

Mitigating Gender Bias in Face Recognition using the von Mises-Fisher Mixture Model

Jean-Rémy Conti
Télécom Paris, IDEMIA

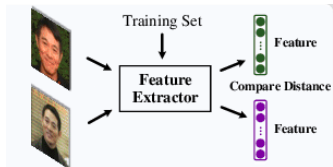
October 12, 2022

- Many applications of Face Recognition: access control, identity verification (smartphones, suspects), social media ...
- Bias with respect to race, gender, age, ...
- Different causes of bias, popular subject in machine learning

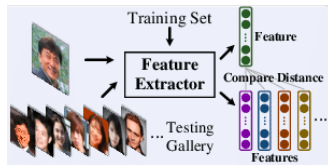
- How to recognize a face ?
- How to measure fairness ?
- How to mitigate gender bias ?

Face Recognition : short introduction

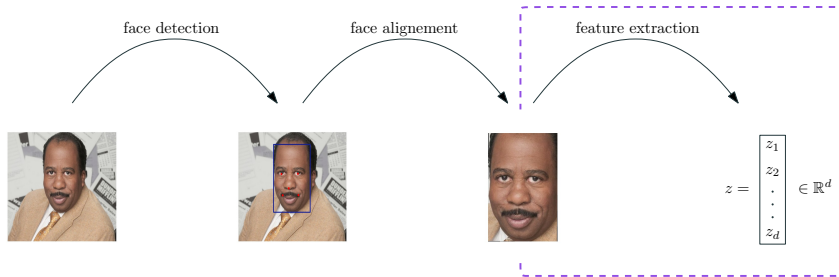
- Face verification:



- Face identification:



The steps in Face Recognition

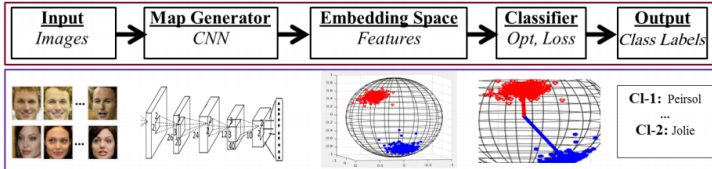


Depuis 2014 : réseau neuronal convolutif

Goal : Make the latent representations from a same identity as close as possible in the latent space.

Ingredients : training set, architecture of neural network (feature extractor), loss function.

Face Recognition training



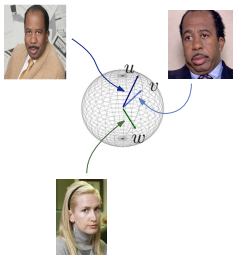
Workflow of Deep Face Recognition training.

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\kappa} \mu_{y_i}^T \mathbf{x}_i}{\sum_{k=1}^C e^{\kappa} \mu_k^T \mathbf{x}_i} \quad \|\mathbf{x}_i\|_2 = \|\mu_k\|_2 = 1$$

Face Recognition (Verification)

Face Recognition systems use face embeddings which are normalized (they lie on the hypersphere \mathbb{S}^{d-1}).

The similarity between two faces is usually measured by the cosine similarity $\langle u, v \rangle = \frac{u^T v}{\|u\|_2 \|v\|_2}$.



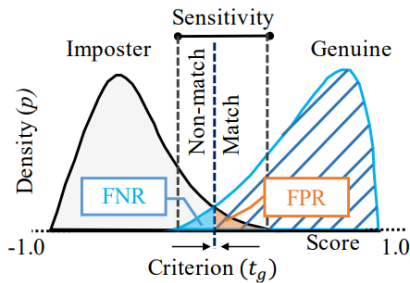
Decision rule : $t \in [-1, 1]$, fixed threshold.

- $\langle u, v \rangle \geq t \Rightarrow$ same identity (genuine),
- $\langle u, w \rangle < t \Rightarrow$ distinct identities (impostor).

Evaluation Metric

Two kinds of errors:

- False Positives : predicting "same identity" for two faces from distinct identities. \rightsquigarrow False Acceptance Rate: $FAR(t)$.
- False Negatives : predicting "distinct identities" for two faces from a same identity. \rightsquigarrow False Rejection Rate: $FRR(t)$.



