

Table of Contents

Contents

1 Bias in Data and Mathematical Models	1
1.1 Basic Definitions	1
1.2 Algorithmic Fairness	1
1.3 The Confusion Matrix and Related Measures	2
1.4 Measures for Algorithmic Fairness	3
2 Origin and Types of Bias	3
3 Conclusions	5
4 Case Study	6

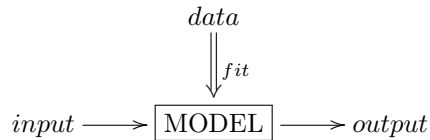
1 Bias in Data and Mathematical Models

1.1 Basic Definitions

What is a Model?

Definition 1 (Mathematical Model). A mathematical model is a description of a system using mathematical concepts and language.

Use of a model in banking is fact-based (data is used to fit the model):



What is Bias in Models?

Definition 2 (Bias in Mathematical Models). Model Bias (also Algorithmic Bias) refers to the systematic and repeatable error in a Model that creates outcomes that are statistically at odds with the reality of the population whose behaviour it is supposed to reflect.

1.2 Algorithmic Fairness

Fairness Assumptions

Assume S the protected feature with an under-privileged group S_u and a privileged group S_p , and predicted outcomes \hat{Y} .

- A. **Independence:** the probability of being in S_p or S_u has nothing to do with $\hat{Y} - \hat{Y}$ is independent of S .
- B. **Separation:** The predicted probability of having a favourable outcome has nothing to do with membership of the protected feature – \hat{Y} is independent of S , given Y .
- C. **Sufficiency:** the prediction should not depend on the protected group – Y is independent of S , given \hat{Y} .

1.3 The Confusion Matrix and Related Measures

A simple and clear way to show what the model does is a “confusion matrix”: simply show how much of the observations are classified correctly by the model.

Useful Concepts for the Confusion Matrix

The following are useful measures for how good a classification model fits its data:

- *Accuracy:* The proportion of predictions that were correctly identified.
- *Precision* (or positive predictive value): The proportion of positive cases that correct.
- *Negative predictive value:* The proportion of negative cases that were correctly identified.
- *Sensitivity* or Recall: The proportion of actual positive cases which are correctly identified.
- *Specificity:* The proportion of actual negative cases which are correctly identified.

Let us use the following definitions:

- Objective concepts (depends only on the data):
 - P : The number of positive observations ($y = 1$)
 - N : The number of negative observations ($y = 0$)
- Model dependent definitions:
 - True positive (TP) the positive observations ($y = 1$) that are by the model correctly classified as positive;
 - False positive (FP) the negative observations ($y = 0$) that are by the model incorrectly classified as positive – this is a false alarm (Type I error);
 - True negative (TN) the negative observations ($y = 0$) that are by the model correctly classified as negative;
 - False negative (FN) the positive observations ($y = 1$) that are by the model incorrectly classified as negative – miss (Type II error).

	Observed pos.	Observed neg.	
Pred. pos.	TP	FP	Pos.pred.val = $\frac{TP}{TP+FP}$
Pred. neg.	FN	TN	Neg.pred.val = $\frac{TN}{FN+TN}$
	Sensitivity	Specificity	Accuracy
	$= \frac{TP}{TP+FN}$	$= \frac{TN}{FP+TN}$	$= \frac{TP+TN}{TP+FN+FP+TN}$
	$= \frac{TP}{TP+FN}$	$= \frac{TN}{FP+TN}$	$= \frac{TP+TN}{TP+FN+FP+TN}$

Table 1: The confusion matrix, where “pred.” refers to the predictions made by the model, “pred.” stands for “predicted,” and the words “positive” and “negative” are shortened to three letters.

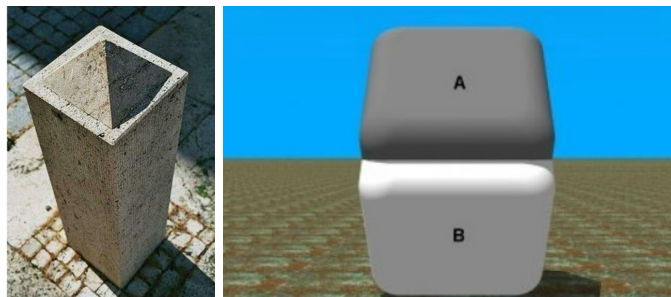
The Definition of the Confusion Matrix

1.4 Measures for Algorithmic Fairness

Fairness Metrics

- A. **Equal Opportunity:** $FNR_p = FNR_u$ (equivalent with $TPR_p = TPR_u$)
- B. **Predictive Equality:** $FPR_p = FPR_u$ (equivalent with $TNR_p = TNR_u$)
- C. **Equalised Odds:** both previous
- D. **Predictive Parity** or outcome test: $Prec_p = Prec_u$ (equal positive predictive value or precision ($Prec = \frac{TP}{TP+FP}$))
- E. **Demographic Parity:** membership of S has no correlation with favourable outcome (\hat{Y}).

