# DATA QUALITY DIMENSIONS FOR FAIR AI

**Camilla Quaresmini**[*]
Department of Electronics, Information and Bioengineering
Politecnico di Milano
Piazza Leonardo da Vinci 32, MI 20133, Italy
camilla.quaresmini@polimi.it

**Giuseppe Primiero**[†]
Department of Philosophy
University of Milan
Via Festa del Perdono 7, MI 20122, Italy
giuseppe.primiero@unimi.it

May 12, 2023

## ABSTRACT

AI systems are not intrinsically neutral and biases trickle in any type of technological tool. In particular when dealing with people, AI algorithms reflect technical errors originating with mislabeled data. As they feed wrong and discriminatory classifications, perpetuating structural racism and marginalization, these systems are not systematically guarded against bias. In this article we consider the problem of bias in AI systems from the point of view of Information Quality dimensions. We illustrate potential improvements of a bias mitigation tool in gender classification errors, referring to two typically difficult contexts: the classification of non-binary individuals and the classification of transgender individuals. The identification of data quality dimensions to implement in bias mitigation tool may help achieve more fairness. Hence, we propose to consider this issue in terms of completeness, consistency, timeliness and reliability, and offer some theoretical results.

*Keywords* Information Quality, Fairness, Mislabeling, Gender bias, Timeliness

## 1 Introduction

Machine Learning models trained on huge amounts of data are intrinsically biased when dealing with people. Common face recognition systems used in surveillance tasks generate false positives labeling innocent people as suspects. Social credit systems link individuals to the state of their social credit, making decisions based on that score. Systems of this kind also confuse biological sex with gender, feeding wrong and discriminatory classifications. In all of those cases, subjects suffer a credibility deficit due to prejudices related to their social identity [Fricker, 2007]: a dark-skinned man could be characterized by a higher risk of recidivism after being arrested; a short-haired skinny young woman – or a long-haired boy with feminine traits – might be the target of transphobic attacks following misgendering. Through the deployment of these technologies, society makes the gap separating rich from poor, cisnormative from non-cisnormative individuals, more constitutive as automatized and standardized.

Already before the explosion of ML algorithms, Friedman and Nissenbaum [1996] offered a framework for understanding three categories of bias in computer systems, assuming the absence of bias as necessary to define their quality. Later on, the emergence of contemporary, data-driven AI systems based on learning has significantly worsened the situation. Buolamwini and Gebru [2018] show cases of algorithmic discrimination based on gender and race. In particular, analysing three commercial gender classification systems, they found gender reductionism and accuracy disparities in the classification of darker/lighter and females/males individuals, showing that dark-skinned females are the most misclassified group. Mehrabi et al. [2019] list different types of bias in AI applications, based on the identification of two underlying potential sources of unfairness in ML systems: data and algorithms. Miceli et al. [2021] investigate the power dynamics and the political nature of data in AI, based on the assumption that data is always biased, being the product of subjective and unbalanced social dynamics. According to a recent MIT study [Northcutt et al., 2021a] the

---

[*]META Social Sciences and Humanities for Science and Technology, Politecnico di Milano
[†]Logic, Uncertainty, Computation and Information Group, University of Milan

most well-known AI datasets are full of labeling errors, something which distorts our perception of progress related to the field.

On this basis, the development and deployment of fairer AI system has been increasingly demanded. Fairness itself is a social construct which depends on the contexts of application [Bellamy et al., 2018], and it is also difficult to choose a definition among those provided in the literature [Mehrabi et al., 2019, Verma and Rubin, 2018], many of which are known to be mutually incompatible [Friedler et al., 2016, Saravanakumar, 2020, Berk et al., 2017, Kleinberg et al., 2016, Chouldechova, 2016]. Such request appear especially relevant in certain application contexts. For example, as examined in Hanna et al. [2021], face is commonly used as a legitimate site of gender classification, and this is operazionalized and automatized in technologies such Automatic Gender Recognition (AGR). AGR is a subfield of Face Recognition (FR), which algorithmically derives gender from faces' physical traits to perform individuals classification [Keyes, 2018, Hamidi et al., 2018]. This technique relies on the assumption that gender identity can be computationally derived from facial traits. Scheuerman et al. [2019] show that the most famous AGR systems are not able to classify non-binary genders, also performing poorly on transgender individuals. This is due to the fact that AGR incapsulates a binary, cisnormative conception of gender, adopting a male/female scheme which invalidates non-binary identities.

An important task, common to technology and philosophy, is therefore the identification of criteria that may help developing fairness conditions for AI systems. This is especially true for systems of machine vision that implement FR algorithms. While a number of mitigation techniques are available to mitigate bias, their primary focus on accuracy is clearly not sufficient to mitigate social discrimination. It remains to be explored which data dimensions (and related biases) do they address, and whether these are sufficient in all cases. In the following of this paper, we offer a contribution in this direction, arguing that the task of bias mitigation tools can be supported by reasoning on quality dimensions that so far have been left ignored. In particular, we offer examples to show how dimensions of consistency, completeness, timeliness and reliability can be used to establish fairer AI systems. This research is in line with the quest for integrating useful empirical metrics on fairness in AI with asking key (conceptual) questions, see Scantamburlo [2021].

The paper is structured as follows. In Section 2 we offer an overview of fainess definitions and bias types relevant for the work. In Section 3 we briefly overiview the technical details of Cleanlab, one of the available bias mitigation tools for AI algorithms. In Section 4 we introduce data dimensions as the criteria to be used for assessing data quality and we argue for reconsidering the relevance of such dimensions in the task of reducing biases in AI algorithms. In Section 5 we present two examples to justify the requirement that certain data quality dimensions be explicitly accounted for in classification tasks and reassess the Cleanlab strategy with this revision in mind. In Section 6 we propose a definition of fair AI that includes such dimensions and formulate some theoretical results on its basis. In Section 7 we conclude illustrating future research lines.

## 2   Fairness and Bias

Fairness is a broad concept, typically modeled for classification tasks in terms of individual and group fairness.

*Individual fairness* is concerned with treating similar individuals in the same way, based on a metric that defines how similar two individuals are in a given context. Formal solutions defining such distance between individuals may include (*fairness through awareness* [Dwork et al., 2011]) or ignore (*fairness through unawareness* [Grgic-Hlaca et al., 2016]) protected attributes during the decision-making process. *Protected attributes* are characteristics for which an individual cannot be discriminated against, by law. Those attributes are potentially sensitive features which can be tied to social bias (e.g. ethnicity, sex, gender, socioeconomic status, religion, ability). Groups of population according to these attributes should have parity in terms of benefits received [Bellamy et al., 2018]. Nonetheless, protected attributes are not universal: some groups may have had historically systematic advantages, i.e. cases of *privileged values* of a protected attribute. This leads to differentiating privileged and unprivileged groups. A *privileged group* is systematically put at an advantage with respect to the beneficial outcome. On the contrary, an *unprivileged group* is systematically put at a disadvantage with respect to the beneficial outcome [Aasheim et al., 2020]. If a decision is the same both in the actual world where the individual belongs to the unprivileged group and in a counterfactual world where the individual belongs to the privileged one, then *counterfactual fairness* is satisfied [Kusner et al., 2018].

