

Through the Fairness Lens: Experimental Analysis and Evaluation of Entity Matching







Nima Shahbazi, Nikola Danevski, <u>Abolfazl Asudeh</u>, Fatemeh Nargesian, Divesh Srivastava





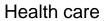
Agenda

- Motivation
- Part I: Fairness Evaluation in the context of EM
- Part II: Evaluation
 - 1. Entity-matching Approaches, datasets
 - 2. A few (representative) Evaluation Results
- Part III: Lessons and Discussions

(Direct) Social applications of EM









Things can go wrong! (esp. for marginalized people)

No-fly List





High-value customers list



How about EM for Union & Join?

Part I: (Group) Fairness Evaluation in the Context of EM

Fairness Measures for EM Evaluation

- Sensitive attributes (e.g., race, gender)
- Demographic Groups (e.g., female)
 - Non- intersectional (e.g., male) v.s. Intersectional groups (e.g. Black male)
 - Non-overlapping (e.g., gender) v.s. Overlapping groups (e.g., Research domain (DB, ML))

Single v.s. Pairwise Fairness

- Pairwise nature of EM (a differentiation from ML classification)
 - Each row in the test set: a record $c = (e_i, e_j, h, y)$
 - Single Fairness: fairness is evaluated for a group g
 - a record (in test-set) is "counted" for group g if at least one of the pairs belong to g
 - Pairwise Fairness: fairness is evaluated for a pair (g,g') group.

(left record) title: lineage tracing for general data warehouse transformations; author: jennifer widom , yingwei cui; venue: VLDBJ year: 2003

(right record) title: data extraction and transformation for the data warehouse; author: case squire; venue: SIGMOD; year: 1995

(left record) title: efficient and cost-effective techniques for browsing and indexing large video databases; author: kien a. hua , jung-hwan oh; venue: SIGMOT; year: 2000 (right record) title: effective timestamping in databases; author: kristian torp , christian s. jensen , richard thomas snodgrass; venue: VLDBJ year: 2000

(left record) title: efficient schemes for managing multiversionxml documents; author: shu-yao chien , carlo zaniolo , vassilis j. tsotras; venue: VLDBJ year: 2002 (right record) title: efficient management of multiversion documents by object referencing; author: shu-yao chien , vassilis j. tsotras , carlo zaniolo; venue VLDB; year: 2001

Name	Description	Equation $(\forall g_i \in \mathcal{G})$
Accuracy Parity (AP)	requires the independence of matchers's accuracy from groups	$Pr(h(e, e') = y g_i) \simeq Pr(h(e, e') = y)$
Statistical Parity (SP)	requires the independence of the matcher from groups	$Pr(h(e, e') = 'M' \mid g_i) \simeq Pr(h(e, e') = 'M')$
¹ True Positive Rate Parity (TPRP)	a.k.a <i>Equal Opportunity</i> ; in the group of true matches requires the independence of match predictions from groups	$Pr(h(e, e') = 'M' g_i, y = 'M') \simeq Pr(h(e, e') = 'M' y = 'M')$
False Positive Rate Parity (FPRP)	in the group of true non-matches, requires the independence of match predictions from groups	$Pr(h(e, e') = `M' g_i, y = `N') \simeq Pr(h(e, e') = `M' y = `N')$
¹ False Negative Rate Parity (FNRP)	in the group of true matches, requires the independence of non-match predictions from groups	$Pr(h(e, e') = 'N' g_i, y = 'M') \simeq Pr(h(e, e') = 'N' y = 'M')$
True Negative Rate Parity (TNRP)	in the group of true non-matches, requires the independence of non-match predictions from groups	$Pr(h(e, e') = 'N' g_i, y = 'N') \simeq Pr(h(e, e') = 'N' y = 'N')$
¹ Equalized Odds (EO)	in both groups of true matches and true non-matches requires the independence of match predictions from groups	$Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} g_i, y = {}^{\backprime}M{}^{\backprime}) \simeq Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} y = {}^{\backprime}M{}^{\backprime})$ $Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} g_i, y = {}^{\backprime}N{}^{\backprime}) \simeq Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} y = {}^{\backprime}N{}^{\backprime})$
¹ Positive Predictive Value Parity (PPVP)	among the pairs predicted as match requires the independence of true matches from groups	$Pr(y = 'M' h(e, e') = 'M', g_i) \simeq Pr(y = 'M' h(e, e') = 'M')$
¹ Negative Predictive Value Parity (NPVP)	among the pairs predicted as non-match, requires the independence of true non-matches from groups	$Pr(y = 'N' h(e, e') = 'N', g_i) \simeq Pr(y = 'N' h(e, e') = 'N')$
¹ False Discovery Rate Parity (FDRP)	e among the pairs predicted as match, requires the independence of true non-matches from groups	$Pr(y = 'N' g_i, h(e, e') = 'M') \simeq Pr(y = 'N' h(e, e') = 'M')$
¹ False Omission Rate Parity (FORP)	among the pairs predicted as non-match, requires the independence of true matches from groups	$Pr(y = 'M' g_i, h(e, e') = 'N') \simeq Pr(y = 'M' h(e, e') = 'N')$
		8