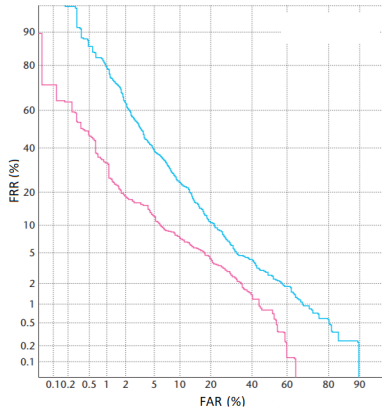
In practice :

1. A threshold $t \in [-1, 1]$ is set to get a deemed acceptable security level α for $\text{FAR}(t)$.
2. The False Rejection Rate is computed at this threshold:

$$\text{FRR}@\text{(FAR} = \alpha) := \text{FRR}(t), \text{ where } \text{FAR}(t) = \alpha.$$

Typically $\alpha = 10^{-1}, 10^{-2}, \dots, 10^{-8}$.

DET/ROC curve



Two typical ROC curves.

Demographic differentials

How to Measure Fairness ?

Context

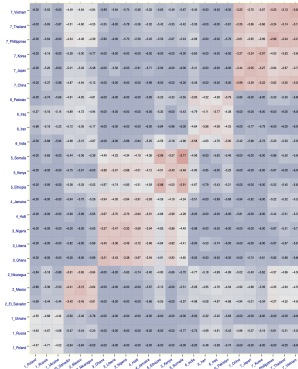
Do Face Recognition algorithms have uniform performance among the population ?

- \mathcal{G} : set of subgroups of the population.
Examples : women, men, young, old ...
- For all $g \in \mathcal{G}$, we can compute $\text{FAR}_g(t)$ and $\text{FRR}_g(t)$, the False Acceptance and False Rejection Rates, specific to subgroup g .

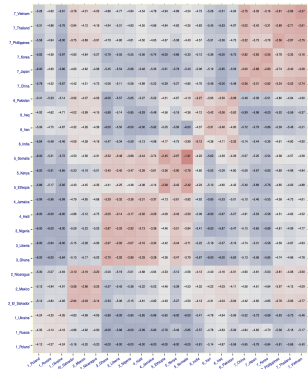
The *National Institute of Standards and Technology* regularly evaluates face recognition algorithms. On their performance ...