*Group fairness* consists in groups defined by protected attributes which receive similar treatments. If individuals from the protected and unprotected groups have the same probability of being assigned the *favorable label*, i.e. the label whose value is considered the denote the favorable outcome, this amounts to *statistical parity* and the value is zero. A negative value indicates a disadvantage for the unprivileged group [Aasheim et al., 2020]. The proportion of favourable labels of the members of a privileged group compared to favorable labels of the members of an unprivileged group is often used to define fairness [Friedman and Nissenbaum, 1996, Aasheim et al., 2020]. Several solutions are available, starting from setting the probability of being assigned a positive outcome equal both for individuals belonging to

Table 1: Data bias.

| Bias type | Definition | Literature |
|---|---|---|
| **Behavioral bias** | User's behavior can be different across contexts | Olteanu et al. [2019] |
| **Exclusion bias** | Systematic exclusion of some data | Fabbrizzi et al. [2021] |
| **Historical bias** | Cultural prejudices are included into systematic processes | Suresh and Guttag [2021] |
| **Time interval bias** | Data collection in a too limited time range | CertNexus [2021] |

Table 2: Label bias.

| Bias type | Definition | Literature |
|---|---|---|
| **Chronological bias** | Distortion due to temporal changes in the world which data are supposed to represent | Fabbrizzi et al. [2021] |
| **Historical bias** | Cultural prejudices are included into systematic processes | Suresh and Guttag [2021] |
| **Misclassification bias** | Data points are assigned to incorrect categories | Centre for Evidence-Based [2022] |

protected and unprotected groups (*equality of opportunity* [Hardt et al., 2016]). If the likelihood of a positive outcome is the same regardless of whether the person is in the protected group, then the predictor satisfies *demographic parity* [Dwork et al., 2011], and if this probability also considers a set of legitimate factors, then *conditional statistical parity* [Corbett-Davies and Goel, 2018] is satisfied. Furthermore, the *equalized odds* [Mehrabi et al., 2019] principle sets the probability of a person in the positive class being correctly assigned a positive outcome equal to the probability of a person in a negative class being incorrectly assigned a positive outcome, both for the protected and unprotected groups.

Despite a unique definition missing in the literature, fairness is often intended to correspond to the avoidance of bias. This can be formulated at two distinct levels: first, identifying and correcting problems in datasets, as a model trained with a mislabeled dataset will provide biased outputs; second, correcting the algorithms, as even in the design of algorithms biases can emerge [Hooker, 2021]. In the present Section we are interested in considering the former case, both for datasets and for their labels. A variety of approaches and contributions is available in the literature focusing on identifying bias in datasets and labels. Here we list the types of bias which are relevant to the present work, see respectively Tables 1 and 2. Albeit not exhaustive, these lists of biases represent a good starting point to investigate quality dimensions required to address them.

*Behavioral bias* cause users to be included or excluded from a dataset representing the population of interest [Olteanu et al., 2019]. This leads to the systematic exclusion of some data from the population, resulting in *exclusion bias*. *Historical bias* can arise both during data collection and labelling. Concerning the dataset construction, it may occur during the selection of training data, leading to exclusion bias on the basis of personal beliefs and cognitive preconceptions.

When data collection is achieved in a too limited time range, the *time interval bias* could impact the outcomes. We will argue below that this type of bias occurs in datasets about people.

Bias may also affect the label set. We talk about *label quality bias* when errors hit the quality of labels.

In particular, *chronological bias* arises when data are no longer representative of the world they must describe, a case we will stress as extremely significant in our analysis below.

*Historical bias* consists of systematic cultural prejudices. When datapoints are mislabelled, being assigned to incorrect categories, *misclassification bias* occurs. This type of bias can be caused by omitted variables, which may happen during the selection of data relevant features.

A crucial task is therefore the development of conceptual strategies and technical tools to mitigate bias emergence in both data and label sets.

## 3   Mitigating Bias

A *bias mitigation algorithm* is a procedure for reducing unwanted bias in training datasets or models, with the aim to improve the fairness metrics. Those algorithms can be classified into three categories [D'Alessandro et al., 2017]:

- pre-processing, when the training data is modified;
- in-processing, when the learning algorithm is modified;
- post-processing, when the predictions are modified.

Several tools are available to audit and mitigate biases in datasets, thereby attempting to implement diversity and to reach fairness. Among the most common are AIF360[3], Aequitas[4] and Cleanlab.[5] While AIF360 and Aequitas focus on examining and reporting bias and discrimination, referring to a precise idea of fairness, Cleanlab's priority is data cleaning. In this article, we refer to the latter one as a testbed for the proposed extension of data quality dimensions suggested below in Section 5. For an overview of the symbols used from now on, see Table 3

Table 3: Symbols used in the present work.

| | |
|---|---|
| $t_n$ | Time index |
| $\mathcal{T} := \{t_1, \ldots, t_n\}$ | Time frame |
| $d$ | Generic datapoint |
| $i, j, l$ | Data Labels |
| $y^*$ | Discrete random variable correctly labelled |
| $\tilde{y}$ | Discrete random variable wrongly labelled |
| $y* \to \tilde{y}$ | A mapping between variables |
| $p_{\mathcal{T}}[(\tilde{y} = i)_{t_n} \mid (y* = j)_{t_{n-m}}]$ | The probability of label $i$ being wrong at time $t_n$, given that label $j$ was correct at time $t_{n-m}$ |
| $C_{\tilde{y},y^*}[i, j, \mathcal{T}]$ | Temporal confident joint, where the correct label can change from $i$ to $j$ in time frame $\mathcal{T}$ |
| $C_{\tilde{y},y*}[i, \mathcal{T}]$ | Temporal confident joint, where the correctness of the same fixed label $i$ can change in time frame $\mathcal{T}$ |
| $\varepsilon$ | Change rate |
| $\hat{p}'(\tilde{y}; x_i; \theta)$ | Predicted probability of label $\tilde{y}$ for variable $x_i$ and model parameters $\theta$ |
| $L$ | Label set |
| $X$ | AI system |
| $L_{t_1} := \{l_1, \ldots, l_n\}$ | Partition of the label set |
| $P$ | Population of interest |
| $p$ | An element from $P$ |
| $d(X)_{\mathcal{T}}$ | A datapoint in system $X$ over time frame $\mathcal{T}$ |
| $y^*(d)$ | A correct label for the datapoint $d$ |
| $\pi$ | Threshold variable |

### 3.1   A Mitigation Tool

Cleanlab is a deep learning framework to find label errors in datasets.[6] It uses Confident Learning (CL), an approach which focuses on label quality with the aim to address uncertainty in dataset labels using three principles:

1. counting examples that are likely to belong to another class using the *confident joint* and probabilistic thresholds to find label errors and to estimate noise;

---

[3]Code available at `https://github.com/Trusted-AI/AIF360`. See Aasheim et al. [2020] for a comparative study on bias mitigation with AIF360.

[4]Code available at `https://github.com/dssg/aequitas`.

[5]Code available at `https://github.com/cleanlab/cleanlab`.

[6]`https://l7.curtisnorthcutt.com/cleanlab-python-package`.

2. pruning noisy data; and

3. ranking examples to train with confidence on clean data.[7]

These three approaches are combined by an initial assumption of a class-conditional noise process, to directly estimate the joint distribution between noisy given labels and uncorrupted unknown ones. For every class, the algorithm learns the probability of it being mislabeled as any other class. This assumption may have exceptions but it is considered reasonable. For example, a "cat" is more likely to be mislabeled as "tiger" than as "airplane". This assumption is provided by the classification noise process (CNP, Angluin and Laird [1987]), which leads to the conclusion that the label noise only depends on the latent true class, not on the data.