N 1 / /

Name	Description	Equation $\forall g_i \in \mathcal{G}$)
Accuracy Parity (AP)	requires the independence of matchers's accuracy from groups $% \left(1\right) =\left(1\right) \left(1\right)$	$Pr(h(e, e') = y g_i) \simeq Pr(h(e, e') = y)$
Statistical Parity (SP)	requires the independence of the matcher from groups	$Pr(h(e, e') = 'M' \mid g_i) \simeq Pr(h(e, e') = 'M')$
¹ True Positive Rate Parity (TPRP)	a.k.a <i>Equal Opportunity</i> ; in the group of true matches requires the independence of match predictions from groups	$Pr(h(e, e') = M' g_i, y = M') \simeq Pr(h(e, e') = M' y = M')$
False Positive Rate Parity (FPRP)	in the group of true non-matches, requires the independence of match predictions from groups	$Pr(h(e, e') = `M' g_i, y = `N') \simeq Pr(h(e, e') = `M' y = `N')$
¹ False Negative Rate Parity (FNRP)	in the group of true matches, requires the independence of non-match predictions from groups	$Pr(h(e, e') = 'N' g_i, y = 'M') \simeq Pr(h(e, e') = 'N' y = 'M')$
True Negative Rate Parity (TNRP)	in the group of true non-matches, requires the independence of non-match predictions from groups	$Pr(h(e, e') = `N' g_i, y = `N') \simeq Pr(h(e, e') = `N' y = `N')$
¹ Equalized Odds (EO)	in both groups of true matches and true non-matches requires the independence of match predictions from groups	$Pr(h(e, e') = {}^{`}M{}^{`} g_i, y = {}^{`}M{}^{`}) \simeq Pr(h(e, e') = {}^{`}M{}^{`} y = {}^{`}M{}^{`})$ $Pr(h(e, e') = {}^{`}M{}^{`} g_i, y = {}^{`}N{}^{`}) \simeq Pr(h(e, e') = {}^{`}M{}^{`} y = {}^{`}N{}^{`})$
¹ Positive Predictive Value Parity (PPVP)	among the pairs predicted as match requires the independence of true matches from groups	$Pr(y = 'M' h(e, e') = 'M', g_i) \simeq Pr(y = 'M' h(e, e') = 'M')$
¹ Negative Predictive Value Parity (NPVP)	among the pairs predicted as non-match, requires the independence of true non-matches from groups	$Pr(y = 'N' h(e, e') = 'N', g_i) \simeq Pr(y = 'N' h(e, e') = 'N')$
¹ False Discovery Rate Parity (FDRP)	among the pairs predicted as match, requires the independence of true non-matches from groups	$Pr(y = 'N' g_i, h(e, e') = 'M') \simeq Pr(y = 'N' h(e, e') = 'M')$
¹ False Omission Rate Parity (FORP)	among the pairs predicted as non-match, requires the independence of true matches from groups	$Pr(y = 'M' g_i, h(e, e') = 'N') \simeq Pr(y = 'M' h(e, e') = 'N')$
		9

