







anr®

# Fair Text Classification with Wasserstein Independence

T. Leteno<sup>1</sup>, A. Gourru<sup>1</sup>, C. Laclau<sup>2</sup>, R. Emonet<sup>1</sup>, C. Gravier<sup>1</sup>

<sup>1</sup> Laboratoire Hubert Curien, UMR CNRS 5516, Saint-Etienne, France <sup>2</sup> Télécom Paris, Institut Polytechnique de Paris, Paris, France

# EALM25 Ethic and Alignment of (Large) Language Models

June 30<sup>th</sup>, 2025





Fairness in Natural Language Processing. (Sun et al., 2019; Bender et al., 2021)

Public-ready Al-powered NLP systems (ChatGPT, Google Bard) expose the public to those bias.

## Example of fairness in classification

Sensitive group A
 Sensitive group B
 Hired
 Rejected



#### Toy example of classification Hiring process (Cannot reject candidates unfairly)

## Example of fairness in classification

Sensitive group A
Sensitive group B
Rejected



With **C**<sub>1</sub>:

$TPR_A = 9/10$	$FNR_A = 1/1$
$TPR_B = 7/10$	$FNR_{B} = 3/1$

Accuracy = 90.9%

Toy example of classification Hiring process (Cannot reject candidates unfairly)

\*TPR = rightly accepted, FNR = wrongly rejected

## Example of fairness in classification

Sensitive group A
Sensitive group B
Hired
Rejected



Toy example of classification Hiring process (Cannot reject candidates unfairly)

With **C**<sub>1</sub>:  $FNR_{A} = 1/10$  $TPR_{A} = 9/10$  $TPR_{B} = 7/10$  $FNR_{B} = 3/10$ Accuracy = 90.9%With  $C_2$ :  $TPR_{A} = 10/10$  $FNR_{A} = 0/10$  $TPR_{B} = 10/10$  $FNR_{B} = 0/10$ 

Accuracy = 88.6% (Loss in accuracy but no legitimate candidate is rejected)

\*TPR = rightly accepted, FNR = wrongly rejected

#### Example of fairness in classification in NLP

Bias for classification tasks :

- Hate speech detection (Baldini, et al., 2022; Huang et al., 2020; Davani et al. 2021)
- Job offers recommendations systems (Classification from biographies) (De-Arteaga et al., 2019)

'He went to medical school' → 'Surgeon' 'She went to medical school' → 'Nurse'

Impact on the jobs offers users receive.

## Motivation and theory



Fairness  $\approx$  Minimizing the Mutual Information (MI) between predictions and sensitive attributes.



Fairness  $\approx$  Minimizing the Mutual Information (MI) between predictions and sensitive attributes.

Notations:x : input texts : sensitive attribute $y / \hat{y} : label / prediction$  $h_y : classifier$ 



Fairness  $\approx$  Minimizing the Mutual Information (MI) between predictions and sensitive attributes.



→ Train a classifier while minimizing the Mutual Information between predictions and sensitive attributes.

#### min MI(ŷ, s)

 $\rightarrow$  Several locks to overcome.

#### Lock 1: Real world scenarios

In many real world scenarios sensitives attributes are not available (Veale et al, 2017). → Privacy concerns, Regulations (e.g. GPDR)

Most debiasing approaches require s.

#### Lock 1: Real world scenarios

In many real world scenarios sensitives attributes are not available (Veale et al, 2017). → Privacy concerns, Regulations (e.g. GPDR)

Most debiasing approaches require s.

Let  $\hat{s}$  be a predicted sensitives attributes by a classifier  $h_s$ ,  $h_s(x) = \hat{s}$ .

#### min MI( $\hat{y}$ , s) $\rightarrow$ min MI( $\hat{y}$ , $\hat{s}$ )

Need for data to train  $h_s$  to predict the *proxy* sensitive attributes :

 $\rightarrow$  subset of data annotated with *s*.

 $\rightarrow$  another source of data with transfer learning.

### Lock 2: Mutual Information Tractability

#### **Mutal Information**

Let  $\alpha$  and  $\beta$  be two random variables with joint distribution  $p(\alpha, \beta)$ , and  $p(\alpha)p(\beta)$  the product of marginal distributions.

 $MI(\alpha, \beta) = KL(p(\alpha, \beta) || p(\alpha)p(\beta))$ 

MI intractable for most real-life scenarios (sensitive to variation, theoretical limitations (McAllester & Stratos, 2020)).

## Lock 2: Mutual Information Tractability

#### **Mutal Information**

Let  $\alpha$  and  $\beta$  be two random variables with joint distribution  $p(\alpha, \beta)$ , and  $p(\alpha)p(\beta)$  the product of marginal distributions.

 $MI(\alpha, \beta) = \mathsf{KL}(p(\alpha, \beta) || p(\alpha)p(\beta))$ 

MI intractable for most real-life scenarios (sensitive to variation, theoretical limitations (McAllester & Stratos, 2020)).