Algorithm	Constrained, Cooperative					
	FMR	= 0.000001	= 0.00001	= 0.00001	= 0.000001	= 0.000001
	Submission Date	VISA Photos	MUGSHOT Photos	MUGSHOT Photos 12+YRS	VISABORDER Photos	BORDER Photos
sunsetime-005	2021-05-24	0.0029 ⁽¹⁷⁾	0.0022 ⁽¹¹⁾	0.0021 ⁽¹¹⁾	0.0023 ⁽¹¹⁾	0.0044 ⁽¹¹⁾
visionlabs-011	2021-10-13	0.0022 ⁽⁷⁾	0.0024 ⁽⁹⁾	0.0026 ⁽⁷⁾	0.0028 ⁽²⁾	0.0053 ⁽²⁾
ntechlab-011	2021-09-13	0.0019 ⁽⁴⁾	0.0024 ⁽⁸⁾	0.0028 ⁽¹⁸⁾	0.0029 ⁽³⁾	0.0055 ⁽³⁾
clearviewai-000	2021-09-22	0.0019 ⁽⁵⁾	0.0024 ⁽⁵⁾	0.0028 ⁽¹¹⁾	0.0030 ⁽⁴⁾	0.0058 ⁽⁵⁾
mendaxiatech-000	2021-09-15	0.0036 ⁽²⁵⁾	0.0029 ⁽²⁶⁾	0.0036 ⁽³²⁾	0.0031 ⁽⁵⁾	0.0057 ⁽⁴⁾
ntechlab-010	2021-04-30	0.0017 ⁽²⁾	0.0024 ⁽¹⁰⁾	0.0029 ⁽²⁰⁾	0.0031 ⁽⁶⁾	0.0058 ⁽⁶⁾
cubox-002	2021-08-24	0.0041 ⁽³⁵⁾	0.0025 ⁽¹¹⁾	0.0025 ⁽⁶⁾	0.0033 ⁽⁷⁾	0.0064 ⁽⁹⁾
visionlabs-010	2021-01-25	0.0024 ⁽⁹⁾	0.0026 ⁽¹⁶⁾	0.0030 ⁽²²⁾	0.0033 ⁽⁸⁾	0.0061 ⁽⁷⁾
toshiba-004	2021-09-27	0.0042 ⁽⁴⁰⁾	0.0025 ⁽¹²⁾	0.0027 ⁽⁹⁾	0.0034 ⁽⁹⁾	0.0063 ⁽⁸⁾
idemia-008	2021-07-07	0.0032 ⁽¹⁹⁾	0.0023 ⁽³⁾	0.0028 ⁽¹⁰⁾	0.0034 ⁽¹⁰⁾	0.0067 ⁽¹¹⁾
kakao-006	2021-10-13	0.0029 ⁽¹⁶⁾	0.0024 ⁽⁴⁾	0.0028 ⁽¹⁶⁾	0.0035 ⁽¹¹⁾	0.0065 ⁽¹⁰⁾
insightface-001	2021-09-27	0.0014 ⁽¹⁾	0.0027 ⁽¹⁹⁾	0.0024 ⁽³⁾	0.0035 ⁽¹²⁾	0.0070 ⁽¹³⁾
paravision-008	2021-07-01	0.0025 ⁽¹¹⁾	0.0024 ⁽⁶⁾	0.0025 ⁽⁵⁾	0.0036 ⁽¹³⁾	0.0070 ⁽¹⁴⁾
insightface-000	2021-03-17	0.0027 ⁽¹⁵⁾	0.0029 ⁽²⁵⁾	0.0030 ⁽²³⁾	0.0038 ⁽¹⁴⁾	0.0077 ⁽²⁰⁾

... And on their performance differentials with respect to some subgroups of the population !



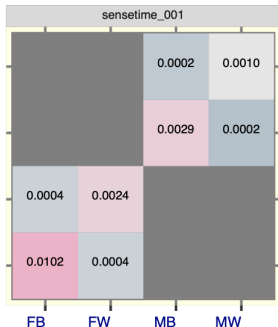
FAR for men.



FAR for women.

NIST reports

... And on their performance differentials with respect to some subgroups of the population !



FAR for ethnicity+gender subgroups.
F: female, M: male, B: black, W: white.

↪ Some algorithms make 10 times more errors on black women than on white men.

How to Measure Fairness ?

Context

Some algorithms make 10 times more errors on black women than on white men¹.

- \mathcal{G} : set of subgroups of the population.
Examples : women, men, young, old ...
- For all $g \in \mathcal{G}$, we can compute $\text{FAR}_g(t)$ and $\text{FRR}_g(t)$, the False Acceptance and False Rejection Rates, specific to subgroup g .

¹Grother et al. *Ongoing face recognition vendor test (frvt) part 3: Demographic effects?* NIST, 2019.

How to Measure Fairness ?

Context

- \mathcal{G} : set of subgroups of the population.
- For all $g \in \mathcal{G}$, we can compute $\text{FAR}_g(t)$ and $\text{FRR}_g(t)$, the False Acceptance and False Rejection Rates, specific to subgroup g .

Our new fairness metrics

1. Two ratios \rightsquigarrow interpretable metrics:

$$\frac{\max_g \text{FAR}_g(t)}{\min_g \text{FAR}_g(t)} \quad \text{and} \quad \frac{\max_g \text{FRR}_g(t)}{\min_g \text{FRR}_g(t)}$$

¹Grother et al. *Ongoing face recognition vendor test (frvt) part 3: Demographic effects?* NIST, 2019.

How to Measure Fairness ?

Context

- \mathcal{G} : set of subgroups of the population.
- For all $g \in \mathcal{G}$, we can compute $\text{FAR}_g(t)$ and $\text{FRR}_g(t)$, the False Acceptance and False Rejection Rates, specific to subgroup g .