3 7

N 1 /

Name	Description	Equation $(\forall g_i \in \mathcal{G})$					
Accuracy Parity (AP)	requires the independence of matchers's accuracy from groups	$Pr(h(e, e') = y g_i) \simeq Pr(h(e, e') = y)$					
Statistical Parity (SP)	requires the independence of the matcher from groups	$Pr(h(e, e') = 'M' \mid g_i) \simeq Pr(h(e, e') = 'M')$					
¹ True Positive Rate Parity (TPRP)	a.k.a <i>Equal Opportunity</i> ; in the group of true matches requires the independence of match predictions from groups	$Pr(h(e, e') = 'M' g_i, y = 'M') \simeq Pr(h(e, e') = 'M' y = 'M')$					
False Positive Rate Parity (FPRP)	Class Imbalance in EM:	$Pr(h(e, e') = `M' g_i, y = `N') \simeq Pr(h(e, e') = `M' y = `N')$					
¹ False Negative Rate Parity (FNRP)	among the $O(n^2)$ pairs,	$Pr(h(e, e') = 'N' g_i, y = 'M') \simeq Pr(h(e, e') = 'N' y = 'M')$					
True Negative Rate Parity (TNRP)	(usually) only $O(n)$ of	$Pr(h(e, e') = 'N' g_i, y = 'N') \simeq Pr(h(e, e') = 'N' y = 'N')$					
¹ Equalized Odds (EO)	them are match	$Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} g_i, y = {}^{\backprime}M{}^{\backprime}) \simeq Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} y = {}^{\backprime}M{}^{\backprime})$ $Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} g_i, y = {}^{\backprime}N{}^{\backprime}) \simeq Pr(h(e, e') = {}^{\backprime}M{}^{\backprime} y = {}^{\backprime}N{}^{\backprime})$					
¹ Positive Predictive Value Parity (PPVP)	among the pairs predicted as match requires the independence of true matches from groups	$Pr(y = 'M' h(e, e') = 'M', g_i) \simeq Pr(y = 'M' h(e, e') = 'M')$					
¹ Negative Predictive Value Parity (NPVP)	among the pairs predicted as non-match, requires the independence of true non-matches from groups	$Pr(y = 'N' h(e, e') = 'N', g_i) \simeq Pr(y = 'N' h(e, e') = 'N')$					
¹ False Discovery Rate Parity (FDRP)	among the pairs predicted as match, requires the independence of true non-matches from groups	$Pr(y = 'N' g_i, h(e, e') = 'M') \simeq Pr(y = 'N' h(e, e') = 'M')$					
¹ False Omission Rate Parity (FORP)	among the pairs predicted as non-match, requires the independence of true matches from groups	$Pr(y = 'M' g_i, h(e, e') = 'N') \simeq Pr(y = 'M' h(e, e') = 'N')$					
		10					

and the second s

N 1 /

Unfairness (Disparity) -- Division

$$F_{\alpha,\beta}^{(d)}(g_i) = \max\left(0, 1 - \frac{Pr(\alpha \mid \beta, g_i)}{Pr(\alpha \mid \beta)}\right)$$

E.g., True Positive Rate Parity (TPRP):

$$F = \max(0.1 - \frac{TPR(g_i)}{TPR(all)})$$

• Rule of thumb: unfairness more than 20% is not acceptable $F \le 0.2$

Part II: Evaluation

Entity Matchers

- Rule-based Matchers: 1
- Non-Neural Matchers: 7
- Neural Matchers: 5

Datasets

- Structured: 4
 - including 2 semi-syntheticsocial datasets
- Dirty: 2
- Textual: 2



We created two social datasets:

- NoFlyCompas: Used COMPAS dataset and created a no-fly list (impact of <u>over-representation</u>)
- FacultyMatch: Created using CSRankings faculty profiles in China and Germany (impact of <u>name similarities</u>)

NoFlyCompas

- W-B population in US ~ 75%-13%. In no-fly-list ~ 50%-50%
 (Black is over-represented)
- Non-neural matchers were dominant, both on overall performance and on fairness
 - Reason: better utilization of the structure of data.
 Also reported in [Mudgal2018]

	TPR		Disp	arity	F	DR	Disparity	
Matcher	Afr.	Cauc.	sub	div	Afr.	Cauc.	sub	div
DEEPMATCHER	0.89	0.86	-0.03	-0.03	0.20	0.18	0.02	0.11
Dітто	0.76	0.82	0.06	0.08	0.31	0.22	0.09	0.41
GNEM	0.84	0.84	0.00	0.00	0.17	0.09	0.08	0.88
HIERMATCHER	0.72	0.74	0.02	0.10	0.22	0.16	0.06	0.38
MCAN	0.54	0.57	0.03	0.05	0.19	0.05	0.14	2.8

A FP result by Ditto

(left record) firstName: James lastName: Brown race: African-American (right record) firstName: Samanthai lastName: Browne race: African-American

