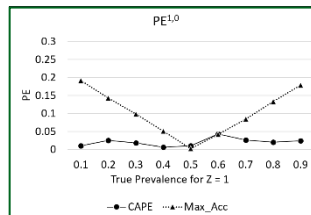
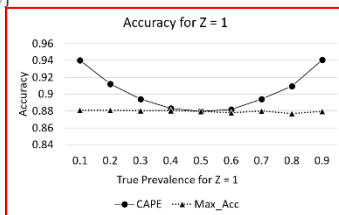
Prediction phase



Results

Theorem 1. Let $\Theta = \{\epsilon, 2\epsilon \dots, k\epsilon\}$, $\epsilon \in (0, 1)$, $k > \frac{1}{\epsilon} - 1$. For a group z , and test dataset d , let the quantifier be such that $|\rho_d^z - \hat{q}^z(d)| \leq \delta_1$, and the classifiers be such that $|\theta_j - \hat{\rho}_j^z| \leq \delta_2$, for small δ_1 and δ_2 . Then, for the best classifier $J = \arg \min_{j \in \{1, \dots, k\}} |\hat{\rho}_j^z - \hat{q}^z(d)|$, CAPE satisfies $|\rho_d^z - \hat{\rho}_d^z| \leq \epsilon + \delta_1 + \delta_2$.

Synthetic Datasets



CAPE outperforms Max_Acc when there is significant shift in prior probabilities. Max_acc degrades linearly with increasing shift, but CAPE performs well throughout.

Real-world Datasets	Method	Accuracy		PE
		Z=1	Z=0	(unfairness)
COMPAS	CAPE	69	64	0.08
	Max_acc	68	55	0.37
MEPS	CAPE	89	79	0.13
	Max_acc	89	76	0.48

CAPE does well on both fairness and accuracy.

More results in the paper: comparison against various state-of-the-art fair algorithms, and on several fairness metrics.