

## Fairness

Machine Learning Model evaluation is not only just calculating loss metrics → Before it goes to production, it's critical to audit training data & evaluate predictions for bias.

### Prediction Bias

→ This value indicates how far apart the average of predictions is from the labels in the dataset

### ML-Bias

→ studying, perceiving or favoring towards some things, people or groups over others. These biases can affect collection & interpretation of data

- ① automation bias
- ② confirmation bias.
- ③ Experiment bias
- ④ group attribution bias
- ⑤ Implicit bias
- ⑥ In-group bias

⑦ Out-group homogeneity bias.

⑧ - Systematic error introduced by a sampling or reporting procedure.

① Coverage bias, ② Non-response bias, ③ participation bias

④ reporting bias, ⑤ Sampling bias, ⑥ Selection bias.

Machine learning models are not inherently objective. ML practitioners train models by feeding them a dataset of training examples, and human involvement in the provision and curation of this data can make a model's prediction to biased.

① Reporting bias →

occurs when the frequency of events, properties or outcomes captured in a dataset does not accurately reflect the real world frequency. This bias can occur because people tend to focus on documenting things that are unusual or especially momentous. These ordinary are not recorded.

Eg:

Sentiment Analysis model is trained to predict whether book reviews are positive or negative based on ~~contents~~

⇒ Majority of reviews in training dataset reflect extreme opinions (either they love or hate it)

## ② Historical Bias

Eg: A city housing dataset from the 1960's contains home price data that reflects discriminatory lending practices in effect during that decade.

## ③ Automation bias

is tendency to favour results generated by automated systems over those generated by non-automated systems. Perspective of user each.

Eg: Spockit manufacturer

## ④ Selecting Bias

if dataset doesn't reflect real-world distribution  
like coverage bias, non-response bias, sampling bias.

→ coverage bias

if data is not selected in a representative fashion.

Eg: model trained to predict future sales of a new product based on phone surveys conducted with a sample of customers who bought the product and left the people/customers who opt to buy competing product were not surveyed.

→ Non response bias

also known participation bias occurs if data ends up being unrepresentative due to participation gaps in data-collection process.

→ Sampling bias

⇒ occurs if proper randomization is not used during data collection.

★ Group attribution bias

is a tendency to generalize what is true of individuals to the entire group to which they belong.

→ In-group bias

is a preference for members of your own group you also belong or for characteristics that you also share.

⇒ Outgroup homogeneity bias is a tendency to stereotype individual members of a group to which you do not belong to their characteristics as more uniform.

★ Influent bias

occurs when assumptions are made on one's own model of things and personnel experience that don't necessarily apply more generally.

Eg: head shake gesture → yes and some form is no

## \* Confirmation bias

biases consciously favors data in ways

that affirm preexisting beliefs & hypothesis

Eg: Aggregation of data.

## \* Experimental bias

builder keeps adding a model until it produces a result that aligns with their original hypothesis

## Identifying Bias

Keep issues of fairness in mind & audit for potential sources of bias.

### Missing feature values

If dataset is missing one or more feature values for large number of examples.

Indicator that certain key characteristics of your dataset are under-represented.

## ⇒ Unrepresented feature bias

It introduces bias when in dataset you see unprecedented feature values that stand out as especially unusual/unusualistics or unusual.

## ⇒ Data Skew

Example dog - adoptability, it's not sufficient to look overall summary.

### Mitigating Bias

Once the source of bias is identified in training data.

There are two strategies

① Augmenting the training data.

② Augmenting the models loss function

### ① Augmenting the training data

The most straightforward way to address the problem is often to collect additional data.

### ② Augmenting the Model's optimization function

“log loss” to penalize incorrect model predictions.

Log loss doesn't take subsequent membership in consideration

Instead we use optimization functions designed to penalize error in fairness-aware fashion.

Tensorflow Model Remediation library provides utilities for applying two different bias-mitigation techniques during Model training

- ⇒ MinDiff → it aims to balance the error for two different slices of data. by adding a penalty for differences in the prediction distributions for the two groups
- ⇒ Counterfactual Logit Pairing aims to ensure that changing a sensitive attribute of a given example doesn't alter the model's prediction for that example.

### Evaluating for Bias

#### Aggregate performance Metrics

- (~~1 precision~~) <sup>precision</sup> ~~precision~~ <sup>precision</sup>
  - ① recall
  - ② accuracy.
- will not catch issues

#### Demographic parity

##### Majority group & Minority group

The best way to check the demographic parity fairness or bias benefit

of demographic parity both majority and minority are selected in same proportion

### Equality of opportunity

### Counterfactual Fairness

when we don't have data for demographic parity & Equality of opportunity group. Here there would be candidates with marked as eligible or not eligible. Those candidates are picked by model then we can say model is Counterfactual Fairness.

#### Benefits

→ It is used to evaluate the predictions of Model

Demographic parity & Equality of opportunity across the groups in aggregate

Majority → Demographic life plan.

Minority → Demographic life plan where  
Men are less considered.