CL [Northcutt et al., 2021b] exactly finds label errors in datasets by estimating the joint distribution of noisy and true labels. The idea is that when the predicted probability of an example is greater than a threshold per class, we confidently consider that example as actually belonging to the class of that threshold, where the thresholds for each class are the average predicted probability of examples in that class. Given

1. $\tilde{y}$ = Discrete random variable $\tilde{y} \in [m]$ takes an observed, noisy label (potentially flipped to an incorrect class);

2. $y*$ = Discrete random variable $y* \in [m]$ takes the unknown (latent), true, uncorrupted label (latent true label);

3. $[m]$ = The set of unique class labels.

CL assumes that for every example it exists a correct label $y*$ and defines a class-conditional noise process mapping $y* \rightarrow \tilde{y}$, such that every label in class $j \in [m]$ may be independently mislabeled as class $i \in [m]$, with probability $p(\tilde{y} = i \mid y* = j)$. So, maps are associations of data to wrong labels. Then CL estimates $p(\tilde{y} \mid y*)$ and $p(y*)$ jointly, evaluating the joint distribution of label noise $p(\tilde{y}, y*)$ between noisy given labels and uncorrupted unknown labels. CL aims to estimate every $p(\tilde{y}, y*)$ as a matrix $Q_{\tilde{y}, y*}$ to find all mislabeled examples $x$ in dataset $X$, where $y* \neq \tilde{y}$. Given as inputs the out-of-sample predicted probabilities $\hat{P}_{k,i}$ and the vector of noisy labels $\tilde{y}_k$, the procedure is divided into three steps:

1. estimation of $\hat{Q}_{\tilde{y}, y*}$ to characterize class-conditional label noise,

2. filtering of noisy examples,

3. training with the errors found.

To estimate $\hat{Q}_{\tilde{y}, y*}$ i.e. the joint distribution of noisy labels $\tilde{y}$ and true labels $y*$, CL counts examples that may belong to another class using a statistical data structure named *confident joint* $C_{\tilde{y}, y*}$. The formal definition of *confident joint* is as follows

$$C_{\tilde{y}, y*}[i][j] := \mid \hat{X}_{\tilde{y}=i, y*=j} \mid \tag{1}$$

where

$$\hat{X}_{\tilde{y}=i, y*=j} := \{ x \in X_{\tilde{y}=i} : \hat{p}(\tilde{y} = j; x, \theta) \geq t_j, j = \underset{l \in [m]: \hat{p}(\tilde{y}=l; x, \theta) \geq t_l}{\arg\max} \hat{p}(\tilde{y} = l; x, \theta) \} \tag{2}$$

In other words, the *confident joint* estimates the set $X_{\tilde{y}=i, y*=j}$ of examples with noisy label *i* which actually have true label *j* by making a partition of the dataset $X$ into bins $\hat{X}_{\tilde{y}=i, y*=j}$, namely the set of examples labeled $\tilde{y} = i$ with *large enough* expected probability $\hat{p}(\tilde{y} = j; x, \theta)$ to belong to class $y* = j$, determined by a per-class threshold $t_j$, where $\theta$ is the model. The threshold $t_j$ is the average expected self-confidence for each class, formally defined as

$$t_j = \frac{1}{\mid X_{\tilde{y}=j} \mid} \sum_{x \in X_{\tilde{y}=j}} \hat{p}(\tilde{y} = j; x, \theta) \tag{3}$$

When $X_{\tilde{y}=i, y*=j}$ is equal to $\hat{X}_{\tilde{y}=i, y*=j}$, $C_{\tilde{y}, y*}$ accurately finds label errors. The confident joint is the unnormalized estimate of the joint distribution of noisy and true labels. Given the confident joint $C_{\tilde{y}, y*}$, $Q_{\tilde{y}, y*}$ is estimated as

$$\hat{Q}_{\tilde{y}=i, y*=j} = \frac{\frac{C_{\tilde{y}=i, y*=j}}{\sum_{j \in [m]} C_{\tilde{y}=i, y*=j}} \cdot \mid X_{\tilde{y}=i} \mid}{\sum_{i \in [m], j \in [m]} \left( \frac{C_{\tilde{y}=i, y*=j}}{\sum_{j' \in [m]} C_{\tilde{y}=i, y*=j'}} \cdot \mid X_{\tilde{y}=i} \mid \right)} \tag{4}$$

---

[7]For details on Rank Pruning algorithm see Northcutt et al. [2017] and `https://github.com/cgnorthcutt/rankpruning`.

The confident joint algorithm hence takes two inputs. First, $\hat{P}$, an $n \times m$ matrix of out-of-sample predicted probabilities $\hat{P}[i][j] := \hat{p}(\tilde{y} = j; x_i, \theta)$. Second, the associated array of noisy labels.

These tools, extremely useful in estimating label error probabilities, have some limitations and it is easy to formulate examples for which their strategy seems unsound. A first problem arises from the Cleanlab initial assumption of the categoricity of data. Gender labelling of facial images is typically binary ("male" and "female", see e.g. in the FERET database[8]), and for each datapoint a classification algorithm calculates the projected probability that an image is assigned to the respective label. Consider though two very noisy cases:

1. images of individuals which might identify as non-binary;
2. images of transgender individuals.

In the former case, the label set becomes incomplete with respect to the dataset; in the second case, the dataset is inconsistent with respect to the label set. Hence, we have an assumption where we claim that for each data there can be datapoints that have either 1) none of the available labels as the correct one, or 2) at different times they can be under both labels. By definition, if we have disjoint labels there can be high accuracy but only on those datapoints which identify themselves in the disjointed categories.

In situations like these, it appears that the dimension of accuracy alone does no longer satisfy the correctness of the classification algorithm. In terms of dimensions, the possibility of an uncategorical datapoint or that of a moving datapoint is no longer only an accuracy problem. Hence, the identification of other data quality dimensions to be implemented in tools for bias mitigation may help achieve more fairness. In the next section we provide an overview of such dimensions to suggest an improvement.

## 4 Data Quality Dimensions

Quality dimensions are intended in two distinct ways: the extension of data, namely their values; and their intension, i.e. the data schemata, see [Batini and Scannapieco, 2006, p.19]. Each dimension is defined informally and qualitatively. Metrics can be then associated as indicators of the dimension's quality. We first provide an overview of the methodological strategies used to define extensions and intensions and provide some definitions of dimensions; then, we offer a detailed analysis of those definitions. Also in this case, we do not aim at an exhaustive review of this extensive research field, rather we indicate potential candidates to implement in the context of bias mitigation strategies and tools.

The literature on data quality dimensions offers different classifications, and several strategies can be adopted to propose sets of dimension definitions, [Batini and Scannapieco, 2006, p.36]. These include: theoretical approaches, focussing on how data might become deficient during production [Wang and Strong, 1996, p.7], and define dimensions through a formal model to include accuracy, reliability, timeliness, completeness; empirical approaches using questionnaires and interviews to build a size list, [Wang and Strong, 1996]; intuitive approaches defining dimensions according to common sense and practical experience, [Batini and Scannapieco, 2006, p.36].

There is no single and objective vision of data quality dimensions, nor a universal definition for each dimension. This is because often dimensions escape or exceed a formal definition. The cause of the large amount of dimensions also lies in the fact that data aim to represent all spatial and temporal and social phenomena of the real world. Dimensions can also be independent from the domain, so of general possible application, or they can be domain dependent, referring to phenomena characteristic of specific domains. Data quality dimensions are also constantly evolving in response to continuous development of new data-driven technologies. In the following and for the purposes of our analysis, we focus on a small number of such dimensions, which we illustrate first informally.