Our new fairness metrics

1. Two ratios \rightsquigarrow interpretable metrics:

$$\text{BFAR}(\alpha) = \frac{\max_g \text{FAR}_g(t)}{\min_g \text{FAR}_g(t)} \quad \text{and} \quad \text{BFRR}(\alpha) = \frac{\max_g \text{FRR}_g(t)}{\min_g \text{FRR}_g(t)}$$

2. The threshold t satisfies $\max_{g \in \mathcal{G}} \text{FAR}_g(t) = \alpha$ instead of $\text{FAR}_{\text{total}}(t) = \alpha$. \rightsquigarrow more robust to a change of evaluation dataset

¹Grother et al. *Ongoing face recognition vendor test (frvt) part 3: Demographic effects?* NIST, 2019.

Bias Mitigation in Face Recognition

Survey of existing methods

Pre-training : reweighting / augmentation

- **Balanced Datasets Are Not Enough**, Wang and al. 2019.
- **How Does Gender Balance In Training Data Affect Face Recognition Accuracy?**, Albiero and al. 2020.

↪ Not yet adapted to Face Recognition.

Adversarial methods during training

- **PASS: Protected attribute suppression system for mitigating bias in face recognition**, Dhar and al. 2021.
- **How Does Gender Balance In Training Data Affect Face Recognition Accuracy?**, Albiero and al. 2020.

↪ Costly in computing time and unstable.

Post-training methods : modification of matching scores

- **Bias mitigation of face recognition models through calibration.**, Salvador and al. 2021.

↪ Does not solve the problem at the root level.

Core idea : add a shallow neural network on the last layer of a pre-trained model, in order to correct its gender bias.

**Mitigating Gender Bias in Face Recognition
Using the von Mises-Fisher Mixture Model**

Jean-Rémy Conti^{*1,2} Nathan Noiry^{*1} Vincent Despiegel² Stéphane Gentric² Stéphan Cléménçon¹

Accepted at ICML 2022 conference.

Geometric Embedding View on Fairness

Women are disadvantaged compared to men in terms of both FAR and FRR.

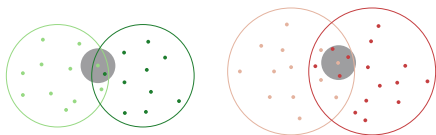


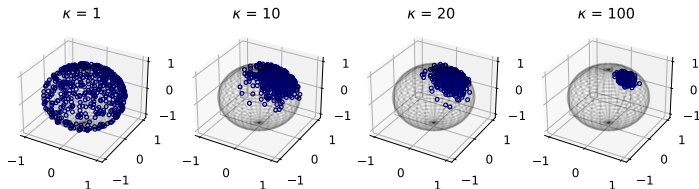
Illustration of the geometric nature of bias. Each point is the embedding of an image. In green: two male identities. In red: two female identities. The overlapping region between two identities is higher for females than for males.

↪ We choose to change the spread of each identity, according to their gender.

vMF distribution

The vMF distribution in dimension d with mean direction $\boldsymbol{\mu} \in \mathbb{S}^{d-1}$ and concentration parameter $\kappa > 0$ is a probability measure defined on the hypersphere \mathbb{S}^{d-1} by the following density:

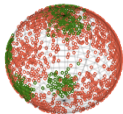
$$V_d(\mathbf{z}; \boldsymbol{\mu}, \kappa) := C_d(\kappa) e^{\kappa \boldsymbol{\mu}^\top \mathbf{z}},$$



500 samples from the vMF distribution in dimension 3.

Statistical Model on the Hypersphere

- females
- males



$$\mathbb{P}(X \in dx) = \sum_{k=1}^K \pi_k \underbrace{C_d(\kappa_k)}_{\text{hyperspherical gaussian}} \exp(\kappa_k \mu_k^T x)$$

K identities

μ_k : centroid of the k -th identity

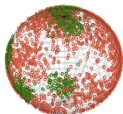
$$\kappa_k = \begin{cases} \kappa_F & \text{if female,} \\ \kappa_M & \text{if male.} \end{cases}$$

↪ We set a mixture of von Mises-Fisher distributions, as a statistical model on the hypersphere \mathbb{S}^{d-1} .