Wasserstein Dependency Measure (Ozair et al., 2019)

 $MI_{W}(\alpha, \beta) = W_{1}(p(\alpha, \beta), p(\alpha)p(\beta))$ 

with  $W_1$  the Wasserstein-1 distance.

#### Lock 2: Mutual Information Tractability

Wasserstein Dependency Measure (Ozair et al., 2019)

 $MI_{W}(\alpha, \beta) = W_{1}(p(\alpha, \beta), p(\alpha)p(\beta))$ 

with  $W_1$  the Wasserstein-1 distance.

min MI( $\hat{y}$ ,  $\hat{s}$ )  $\rightarrow$  min MI<sub>w</sub>( $\hat{y}$ ,  $\hat{s}$ )

## Lock 3: Differentiability of the argmax operator

#### min MI<sub>w</sub>(ŷ, ŝ)

→ Wasserstein-1 distance approximated by a neural network called Critic taking as inputs  $\hat{y}$  and  $\hat{s}$ .

Inspired by the Wasserstein-GAN literature(Arjovsky et al., 2017).

However,  $\hat{\mathbf{y}} = \operatorname{argmax}(\mathbf{h}_{\mathbf{y}}(\mathbf{x})) \rightarrow \operatorname{argmax}$  operation is not differentiable.

## Lock 3: Differentiability of the argmax operator

#### min MI<sub>w</sub>(ŷ, ŝ)

→ Wasserstein-1 distance approximated by a neural network called Critic taking as inputs  $\hat{y}$  and  $\hat{s}$ .

Inspired by the Wasserstein-GAN literature(Arjovsky et al., 2017).

However,  $\hat{\mathbf{y}} = \operatorname{argmax}(\mathbf{h}_{\mathbf{y}}(\mathbf{x})) \rightarrow \operatorname{argmax}$  operation is not differentiable.

Let :

- $z_y$  be the hidden representations of classifier  $h_y$ .
- $z_s$  be the hidden representations of classifier  $h_s$  (classifier predicting s).

min  $MI_w(\hat{y}, \hat{s}) \rightarrow min MI_w(z_y, z_s)$ 

#### min MI(ŷ, s)

**Objective:** Training a classification model while minimizing  $MI(\hat{y}, s)$ .

```
Unobserved sensitive
attributes min MI(ŷ, ŝ)
min MI(ŷ, s)
```

**Objective:** Training a classification model while minimizing  $MI(\hat{y}, \hat{s})$ .



**Objective:** Training a classification model while minimizing  $MI_W(\hat{y}, \hat{s})$ .



**Objective:** Training a classification model while minimizing  $MI_W(z_y, z_s)$ .



Our approximations upper-bound the original measures:

- s  $\rightarrow$  ŝ, creates an approximation due to the error of the model predicting ŝ.
- $(\hat{y}, \hat{s}) \rightarrow (z_y, z_s)$  creates a approximation due to the error introduced by the softmax.
- Wasserstein Dependency Measure is related with common fairness metrics.

## Architecture

**Objective:** Training a classification model while minimizing the Wasserstein Dependency Measure between the representations.















## Experiments



#### 1. Bias in Bios dataset (De-Arteaga et al., 2019)

"Mr. Miserez devotes a substantial portion of his practice to representing healthcare entities in the defense of RAC, Medicare, Medicaid and other third party payor audits. He has significant experience defending hospitals, home health agencies, and other healthcare providers (...)"

Label Attorney (occupation) Sensitive attribute Male (gender)

#### 2. Moji dataset (Blodgett et al., 2016)

"Why am I struggling at work right now"

Label Negative sentiment Sensitive attribute Standard-American English (type of language)

"Sleep Real Good Cuss Ain Got No Worries"

Label Positive sentiment Sensitive attribute African-American English (type of language)

#### **Evaluation criteria**

Performance on the classification task: accuracy.

Fairness: difference of true positive rate between each sensitive group.

Performance-fairness trade-off: distance to optimum (utopia point).



#### Main results

Model	Accuracy ↑	Fairness ↑	DTO ↓	Model	Accuracy ↑	Fairness †	DTO ↓
*CE	72.3 ± 0.5	$61.2 \pm 1.4$	31.0	*CE	82.3 ± 0.2	$85.1 \pm 0.8$	5.67
INLP	73.3 ± 0.0	85.6 ± 0.0	7.02	INLP	82.3 ± 0.0	$88.6 \pm 0.0$	2.44
Adv	$75.6 \pm 0.4$	$90.4 \pm 1.1$	1.71	Adv	$81.9 \pm 0.2$	90.6 ± 0.5	1.80
Gate	76.2 ± 0.3	90.1 ± 1.5	1.90	Gate	83.7 ± 0.2	$90.4 \pm 0.9$	0.20
FairBatch	$75.1 \pm 0.6$	90.6 ± 0.5	1.78	FairBatch	82.2 ± 0.1	89.5 ± 1.3	1.86
EOGLB	75.2 ± 0.2	$90.1 \pm 0.4$	2.15	EOGLB	$81.7 \pm 0.4$	$88.4 \pm 1.0$	2.97
Condp	75.8 ± 0.3	88.1 ± 0.6	3.92	Condp	82.1 ± 0.2	$84.3 \pm 0.8$	6.50
Coneo	$74.1 \pm 0.7$	84.1 ± 3.0	8.17	Coneo	$81.8 \pm 0.3$	$85.2 \pm 0.4$	5.72
WFC	$75.2 \pm 0.1$	91.4 ± 0.3	1.17	WFC	$82.4 \pm 0.1$	89.0 ± 0.3	2.06