Visual Biases are systematic mis-interpretations



2 Origin and Types of Bias

General Causes of Bias in Data

- **Confirmation Bias:** bias governs what we search for and we find what we search for (e.g. belief preservation via social media)
- **Selection Bias:** sample is not representative
- **Historical Bias:** socio-cultural prejudices are reflected in data (e.g. Google image search)
- **Survival Bias:** the winner takes it all and losers are forgotten (e.g. hedge fund performance statistics)
- **Availability Bias:** violence has systematically decreased over the millennia, however we might think of modern times to be more violent. (focus on through the cycle data, different perspectives)
- **Outlier Bias:** thinking about a startup one thing of Google, Amazon, Facebook, etc. These are the outliers. (use median instead of average, identify outliers, etc.)

Process Related Bias

- **Reporting bias:** selective reporting of some data: citation bias, language bias (ignore reports in other languages), duplicate publication bias (copied data found twice), location bias (some studies are hard to find), publication bias (secret or not popular data is hard to find), outcome reporting bias (e.g. company does not report with much bravado when results are bad), time lag bias
- **Automation bias:** we prefer automated systems to provide data
- **Selection bias:** data is not representative (e.g. sampling bias, convergence bias (only people who got a loan are in our database), participation bias, etc.)
- **Overgeneralisation bias:** no black swans in the data
- **Group Attribution bias:** generalisation of stereotypes, in-group bias (preference for members in the group), out-group bias (stereotype for other groups), etc.
- **Implicit bias:** search for conforming information

What causes Bias in Models?

Bias can appear due to

A. biases in the data

- less representative group
- the data is not a good representation of the reality (e.g. survivor bias, previous societal bias or limitations, etc.)

- B. the data processing pipeline (extraction, cleaning, transformation, binning)
- C. the model development process (including selection of algorithm, variables, etc.)
- D. the particulars of model implementation

3 Conclusions

Fairnes is a matter of perspective

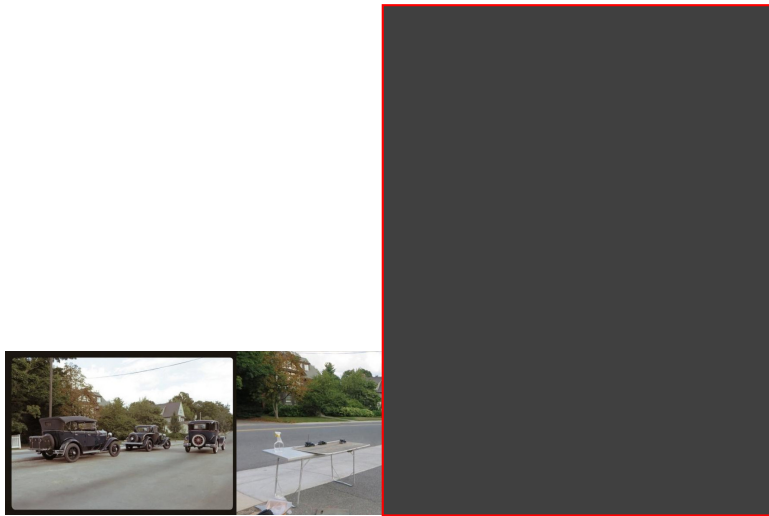


Figure 1: Fairness is a matter of perspective

Bias is learnt



4 Case Study

Case study on addressing bias: car insurance



References

- [1] Sray Agarwal and Shashin Mishra. *Responsible AI: Implementing Ethical and Unbiased Algorithms*. Springer, 2021.
- [2] Mark Coeckelbergh. *AI ethics*. Mit Press, 2020.
- [3] Philippe J.S. De Brouwer. *The Big R-Book: From Data Science to Learning Machines and Big Data*. New York: John Wiley & Sons, Ltd, 2020. ISBN: 978-1-119-63272-6.
- [4] Markus Dirk Dubber, Frank Pasquale, and Sunit Das. *The Oxford handbook of ethics of AI*. Oxford Handbooks, 2020.
- [5] Moritz Hardt, Solon Barocas, and Arvind Narayanan. 'Fairness in Machine Learning: Limitations and Opportunities'. In: *Solon Barocas Moritz Hardt and Arvind Narayanan* (2018).