Faculty-Match

	TPR		Disp	arity	PI	PV	Disparity	
Matcher	cn	de	sub	div	cn	de	sub	div
DEEPMATCHER	0.48	0.72	0.23	0.50	0.79	0.87	0.08	0.11
Ditto	0.59	0.85	0.26	0.44	0.77	0.94	0.17	0.22
GNEM	0.78	0.90	0.12	0.15	0.83	0.92	0.08	0.11
HierMatcher	0.47	0.78	0.31	0.66	0.78	0.89	0.11	0.14
MCAN	0.40	0.70	0.30	0.75	0.86	0.94	0.08	0.09
DTMatcher	0.95	0.90	-0.05	-0.05	0.89	0.98	0.09	0.10
LinRegMatcher	0.33	0.23	-0.09	-0.43	0.44	0.96	0.52	1.18
LogRegMatcher	0.95	0.88	-0.07	-0.08	0.93	1.0	0.07	0.07
NBMATCHER	0.99	0.99	0.00	0.00	0.03	0.58	0.55	18.3
RFMATCHER	0.96	0.89	-0.06	-0.08	0.98	0.99	0.01	0.01
SVMMATCHER	0.95	0.87	-0.07	-0.09	0.94	0.99	0.05	0.05

A FP example (Ditto)

(left record) fullName: Qingming Huang country: cn (right record) fullName: Qing-Hu Huang country: cn

A FN example (Ditto)

(left record) fullName: LinLin Shen country: cn (right record) fullName: Linlin phen country: cn

Highlights from Extended Experiments

- Neural matchers were more unfair on structured data
 - Relying on pre-trained language models and embeddings

(left record) song: Tequila Loves Me; artist: K. Chesney (right record) song: Likes Me; artist: K. Chesney

 not fully considering the dataset structures

(left record) **title**: efficient schemes for managing multiversionxml documents; **author**: shu-yao chien , carlo zaniolo , vassilis j. tsotras; **venue**: VLDBJ; **year**: 2002 (right record) **title**: efficient management of multiversion documents by object referencing; **author**: shu-yao chien , vassilis j. tsotras , carlo zaniolo; **venue**: VLDB; **year**: 2001

- Neural matchers failed on Textual data
- Putting a high-weight on proxy attributes causes unfairness

```
(left record) title: guest editorial; author: alon y. halevy; venue: VLDBJ; year: 2002 (right record) title: guest editorial; author: vijay atluri, anupam joshi, yelena yesha; venue: VLDBJ; year: 2003
```

 Reason: Different "behaviors" between groups; Insufficient coverage for "minorities"

Sensitivity to the choice of matching threshold

		Non-neural						Neural				
	Dataset	DTMATCHER	LinRegMatcher	LogRegMatcher	NBMATCHER	RFMATCHER	SVMMATCHER	D ееРМАТСНЕR	DITTO	GNEM	HIERMATCHER	Mcan
	iTunes-Amazon	0	0	2.4	0	2.2	2.4	3.9	9.3	1	6.9	2.4
TPRP	CAMERAS	1	0	8.4	2.8	8.7	7.1	3.3	2.8	1	2.6	3.6
TP	ДВГБ-ЧСМ	0	0	0	0	0	0	0	2	0	0	0
	DBLP-SCHOLAR	0	0	0	0	0	1	2.4	2	0	2.2	2.4
_	iTunes-Amazon	0	0	0	0	2	0	1.7	5.2	0	2	1.4
\mathbf{ppvp}	Cameras	1	0	5.8	4.5	4.6	3.7	3.4	2.4	1.7	4.6	3.6
PP	ДВГБ-ЧСМ	0	0	0	0	0	2.6	0	0	0	0	0
	DBLP-SCHOLAR	0	0	1	1	1	1.4	1	1.4	0	2.4	1

Part III: Lessons and Discussions

- <u>Call for action</u> to collect EM benchmarks on societal applications
- Inherent issues in social data
 - Over/under-representation
 - Textual (name) similarity

We created two social datasets:

- NoFlyCompas: impact of <u>over-representation</u>
- FacultyMatch: impact of <u>name similarities</u>

Guide for EM Researchers and Practitioners

		Rules of Thumb					
	- Non-neural matchers are preferred - Obtain attributes with min correlation with sensitive attributes - Minimize Representation bias in training data - Make sure the model is not putting high weights on only a few attributes						
	Textual& durty datasets	 Neural matchers are preferred Obtain additional (unbiased) features Use unbiased pretrained models Minimize Representation bias in training data Considering their sensitivity, try out different matching thresholds and select the most fair/accurate one 					
\llbracket		ness measure: TPRP and PPVP are usually preferred (see § 3.5 and § 5.3.2)					
	Use an ensemble of matchers (for single sensitive attributes with exclusive values): construct a set of matchers; for each group use the matcher with best performance on it (using separate test sets for each group)						

Thank you!