#### a) Moji dataset

b) Bias in Bios dataset

Table : Model's evaluation. For baselines, results are drawn from (Shen et al., 2022a). We report the mean ± standard deviation over 5 runs. \* indicates the model without fairness consideration.

**Proposition** Disentanglement of representations to train a fair classifier.

Performance Competitive or better than state-of-the-art baselines

Advantages

Applicable when datasets lack of sensitive attribute annotations.

Generalizable to other encoder-decoder architecture and to other sensitive attributes (continuous).



## More details in the paper Fair Text Classification with Wasserstein Independence.

#### Contact

thibaud.leteno@univ-st-etienne.fr antoine.gourru@univ-st-etienne.fr



DIKÉ project website

## Bibliography

Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, pages 1630–1640.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In FAccT, pages 610–623.

Sherjil Ozair, Corey Lynch, Yoshua Bengio, Aaron Van den Oord, Sergey Levine, and Pierre Sermanet. 2019. Wasserstein dependency measure for representation learning. Advances in Neural Information Processing Systems, 32.

Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein generative adversarial networks. In ICML, pages 214–223. PMLR

Su Lin Blodgett, Lisa Green, and Brendan O'Connor. 2016. Demographic dialectal variation in social media: A case study of african-american english. ArXiv preprint arXiv:1608.08868.

Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In FaccT, pages 120–128.

Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. Marked personas: Using natural language prompts to measure stereotypes in language models. ArXiv preprint arXiv:2305.18189.

Svetlana Kiritchenko and Saif M. Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems.

## Bibliography

Soares, Ioana Baldini, et al. "Your fairness may vary: pretrained language model fairness in toxic text classification." Annual Meeting of the Association for Computational Linguistics. 2022.

Huang, Xiaolei, et al. "Multilingual Twitter Corpus and Baselines for Evaluating Demographic Bias in Hate Speech Recognition." Proceedings of the Twelfth Language Resources and Evaluation Conference. 2020.

Davani, Aida Mostafazadeh, et al. "Improving Counterfactual Generation for Fair Hate Speech Detection." Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021). 2021.

Veale, M., and R. Binns. "Fairer machine learning in the real world: Mitigating discrimination without collecting sensitive data." Big Data and Society 4.2 (2017).

#### Transfer of sensitive attributes



Figure : simplified representation of the method architecture.

#### Transfer of sensitive attributes

![](_page_39_Figure_1.jpeg)

#### Figure : representation of the transfer of sensitive attributes.

#### Transfer of sensitive attributes

Dataset	Accuracy ↑	Fairness ↑	DTO ↓	Leakage↓	
Bios	82.4 ± 0.1	89.0 ± 0.3	2.06	96.5 ± 0.5	
EEC	82.2 ± 0.4	$88.9 \pm 0.4$	2.26	97.5 ± 0.3	Table : Comparison between several
MP	82.4 ± 0.3	88.9 ± 0.4	2.14	96.4 ± 0.5	for prediction on Bias in Bios.

#### Other datasets

Equity Evaluation Corpus (EEC): synthetic dataset for sentiment analysis with gender bias. (Kiritchenko and Mohammad, 2018)

Marked Personas dataset (MP): description of personas generated by NLP systems. (Cheng et al., 2023)

#### Wasserstein dependency measure

min MI( $\hat{y}, \hat{s}$ )  $\rightarrow$  min MI<sub>w</sub>( $\hat{y}, \hat{s}$ )

Wasserstein distance approximation (Arjovsky et al., 2017)

Use of a neural network called Critic  $C_{\boldsymbol{\omega}}.$ 

•  $W_1(p(\alpha, \beta), p(\alpha)p(\beta)) = \sup_{\omega, \|C_{\omega}\|_{L} \leq 1} E_{\alpha,\beta\sim p(\alpha,\beta)}[C_{\omega}(\alpha, \beta)] - E_{\alpha\sim p(\alpha),\beta\sim p(\beta)}[C_{\omega}(\alpha, \beta)]$ 

where  $||C_{\omega}||_{L}$  is the set of all 1-Lipschitz functions.

• Optimization : max  $E_{\hat{y},\hat{s}\sim p(\hat{y},\hat{s})}[C_{\omega}(\hat{y},\hat{s})] - E_{\hat{y}\sim p(\hat{y}),\hat{s}\sim p(\hat{s})}[C_{\omega}(\hat{y},\hat{s})]$