*Accuracy* is defined as the closeness between a value $v$ and a value $v'$, where the latter is the correct representation of the real-life phenomenon that $v$ aims to represent [Batini and Scannapieco, 2006, p.20]. Given a proper name, for example Abraham, the value $v' = Abraham$ is correct, while $v = Abraam$ is wrong. Accuracy can be syntactic or semantic. Syntactic accuracy is the closeness between the value $v$ and the elements of the corresponding domain $D$. In this case, one is interested in understanding if $v$ is among the values of $D$ and there is no need to compare $v$ with a real value $v'$. The $v$ value can be any value in the domain $D$. Instead, semantic accuracy measures the closeness of the value $v$ to a real value $v'$ comparing the two values. In this case, accuracy corresponds with the concept of correctness. Semantic accuracy seems to be more difficult to calculate than its syntactic counterpart. To check semantic accuracy we can search the same data in different data sources, and then compare those data to find the correct ones.

*Completeness* is the level at which data have the sufficient breadth, depth, and scope for their task, see Wang and Strong [1996]. In Pipino et al. [2002] three types of completeness are identified. Schema completeness is the extent to which

---

[8] https://www.nist.gov/itl/products-and-services/color-feret-database.

concepts and related properties are not lacking in the schema. Column completeness is the measure of the missing values for a specified property. Population completeness compares the lacking values with the reference population. Intuitively, the degree of completeness of a table depends on how well the table represents the corresponding real world. Sometimes it is possible to find null values, that are missing values. So the value is not available for some reason, but it exists in the real world. This value is missing because it may be unknown, or it may not exist at all, or it may exist but it is not known whether or not it actually exists [Batini and Scannapieco, 2006, p.24].

Data are characterized by evolution over the time, we have to consider the *temporal* dimension. As completeness dimension is traditionally static, the notion of completability is needed to consider also the temporal dynamics of completeness [Batini and Scannapieco, 2006, p.27]. The completability dimension provides information on how quickly the completeness degree grows over time. Based on the time dimension, data are classified as follows: stable data, which are unlikely to change; data that change in the long term, that have a very low frequency of modification but which can still change (for example addresses); data that change frequently, even continuously such as real-time information on traffic. There are three time-related dimensions about these three different types of data [Batini and Scannapieco, 2006, p.29]:

1. Currency concerns how quickly data is updated, measuring with respect to the metadata of the last update.

2. Volatility is the frequency with which the data is updated over time. A metric to measure the volatility is the length of time during which data are valid.

3. Timeliness expresses how current the data are for the activity to be performed.

Batini and Scannapieco [2006] emphasize that timeliness implies that data are not only current, but also in time for events that correspond to their use. In fact, it is possible to have current data which, however, are useless because they are late for the specific task.

*Consistency* is the coherence dimension. It amounts to checking whether or not the semantic rules defined on a set of data elements have been respected.

The variety and diversity of quality dimensions for data has suggested that a category of higher-order may be identified to establish and collect all properties that make data of good quality. This has been often formulated as *fit-for-purposeness* or *fitness-for-use*, as a general definition of Information Quality (IQ), see Wang and Strong [1996], Kahn et al. [2002] and [Batini and Scannapieco, 2006, p.221]. However, this concept appears sometimes too broad, making it hard to consider in absence of more specific dimensions.

[Illari and Floridi, 2014, p.284–285] call this the "relevance problem", because all dimensions are infected by the detection of information with respect to the intended purposes. Ideally, high quality information is suitable both for the purpose for which it was processed, and easily reusable for new purposes. In the first case we talk about purpose-depth, in the second is purpose-scope, see [Illari and Floridi, 2014, p.312]. There are multiple senses in which we speak of information or data being fit for purpose. Far from constituting a relativistic stance, it is an explicit choice of the purpose by which to examine the value of some information. While the purpose is obviously chosen by the user, the evaluation of the metrics of each dimension is an objective matter. In this sense, a dataset can be of high quality for one application, but at the same time of low quality for another [Batini and Scannapieco, 2006, p.221]. Hence, the common definition considers the quality of the data as the suitability for present and future uses, what has been called "multiple purposes response" [Illari and Floridi, 2014, p.285].

The separation between purpose-depth and purpose-scope can be formulated in terms of a bi-categorical approach. First, a distinction must be made between the purpose/s for which some information is produced ($P$-purpose), and the potentially unlimited purpose/s for which the same information may be consumed ($C$-purpose), see [Illari and Floridi, 2014, p.313]. Accordingly, in terms of a Cartesian space, we have the $P$-purpose on the $x$-axis and the $C$-purpose on the $y$-axis. For any information $I$, one must have two values in order to be represented in that Cartesian space. Second, we can compare the quality of some information with respect to $P$-purpose /$C$-purpose, identifying in this way potential discrepancies. The result allows to link IQ to a specific purpose, instead of talking about IQ as fit-for-purpose in absolute terms. Considering other potential future purposes is possible only in the case of compatibility between the new purpose and the original one, where compatible means that there is a link between the two purposes, or the nature/context of data is similar. However, this seems to preclude an update of the data, which necessarily affects the accuracy of the same. In fact, sometimes the better information fits its original purpose, the less likely it is re-purposable [Illari and Floridi, 2014, p.313].

To consider degradation of the purpose of data, the difference between $P$-purpose and $C$-purpose must be assumed and this difference introduces the consideration of time, and in particular how data might become outdated. While prior works have mainly focused on assessing fairness in static contexts, there are also attempts in the literature to consider

the dynamic nature of fairness, e.g. D'Amour et al. [2020] evaluates the environment's state over time, considering agents' decisions that cause the underlying population of individuals to evolve.

In the next Section we will argue that the timeliness dimension can be taken as basis for other categories of data quality.

# 5 Extending Data Dimensions for Fair AI

The aim of this section is to suggest improvements on errors identification in the classification of datapoints, using the gender attribute as an illustrative case. We suggest the extension of classification with dimensions of completeness, consistency and timeliness and then return to Cleanlab to illustrate how this extension could be practically implemented.

## 5.1 Incomplete Label Set

Consider the first example of a datapoint which represents a non-binary individual. Being out of gender binarism, the individual identity is not correctly recognized. The conceptual solution would be to simply assume the label set as incomplete. This means that the bias origin is in the pre-processing phase, and a possible strategy is to extend the partition of the labels adding categories as appropriate, e.g. "non binary" or "queer". The problem is here reduced to the consideration of the completeness of the label set.

Scheuerman et al. [2019] can be considered a first step in this direction. It presents a visual database, containing 2450 pictures of people organized in a folksonomy through 7 different genders. Images are scraped from Instagram accounts, and gender labels are matched to hashtags chosen by individuals themselves: #man, #woman, #nonbinary, #genderqueer, #transman, #transwoman, and #agender. The True Positive Rate (TPR) for each hashtag on each system is calculated, denoting the system accuracy. Even if each system is trained on a different database, none of the training sets includes individuals who are not cisnormative. We note that, as systems with gender classification only return binary labels, the non-binary genders are completely unable to be classified.

## 5.2 Inconsistent Labelling of a Datapoint

Consider now a datapoint whose identity shifts over time, because it is a fluid datapoint by definition. Identity is not static: it may move with respect to the labels we have, leading the datapoint to be configured in a label or in a different one during a selected time range. In this case, any extension of the label set is misleading, or at least insufficient. Here we cannot just add more categories, but we have to find a logical solution to changing the label of the same datapoint at different timepoints.

