

Pragmatic Normalized Group-Fairness and Fairness-vs-Performance Trade-Off in Binary Classifiers

Riccardo Rovatti

riccardo.rovatti@unibo.it

23/01/2026

Classical approach

- Group Fairness in binary classifiers
- Let's ask an «expert»

Name		f
statistical parity	F_1	$\frac{TP+FP}{TP+TN+FP+FN}$
conditional accuracy	F_2	$\frac{TP}{TP+FN}$
false positive/negative rate	F_3	$\frac{FP}{TN+FP}$
calibration	F_4	$\frac{TP}{TP+FP}$
false discovery/omission rate	F_5	$\frac{FP}{TP+FP}$
treatment	F_6	$\frac{FP}{FN}$

An example

- ISIC challenge 2020
- Train+Val: 29985 skin lesions – malignant(P 543)/benign(N 29442)
- Test: 9912 skin lesions – malignant (P 255)/benign (N 9657)
- Resnet-50 trained on Image net + 2 fully connected layers to generate the logits
- Unbalanced binary cross-entropy
- ADAM + Reduce learning rate on plateau + batch 32
- 15 epochs – best val loss

Classical analysis

Let's make some calculation

Q = female patients

\mathcal{R} = male patients

		predicted	
		P	N
true	P	69	22
	N	1076	3224

		predicted	
		P	N
true	P	126	38
	N	1487	3870

Fairness normalization

- Hp: focus functions are scale invariant \rightarrow TP, TN, FP, FN are rates instead of counts

$$TP + TN + FP + FN = 1$$

$$0 \leq TP, TN, FP, FN \leq 1$$

Name		f
statistical parity	F_1	$\frac{TP+FP}{TP+TN+FP+FN}$
conditional accuracy	F_2	$\frac{TP}{TP+FN}$
false positive/negative rate	F_3	$\frac{FP}{TN+FP}$
calibration	F_4	$\frac{TP}{TP+FP}$
false discovery/omission rate	F_5	$\frac{FP}{TP+FP}$
treatment	F_6	$\frac{FP}{FN}$

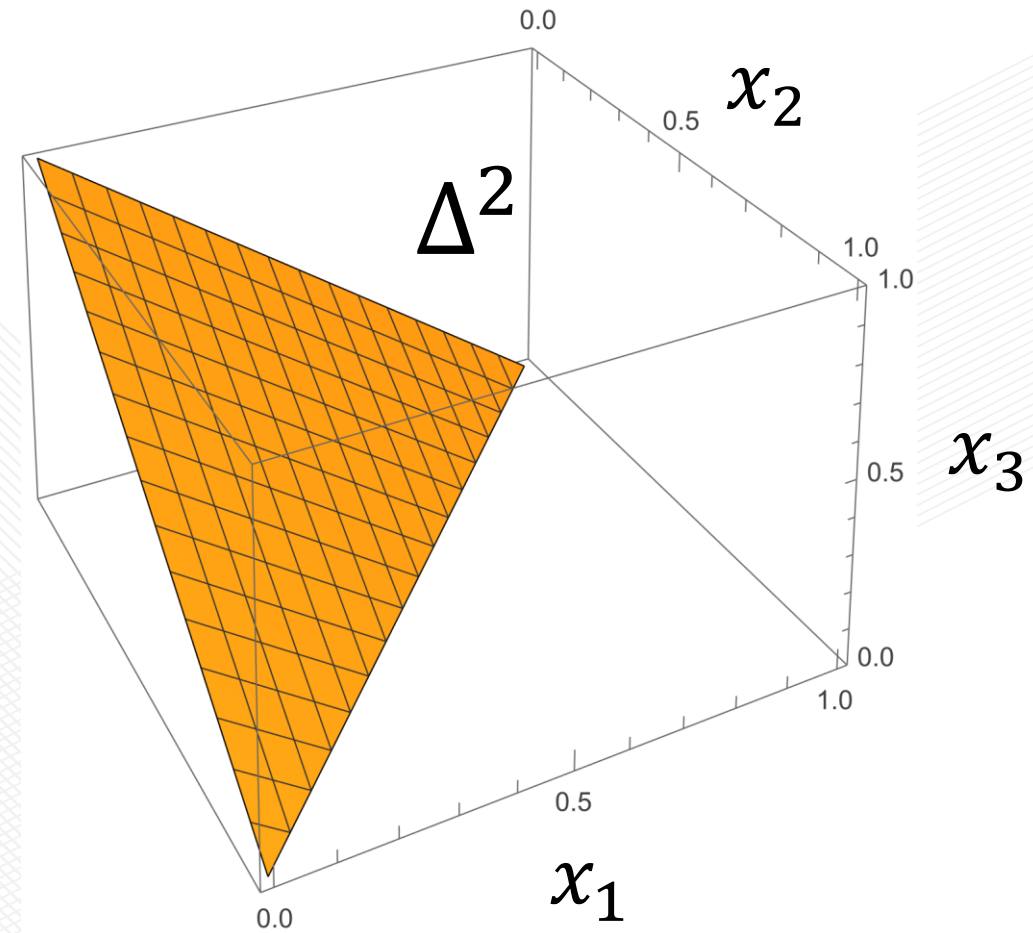
Fairness normalization

- Hp: focus functions are scale invariant \rightarrow TP, TN, FP, FN are rates instead of counts