The parameter κ is the inverse of the variance of a gaussian constrained to live on \mathbb{S}^{d-1} .

Statistical Model on the Hypersphere

- females
- males



$$\mathbb{P}(X \in dx) = \sum_{k=1}^K \pi_k \underbrace{C_d(\kappa_k)}_{\text{hyperspherical gaussian}} \exp(\kappa_k \mu_k^T x)$$

K identities

μ_k : centroid of the k -th identity

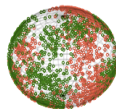
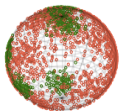
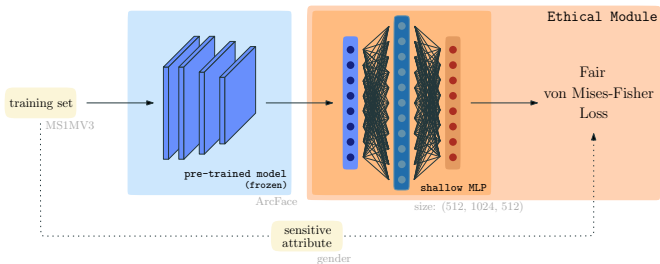
$$\kappa_k = \begin{cases} \kappa_F & \text{if female,} \\ \kappa_M & \text{if male.} \end{cases}$$

With hyperparameters κ_F and κ_M , the negative log-likelihood of the statistical model is the *Fair von Mises-Fisher loss*:

$$\mathcal{L}_{\text{FvMF}}(\Theta, \{\mu_k\}) = -\frac{1}{N} \sum_{i=1}^N \log \left[\frac{C_d(\kappa_{y_i}) e^{\kappa_{y_i} \mu_{y_i}^T \mathbf{z}_i}}{\sum_{k=1}^K C_d(\kappa_k) e^{\kappa_k \mu_k^T \mathbf{z}_i}} \right],$$

where $\mathbf{z}_i = f_{\Theta}(\mathbf{x}_i)$ is the embedding of the image \mathbf{x}_i .

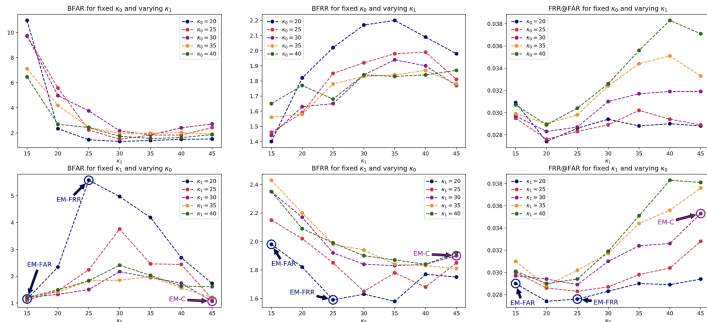
The Ethical Module



○ females
○ males

Results

BFAR and BFRR trends are correlated with κ_H and κ_F .



New SOTA for correcting the gender bias of pre-trained models
(3 methods: EM-FAR, EM-FRR, EM-C).

MODEL	10^{-4}			10^{-3}		
	FRR@FAR (%)	BFRR	BFAR	FRR@FAR (%)	BFRR	BFAR
ARCFACE	0.078	10.27	4.72	0.059	4.17	1.81
ARCFACE + PASS-G	0.315	4.54	6.51	0.107	5.22	2.11
ARCFACE + EM-FAR	0.151	11.22	2.11	0.072	9.16	1.19
ARCFACE + EM-FRR	0.100	5.89	33.65	0.058	4.11	5.24
ARCFACE + EM-C	0.164	9.18	2.44	0.081	5.15	1.20

Advantages

- Can be applied to any pre-trained model,
- Very fast training,
- Takes advantage of the performance of SOTA pre-trained networks,
- Interpretability: minimizing the Fair von Mises-Fisher loss is equivalent to maximizing the true likelihood of a Gaussian mixture model,
- The sensitive attribute (here, the gender) is only used during the training phase of the model, not afterwards.

Thanks for your attention !

For more information, please reach out to:

jean-remy.conti@telecom-paris.fr

or check out our paper