The problem here seems not to be related to the completeness of the labels as an issue to be solved at the pre-processing stage, but rather it concerns the training phase of the algorithm: data from individuals who identify with different labels at different times must also be passed to these algorithms, labeled with a tag corresponding to the gender term in which the subjects identify themselves at the current moment, hence correcting any previous labelling.

## 5.3 Enter Time

The two problems above can be formulated adding to completeness and consistency the dimension of temporality. Consider as a thought experiment a 1940-1950 dataset of white men in the Midwest: the probability that over time someone in this population changes his gender identity is likely to be close to zero, i.e. the probability of a label change for a corresponding classification task is minimal over time and most likely does not increase; on the other hand, in a dataset of young people in Berlin from the 90s, this probability increases in a shorter temporal range, i.e. the probability of someone changing the gender they identify with (and thus the correct label to attach to the corresponding datapoints) might be greater than zero, and might even increase in a (possibly shorter) time frame. Thus, an important starting point is represented by adding the dimension of timeliness, which concerns the degree to which data represent reality within a certain defined time range for a given population.

We suggest here considering the labelling task within a given time frame, whose length depends on the dataset and the classification task over the pairing of datapoints to labels, to measure a probability of a label-change over time. Intuitively, if the analysis is performed less than a certain number of timestamps away from the last data labeling, then we consider the labeling still valid. Otherwise, a new analysis with respect to both completeness of the dataset and label set must be performed. Technically, this means associating temporal parameters to labels and to compute the probability that a given label might change over the given time frame. The probability of a label being correct (its accuracy) decreases within the respective temporal window. In particular, reasoning on the temporal evolution of the dataset could allow us to model the evolution of the label partitions. Two fundamental theses are suggested for evaluation:

1. the correctness of the task does no longer assume static completeness of the label set, i.e. given the label set is complete at time $t_n$, it can be incomplete at time $t_{n+m}$;

2. the labelling does no longer assume static perseverance of the labels, that is, given a label $i$ that is correct at a time $t_n$ for a datapoint $d$, it could be incorrect at a later time, and conversely if it is incorrect it could become correct.

## 5.4 Back to Cleanlab

Considering a possible implementation in Cleanlab able to account for such differences, one requires renouncing the starting assumption on the categoricity of the data, instead one assumes that the probability of assigning a label may change over time. This can be formulated in two distinct ways.

First, the probability value of a given label $i$ being wrong, given a label $j$ is correct (their distance) may change over time. The task is now to give a mapping of all the label-variable pairs, i.e. given a mapping $y* \to \tilde{y}$ between variables, where $y*$ is the correct label and $\tilde{y}$ the wrong one, compute the probability over the time frame $\mathcal{T} := \{t_1, \ldots, t_n\}$

$$p_{\mathcal{T}}[(\tilde{y} = i)_{t_n} \mid (y* = j)_{t_{n-m}}] \tag{5}$$

such that label $i$ is wrong at time $t_n$, given that label $j$ was correct at time $t_{n-m}$. This probability can increase or decrease, depending on the dataset and on the label set, as informally illustrated by the thought experiment at the beginning of this Section. For the definition of the confident joint, this means taking the evaluation of all the elements that have an incorrect label $i$ when their correct label is $j$, and then associate the wrong label to a time $t_n$ and the correct label to a previous time. This estimate must be made on all time points, so for every $m < n$. Given a timepoint $n$ at which the label is wrong, the estimate on all pairs of probabilities for that point with a previous point in which another label can be correct has to be computed

$$C_{\tilde{y},y*}[i, j, \mathcal{T}] := \sum_{1 \leq m < n \in \mathcal{T}}^{n \in \mathcal{T}} | \hat{X}_{\tilde{y}=i_{t_n}, y*=j_{t_{n-m}}} | \tag{6}$$

Second, given a mapping $y* \to \tilde{y}$ between variables, where $y*$ is the correct label and $\tilde{y}$ the wrong one, what is the probability

$$p_{\mathcal{T}}[(\tilde{y} = i)_{t_n} \mid (y* = i)_{t_{n-m}}]$$

such that label $i$ is wrong at time $t_n$, given that the same label $i$ was correct at time $t_{n-m}$? In this case, the same label is fixed and the probability that it becomes incorrect can be calculated. The definition of confident joint becomes

$$C_{\tilde{y},y*}[i, \mathcal{T}] := \sum_{1 \leq m < n \in \mathcal{T}}^{n \in \mathcal{T}} | \hat{X}_{\tilde{y}=i_{t_n}, y*=i_{t_{n-m}}} | \tag{7}$$

To illustrate the point we consider a toy example. Compute

$$p(\tilde{y} = i \mid y* = j) = \frac{p(y* = j \mid \tilde{y} = i) \cdot p(\tilde{y} = i)}{p(y* = j)} \tag{8}$$

i.e. the error rate of $y* = male$ has to be determined. First, a confusion matrix is constructed to analyze errors. So, suppose to have a dataset of 10 datapoints, see Figure 1.

|  |  | Actual | |
| --- | --- | --- | --- |
|  |  | $y* = male$ | $y* = female$ |
| Predicted | $\tilde{y} = male$ | 4 | 2 |
|  | $\tilde{y} = female$ | 1 | 3 |

Figure 1: Confusion matrix at time $n$.

The matrix is $\begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix}$, from which $p(y* = j) = 5/10$ and $p(\tilde{y} = i) = 4/10$. So there are 5 women, of which 2 are incorrectly labeled "male" and 3 are correctly labeled "female", and 5 men of which 1 is incorrectly labeled "female" and 4 are correctly labeled "male". Replacing the values in Equation 8

$$p(\tilde{y} = i \mid y* = j) = \frac{p(y* = j \mid \tilde{y} = i) \cdot p(\tilde{y} = i)}{p(y* = j)} = \frac{\frac{p(\tilde{y}=i \wedge y*=j)}{p(\tilde{y}=i)} \cdot p(\tilde{y} = i)}{p(y* = j)} = \frac{\frac{1}{10}}{\frac{5}{10}} = 0.2$$

In this example, the obtained value represents the error rate of the "male" label, i.e. the probability of a male datapoint being labeled "female". Looking at the diagonals, the true positive rate TPR = 70% and the false positive rate FPR = 30%.

Consider now the same dataset at a later time $t_{n+m}$, see Figure 2.

|  | Actual | |
|---|---|---|
|  | $y* = male$ | $y* = female$ |
| $\tilde{y} = male$ | 2 | 3 |
| $\tilde{y} = female$ | 3 | 2 |

Figure 2: Confusion matrix at time $n + m$.

The labels might have changed. The matrix is $\begin{bmatrix} 2 & 3 \\ 3 & 2 \end{bmatrix}$, from which $p(y* = j) = 5/10$ and that $p(\tilde{y} = i) = 5/10$. Now there are 5 women, of which 3 are incorrectly labeled "male" and 2 are correctly labeled "female", and 5 men of which 3 are incorrectly labeled"female" and 2 are correctly labeled "male". Replacing again the values in 8

$$p'(\tilde{y} = i \mid y* = j) = \frac{p(y* = j \mid \tilde{y} = i) \cdot p(\tilde{y} = i)}{p(y* = j)} = \frac{\frac{p(\tilde{y}=i \wedge y*=j)}{p(\tilde{y}=i)} \cdot p(\tilde{y} = i)}{p(y* = j)} = \frac{\frac{3}{10}}{\frac{5}{10}} = 0.6$$

In this case the true positive rate TPR = 40% and the false positive rate FPR = 60%.