Δ^3

$$TP + TN + FP + FN = 1$$

$$0 \leq TP, TN, FP, FN \leq 1$$



Fairness normalization

$$\begin{aligned} X &= (x_1, x_2, x_3, x_4) \\ Y &= (y_1, y_2, y_3, y_4) \end{aligned} \quad v(X, Y) = |f(X) - f(Y)|$$

$$\begin{aligned} Y &= (\text{TP}^{\mathcal{R}}, \text{TN}^{\mathcal{R}}, \text{FP}^{\mathcal{R}}, \text{FN}^{\mathcal{R}}) \\ X &= (\text{TP}^{\mathcal{Q}}, \text{TN}^{\mathcal{Q}}, \text{FP}^{\mathcal{Q}}, \text{FN}^{\mathcal{Q}}) \end{aligned} \quad v(X, Y) = \underline{v}$$

$$\phi_{\mathcal{Q}, \mathcal{R}, f} = 9 \int_{\{X \in \Delta^3, Y \in \Delta^3 \mid v(X, Y) \geq \underline{v}\}} dX dY$$

Fairness normalization

$$\phi_{\mathcal{Q}, \mathcal{R}, f} = 9 \int_{\{X \in \Delta^3, Y \in \Delta^3 \mid v(X, Y) \geq \underline{v}\}} dX dY$$

Fairness in $[0,1]$

Le possibili
coppie di
matrici di
confusione...

...che risultano
meno «fair»

Special cases

$$\phi_{Q, \mathcal{R}, f} = 9 \int_{\{X \in \Delta^3, Y \in \Delta^3 | v(X, Y) \geq \underline{v}\}} dX dY$$

- Statistical parity $\phi = \frac{1}{5} (5 - 12\underline{v} + 20\underline{v}^3 - 15\underline{v}^4 + 2\underline{v}^6)$
- Conditional accuracy & Calibration $\phi = (1 - \underline{v})^2$
- Treatment $\phi = 2 \frac{\underline{v} - \ln(1 + \underline{v})}{\underline{v}^2}$

Fairness normalization

$$f' = \alpha f + \beta$$

$$v' = |f'^{\mathcal{B}} - f'^{\mathcal{G}}| = |\alpha| |f^{\mathcal{B}} - f^{\mathcal{G}}| = |\alpha| v$$

$v' \geq \underline{v}'$ è equivalente a $v \geq \underline{v}$

Name		f
statistical parity	F_1	$\frac{TP+FP}{TP+TN+FP+FN}$
conditional accuracy	F_2	$\frac{TP}{TP+FN}$
false positive/negative rate	F_3	$\frac{FP}{TN+FP}$
calibration	F_4	$\frac{TP}{TP+FP}$
false discovery/omission rate	F_5	$\frac{FP}{TP+FP}$
treatment	F_6	$\frac{FP}{FN}$

Flipping: an abstract model for fairness improvement

- Select a subgroup in each group
- Flip the decisions in those subgroups with certain probabilities depending
 - On the group
 - On the original prediction

$$\begin{aligned}\hat{TP}^Q &= TP^Q - \nu_P^Q TP^{Q'} + \nu_N^Q FN^{Q'} \\ \hat{TN}^Q &= TN^Q - \nu_N^Q TN^{Q'} + \nu_P^Q FP^{Q'} \\ \hat{FP}^Q &= FP^Q - \nu_P^Q FP^{Q'} + \nu_N^Q TN^{Q'} \\ \hat{FN}^Q &= FN^Q - \nu_N^Q FN^{Q'} + \nu_P^Q TP^{Q'}\end{aligned}$$

$$\rightarrow \hat{\phi}_{Q, \mathcal{R}, f}$$

Flipping: an abstract model for fairness improvement

Administer flipping rates to maximize fairness

$$\begin{aligned} \max_{\nu_P^Q, \nu_N^Q, \nu_P^R, \nu_N^R} & \hat{\phi}_{Q, R, f} \\ \text{s.t.} & \hat{p}^Q \geq p_{\min}^Q \\ \text{s.t.} & \hat{p}^R \geq p_{\min}^R \end{aligned}$$

Keep performance above a minimum level

Conclusion

Normalized Fairness

- Pragmatic
- Simple
- Allows (fair!) comparison between different criteria
- Allows studying trade-off between criteria
- Allows studying trade-off between fairness and performance