To understand how the error rate changes, the difference between the two matrices has to be considered

$$\begin{bmatrix} 2 & 3 \\ 3 & 2 \end{bmatrix} - \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} -2 & 1 \\ 2 & -1 \end{bmatrix} \tag{9}$$

The change rate can be computed as

$$\varepsilon = \hat{p}'(\tilde{y}; x_i; \theta) - \hat{p}(\tilde{y}; x_i; \theta) = 0.6 - 0.2 = 0.4 \tag{10}$$

Now $p_{\mathcal{T}}[(\tilde{y} = i)_{t_n} \mid (y* = j)_{t_{n-m}}]$ can be written as $p_{\mathcal{T}}[(\tilde{y} = i)_{t_{n+m}} \mid (y* = j)_{t_n}]$. Thus, at a time $t_n$ we have $p_{t_n}(y* = j) = 1 - p(\tilde{y} = i)_{t_n})$. At a subsequent time $t_{n+m}$ we have $p_{t_{n+m}}(y* = j) = 1 - p(\tilde{y} = i)_{t_{n+m}}$. Equation 8 can be computed with respect to time as

$$p_{\mathcal{T}}[(\tilde{y} = i)_{t_{n+m}} \mid (y* = j)_{t_n}] = \frac{p[(y* = j) \mid (\tilde{y} = i)]_{t_{n+m}} \cdot [p(\tilde{y} = i)t_n \pm \varepsilon]}{p(y* = j)_{t_n}} =$$

$$\frac{\frac{[p(y*=j \wedge \tilde{y}=i)]_{t_{n+m}}}{p(\tilde{y}=i)_{t_{n+m}}} \cdot [p(\tilde{y} = i)_{t_n} \pm \varepsilon]}{p(y* = j)_{t_n}} = \frac{\frac{3}{10}}{\frac{5}{10}} \cdot [\frac{4}{10} + (1 - \varepsilon)]}{\frac{5}{10}} = 0.288 \tag{11}$$

This value represents the (highest) probability that a given label is wrong at a given time, provided it was correct at some previous time. Indirectly, this also expresses the probability that the labelling set is applied to a dataset containing a point for which the labelling becomes inconsistent over time.

## 6 Temporal-based Fairness in AI

We have argued that a more general discussion on the data dimensions to be adopted in bias mitigation tools is needed, and in particular that the dimension of timeliness is crucial. In this Section we summarise our proposal and offer non exhaustive criteria for fairness in AI based on such temporal approach along with some basic theoretical results valid within such approach.

The first metric that has been addressed in this work is completeness as applied to the label set. In a world where gender classification is actually changing, the present strategy includes the completeness dimension in the quality assessment, verifying that the label set is complete with respect to the ontology of the world at the time this assessment is made. The solution here is to extend the label set as desired adding new labels for the classification task, as already suggested in Scheuerman et al. [2019]. Additionally, we suggest an explicit temporal parametrization: completeness can be considered as a relationship between a label set and an individual $p$ belonging to a certain population $P$, where $p$ is any domain item that enters $P$ at a time $t$. We must ensure that a correct label $l$ exists for each datapoint in the dataset at each time.

**Definition 1** (Completeness of a Label Set). *A label set $L$ for a classification algorithm in a AI system $X$ is considered complete over a time frame $\mathcal{T} : \{t_1, \ldots, t_n\}$ denoted as $Compl_{\mathcal{T}}(L(X))$ iff given two partitions $L_{t_1} := \{l_1, \ldots, l_n\}$ and $L_{t_n} := \{l'_1, \ldots, l'_n\}$, where possibly $L_{t_1} \cap L_{t_n} \neq \emptyset$ for all $(p \in P)_{\mathcal{T}}$ s.t. $p \in d(X)_{\mathcal{T}}$ there is $l \in L_{t_1} \cup L_{t_n}$ s.t. $y^*(d) = l$.*

In other words, the completeness of a dataset over a time frame is granted if for every datapoint representing an element in the population of interest there exists at any two possibly consecutive points in time a correct label for it.

Next, we considered consistency of the label set with respect to datapoints possibly shifting in categorization. The method here again has been to reduce consistency to timeliness. It has been suggested that the value of the probability of an inconsistency arising from a correct label change may be computed. Accuracy, albeit the most used metric for evaluating classification models' performances due to its easy calculability and interpretation, is reductive, trivial and incorrect in some contexts. For example, if the distribution of the class is distorted, accuracy is no longer a useful, nor a relevant metric. Even worse, sometimes greater accuracy leads to greater unfairness [Nielsen, 2020, p.59]: some labels like race or gender may allow models to be more predictive, although it seems to be often controversial to use such categories to increase predictive performance. As mentioned above, *fairness through unawareness* aims precisely at forgetting during the decision process those protected attributes which might have such a controversial nature.

On this basis we have suggested considering temporal accuracy as a function of the error rate over time. In particular, accuracy and error rate are complementary notions, i.e. $ACC = 1 - ER$ and $ER = 1 - ACC$. In fact, classification accuracy is calculated as the proportion of test set's examples which were predicted correctly over all predictions made on the same set, and the error rate as the proportion of test set's examples which were predicted incorrectly, over all predictions made on the same set

$$ACC = \frac{(TP + TN)}{(TP + FN + FP + TN)}$$

and

$$ER = \frac{(FP + FN)}{(TP + FN + FP + TN)}$$

There exists a functional relationship between accuracy and the change rate. In fact, if $\varepsilon = ER_2 - ER_1$ and $ER = 1 - ACC$, then $\varepsilon$ can be expressed as

$$\begin{aligned} \varepsilon &= (1 - ACC_2) - (1 - ACC_1) \\ &= 1 - ACC_2 - 1 + ACC_1 \\ &= \mid ACC_1 - ACC_2 \mid \end{aligned} \quad (12)$$

The value of $\varepsilon$ can be verified considering the confusion matrix $\begin{bmatrix} TP & FP \\ FN & TN \end{bmatrix}$. In our example from the previous Section, $\varepsilon$ is computed against "male" accuracy. The TP value of the matrix has to be considered, so those datapoints that are

correctly labeled "male". If matrix 1 is $\begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix}$ and matrix 2 is $\begin{bmatrix} 2 & 3 \\ 3 & 2 \end{bmatrix}$, then $ACC_1 = \frac{4}{5} = 0.8$ and $ACC_2 = \frac{2}{5} = 0.4$. The value of $\varepsilon$ can now be computed as

$$\varepsilon_{\mathcal{T}=\{t_1,\dots,t_n\}}(X) := \mid ACC_{t_1} - ACC_{t_n} \mid = \mid 0.4 - 0.8 \mid = 0.4 \tag{13}$$

Thus, the value of $\varepsilon$ over a time frame expresses the difference between accuracy at $t_1$ and at $t_n$.

The ability to compute the variance in the error rate across time is functional to determine the reliability of our AI systems. This metric is linked to the notion of accuracy, as it is considered as a measure of data correctness, see Batini and Scannapieco [2006]. For example in Wang and Strong [1996] and Batini et al. [2015], reliability is even contained in the definition of accuracy itself: data must be reliable to satisfy the accuracy dimension. Overall, it seems that reliability is not actually controlled beyond physical reliability, as in the literature on data quality there is no formal definition to compute it. However, following Eurostat [2015, 2020], Black and Van Nederpelt [2020] the previously provided temporal approach is again useful: evaluating reliability is based on the revisions which show how close the initial estimate of accuracy is to the following ones. In this sense, reliability can be reduced to accuracy over time in terms of a threshold on the error rate:

**Definition 2** (Reliability of a classification algorithm)**.** *A classification algorithm in a AI system $X$ is considered reliable over a time frame $\mathcal{T} := \{t_1, \dots, t_n\}$ denoted as $Rel_{\mathcal{T}}(X)$ iff $\varepsilon_{\mathcal{T}}(X) < \pi$, for some safe value $\pi$.*

The change rate $\varepsilon$ we have computed shows how much the system's accuracy deteriorates. If it exceeds a fixed safe value $\pi$, the system is no longer accurate. Plain accuracy is the numerical measure at some time $t \in \mathcal{T} := \{t_1, \dots, t_n\}$. If this value does not deteriorate over a certain fixed threshold, the system is considered reliable, and therefore accurate with respect to time.

The two previous definitions offer non-exhaustive criteria for the identification of fair AI systems:

**Definition 3** (Fairness for AI systems)**.** *$Fair_{\mathcal{T}}(X)$ only if $Rel_{\mathcal{T}}(X)$ and $Compl_{\mathcal{T}}(L(X))$.*

Hence we claim that fairness requires the system's ability to give reliable and correct outcomes over time. While we do not consider these properties sufficient, we believe they are necessary . On this basis, we can formulate two immediate theoretical results:

**Theorem 1.** *Given a label set $L$ complete at time $t$, a classification algorithm guarantees a fair classification at time $t' > t$ if and only if the change rate determined with respect to $L$ is $\epsilon < \pi$.*

*Proof.* Assume $Compl_t(L(X))$, then for $Fair_{t'}(X)$ we need to show $Rel_{t'}(X)$ for $t' > t \in \mathcal{T}$. Assume $\epsilon > \pi$, then by Definition 2 reliability is not satisfied; hence, if $Rel_{\mathcal{T}}(X)$, it must be the case that $\epsilon < \pi$. $\qquad\square$

**Theorem 2.** *Given a fixed change rate $\epsilon < \pi$, a classification algorithm with fair behaviour at time $t$ remains fair at time $t' > t$ if and only if the change to make the label set complete at time $t'$ does not exceed an $\epsilon'$ such that $\epsilon + \epsilon' > \pi$.*

*Proof.* Consider $Fair_t(X)$ with change rate $0 < \pi$ as a base case, then by Definition 3 $Rel_t(X)$ and $Compl_t(L(X))$. Now consider $t' > t$ and a required change $\epsilon'$ in $Compl_t(L(X))$ such that $Rel_{t'}(X)$ holds. This obviously holds only if $0 + \epsilon' < \pi$. Generalize for any $\epsilon > 0$. $\qquad\square$

Note that in these result the value of $\epsilon$, respectively $\epsilon'$, is a proxy for how much the world has changed at $t'$ with respect to $Compl_t(L(X))$.

What types of bias can be identified under this theoretical understanding of fairness? In the context of an incomplete label set, a detected label bias can originate from an exclusion bias in data, which can also result from a time interval bias. In the case of label-changing datapoints a chronological bias occurs. Then, misclassification bias can be reduced to the two previous types. In the context of use, emergent bias can arise as a result of changes in societies and cultures. It might appear in data as chronological, historical or behavioral bias. Here, a different value bias occurs for example when the users are different from the assumed ones during the system's development. This is the case of ontology switching, to which a label set must adapt. Those types of bias could all be mitigated by implementing the proposed framework.

# 7 Conclusion

We presented some recommendations for AI systems design, focusing on timeliness as a founding dimension for developing fairer and more inclusive classification tools. Despite the crucial importance of accuracy metrics as shown by significant works such as Buolamwini and Gebru [2018] and Angwin et al. [2016], the problem of unfairness in AI systems is much broader and more foundational. This can be expressed in terms of data quality: AI systems are limited in that they maximize accuracy, and even if systems become statistically accurate some problems remain unsolved.

Technological problems have become cultural ones. And, as the work of classification is always a reflection of culture, the completeness of the label set and the (constrained) consistency of labelling have an epistemological value: constructing AIs requires us to understand society, and society reflects an ontology of individuals. Misgendering is first of all an ontological error [Keyes, 2018].

We suggested that timeliness is crucial for the definition of gender identity. If we are ready to consider gender as a property that shifts over time [Ruberg and Ruelos, 2020, p.6], and that can also be not singular, as an individual may identify under more than one - not mutually exclusive - labels, then a change of paradigm is required. Current classification methods, in fact, are deeply rooted on three assumptions on gender:

1. Binarism, gender can be classified only in two mutually exclusive categories;
2. Staticity, gender identity category remains unaltered once being assigned;
3. Derivability from physical traits, gender can be correctly inferred from visual clues.

These design limitations must be addressed if fairer classifications and more inclusive models of gender are to be designed.

Further work in this direction includes: an implementation and empirical validation of the proposed model; the formulation of further theoretical results, e.g. in relation to impossibility results for individual and group fairness; the design of an extension to compute the probability of incorrect labels becoming correct over time (i.e. the dual case of what presently addressed); the extension to the temporal considerations introduced in the present work of formal languages designed to facilitate checking correctness of AI models against desirable, transparently obtained classification outputs as e.g. in D'Asaro and Primiero [2021], Primiero and D'Asaro [2022], Termine et al. [2021], D'Asaro et al. [2023].

# References

T. H. Aasheim, K. Hufthammer, S. Ånneland, H. Brynjulfsen, and M. Slavkovik. Bias mitigation with aif360: A comparative study. In *Proceedings of the NIK-2020 Conference*, 2020. URL `http://www.nik.no/`.

D. Angluin and P. D. Laird. Learning from noisy examples. *Mach. Learn.*, 2(4):343–370, 1987.

J. Angwin, J. Larson, S. Mattu, and L. Kirchner. Machine bias. *ProPublica*, 2016.

C. Batini and M. Scannapieco. *Data Quality: Concepts, Methodologies and Techniques*. Springer, 2006.

C. Batini, A. Rula, M. Scannapieco, and G. Viscusi. From data quality to big data quality. *Journal of Database Management*, 26:60–82, 2015. doi:10.4018/JDM.2015010103.

R. K. E. Bellamy, K. Dey, M. Hind, S. C. Hoffman, S. Houde, K. Kannan, P. Lohia, S. Martino, J. and. Mehta, A. Mojsilovic, S. Nagar, K. N. Ramamurthy, J. Richards, D. Saha, P. Sattigeri, M. Singh, K. R. Varshney, and Y. Zhang. Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias, 2018.

Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in criminal justice risk assessments: The state of the art, 2017. URL `https://arxiv.org/abs/1703.09207`.

A. Black and P. Van Nederpelt. *Dimensions of Data Quality (DDQ)*. DAMA NL Foundation, 2020.

J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency, FAT 2018, 23-24 February 2018, New York, NY, USA*, volume 81, pages 77–91. PMLR, 2018.

University of Oxford Centre for Evidence-Based. Catalogue of bias, 2022. `https://catalogofbias.org/biases/`.

CertNexus. Promote the ethical use of data-driven technologies, 2021. `https://www.coursera.org/learn/promote-ethical-data-driven-technologies/lecture/5Ufbp/data-collection-bias`.

Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments, 2016.

Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *CoRR*, abs/1808.00023, 2018. URL `http://arxiv.org/abs/1808.00023`.

Alexander D'Amour, Hansa Srinivasan, James Atwood, Pallavi Baljekar, D. Sculley, and Yoni Halpern. Fairness is not static: deeper understanding of long term fairness via simulation studies. pages 525–534, 01 2020. doi:10.1145/3351095.3372878.

Fabio Aurelio D'Asaro and Giuseppe Primiero. Probabilistic typed natural deduction for trustworthy computations. In Dongxia Wang, Rino Falcone, and Jie Zhang, editors, *Proceedings of the 22nd International Workshop on Trust in Agent Societies (TRUST 2021) Co-located with the 20th International Conferences on Autonomous Agents and Multiagent Systems (AAMAS 2021), London, UK, May 3-7, 2021*, volume 3022 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2021. URL `http://ceur-ws.org/Vol-3022/paper3.pdf`.

Fabio Aurelio D'Asaro, Francesco Genco, and Giuseppe Primiero. Checking trustworthiness of probabilistic computations in a typed natural deduction system, 2023.

Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard S. Zemel. Fairness through awareness. *CoRR*, abs/1104.3913, 2011. URL `http://arxiv.org/abs/1104.3913`.

B. D'Alessandro, C. O'Neil, and T. LaGatta. Conscientious classification: A data scientist's guide to discrimination-aware classification. *Big Data*, 5(2):120–134, 2017.

Eurostat. *ESS Handbook for Quality Reports*. European Statistical System, Brussels, Belgium, 2015.

Eurostat. *European Statistical System handbook for quality and metadata reports 2020 edition*. European Statistical System, Luxembourg, 2020.

S. Fabbrizzi, S. Papadopoulos, E. Ntoutsi, and I. Kompatsiaris. A survey on bias in visual datasets. *CoRR*, abs/2107.07919, 2021.

Miranda Fricker. *Epistemic Injustice: Power and the Ethics of Knowing*. New York: Oxford University Press, 2007.

Sorelle A. Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. On the (im)possibility of fairness. *CoRR*, abs/1609.07236, 2016. URL `http://arxiv.org/abs/1609.07236`.

B. Friedman and H. Nissenbaum. Bias in computer systems. *ACM Trans. Inf. Syst.*, 14(3):330–347, 1996.

Nina Grgic-Hlaca, Muhammad Bilal Zafar, Krishna P. Gummadi, and Adrian Weller. The case for process fairness in learning: Feature selection for fair decision making. 2016.

Foad Hamidi, Morgan Klaus Scheuerman, and Stacy M. Branham. Gender recognition or gender reductionism? the social implications of embedded gender recognition systems. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, CHI '18, page 1–13, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450356206. doi:10.1145/3173574.3173582. URL `https://doi.org/10.1145/3173574.3173582`.

Alex Hanna, Madeleine Pape, and Morgan Klaus Scheuerman. Auto-essentialization: Gender in automated facial analysis as extended colonial project. *Big Data and Society*, 8(2), 2021. doi:10.1177/20539517211053712.

Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. *CoRR*, abs/1610.02413, 2016. URL `http://arxiv.org/abs/1610.02413`.

S. Hooker. Moving beyond algorithmic bias is a data problem. *Patterns*, 2(4), 2021.

P. Illari and L. Floridi. *The Philosophy of Information Quality*. Springer International Publishing, 2014.

B. K. Kahn, D. M. Strong, and R. Y. Wang. Information quality benchmarks: Product and service performance. *Commun. ACM*, 45(4):184–192, 2002.

O. Keyes. The misgendering machines: Trans/hci implications of automatic gender recognition. *Proceedings of the ACM on Human-Computer Interaction*, 2:1–22, 2018.

Jon M. Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. *CoRR*, abs/1609.05807, 2016. URL `http://arxiv.org/abs/1609.05807`.

Matt J. Kusner, Joshua R. Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness, 2018.

Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *CoRR*, abs/1908.09635, 2019. URL `http://arxiv.org/abs/1908.09635`.

M. Miceli, J. Posada, and T. Yang. Studying up machine learning data: Why talk about bias when we mean power? *CoRR*, abs/2109.08131, 2021.

Aileen Nielsen. *Practical Fairness*. O'Reilly Media, Inc., 2020. ISBN 9781492075738. URL `http://gen.lib.rus.ec/book/index.php?md5=F9752B2F9693C98855A51504FE224DF6`.

C G. Northcutt, W. Tailin, and I. L. Chuang. Learning with confident examples: Rank pruning for robust classification with noisy labels. In *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence*, Sydney, Australia, 2017. AUAI Press.

C. G. Northcutt, A. Athalye, and J. Mueller. Pervasive label errors in test sets destabilize machine learning benchmarks, 2021a. Preprint at `https://arxiv.org/pdf/2103.14749.pdf`.

C. G. Northcutt, L. Jiang, and I. L. Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research (JAIR)*, 70:1373–1411, 2021b.

A. Olteanu, C. Castillo, F. Diaz, and E. Kiciman. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in Big Data*, 2, 2019.

L. Pipino, Y. W. Lee, and R. Y. Wang. Data quality assessment. *Commun. ACM*, 45(4):211–218, 2002.

Giuseppe Primiero and Fabio Aurelio D'Asaro. Proof-checking bias in labeling methods. In Guido Boella, Fabio Aurelio D'Asaro, Abeer Dyoub, and Giuseppe Primiero, editors, *Proceedings of 1st Workshop on Bias, Ethical AI, Explainability and the Role of Logic and Logic Programming (BEWARE 2022) co-located with the 21th International Conference of the Italian Association for Artificial Intelligence (AI\*IA 2022), Udine, Italy, December 2, 2022*, volume 3319 of *CEUR Workshop Proceedings*, pages 9–19. CEUR-WS.org, 2022. URL `https://ceur-ws.org/Vol-3319/paper1.pdf`.

B. Ruberg and S. Ruelos. Data for queer lives: How lgbtq gender and sexuality identities challenge norms of demographics. *Big Data and Society*, 7, 2020. doi:10.1177/2053951720933286.

K. K. Saravanakumar. The impossibility theorem of machine fairness - A causal perspective. *CoRR*, abs/2007.06024, 2020. URL `https://arxiv.org/abs/2007.06024`.

Teresa Scantamburlo. Non-empirical problems in fair machine learning. *Ethics Inf. Technol.*, 23(4):703–712, 2021. doi:10.1007/s10676-021-09608-9. URL `https://doi.org/10.1007/s10676-021-09608-9`.

M. K. Scheuerman, J. Paul, and J. Brubaker. How computers see gender: An evaluation of gender classification in commercial facial analysis services. *Proceedings of the ACM on Human-Computer Interaction*, 3:1–33, 2019. doi:10.1145/3359246.

H. Suresh and J. Guttag. A framework for understanding sources of harm throughout the machine learning life cycle. *Equity and Access in Algorithms, Mechanisms, and Optimization*, 2021.

Alberto Termine, Giuseppe Primiero, and Fabio Aurelio D'Asaro. Modelling accuracy and trustworthiness of explaining agents. In Sujata Ghosh and Thomas Icard, editors, *Logic, Rationality, and Interaction - 8th International Workshop, LORI 2021, Xi'ian, China, October 16-18, 2021, Proceedings*, volume 13039 of *Lecture Notes in Computer Science*, pages 232–245. Springer, 2021. doi:10.1007/978-3-030-88708-7_19. URL `https://doi.org/10.1007/978-3-030-88708-7_19`.

Sahil Verma and Julia Rubin. Fairness definitions explained. In *Proceedings of the International Workshop on Software Fairness*, FairWare '18, page 1–7, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450357463. doi:10.1145/3194770.3194776. URL `https://doi.org/10.1145/3194770.3194776`.

R. Y. Wang and D. M. Strong. Beyond accuracy: What data quality means to data consumers. *J. Manag. Inf. Syst.*, 12 (4):5–33, 1996.