

School of Computer Science and Information Technology  
Lucerne University of Applied Sciences and Arts (Switzerland)

# DEMOGRAPHIC BIASES IN DERMATOLOGY MODELS

TODO: subtitle

## BACHELOR THESIS

presented to School of Computer Science and Information Technology of Lucerne  
University of Applied Sciences and Arts (Switzerland) in consideration for the award of  
the academic grade of *Bachelor in Computer Science*.

by

**Nadja Stadelmann**

from

Lucerne (Switzerland)

# Declaration

Bachelor Thesis at Lucerne University of Applied Sciences  
and Arts  
School of Computer Science and Information Technology

Title of Bachelor Thesis:	Demographic Biases inDermatology Models
Name of Student:	Nadja Stadelmann
Degree Program:	Bachelor in Computer Science
Year of Graduation:	2025
Main Advisor:	Dr. Ludovic Amruthalingam
External Expert:	Dr. Jürg Schelldorfer
Industry partner/provider:	Applied AI Research Lab

## Code/Thesis Classification

- ☒ Public (Standard)  
☐ Private

## Declaration

I hereby declare that I have completed this thesis alone and without any unauthorized or external help. I further declare that all the sources, references, literature and any other associated resources have been correctly and appropriately cited and referenced. The confidentiality of the project provider (industry partner) as well as the intellectual property rights of the Lucerne University of Applied Sciences and Arts have been fully and entirely respected in completion of this thesis.

Place/Date, Signature \_\_\_\_\_

## Submission of the Thesis to the Portfolio Database

Confirmation by the student

I hereby confirm that this bachelor thesis has been correctly uploaded to the Portfolio Database in line with the code of practice of the University. I rescind all responsibility and authorization after upload so that no changes or amendments to the document may be undertaken.

Place/Date, Signature \_\_\_\_\_

## Expression of thanks and gratitude

Thanks to my family, relatives and friends for all the support given to finish this thesis. **TODO: add thanks and gratitude** Ludovic Amruthalingam Simone Lionetti - deputy Ludovic Pascal Baumann - LaTeX Philippe Gottfrois - information and work on PASSION project

Nadja Stadelmann, 2025

*Intellectual property of the degree programs of the Lucerne University of Applied Sciences and Arts, FH Zentralschweiz, in accordance with Student Regulations: Studienordnung*

# Summary

TODO: Your abstract here. The content of your thesis in brief.

# Contents

<b>1</b>	<b>Problem Statement</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Objective . . . . .	2
<b>2</b>	<b>State of Research</b>	<b>3</b>
2.1	PASSION for Dermatology . . . . .	3
2.2	Bias . . . . .	7
2.3	Fairness . . . . .	23
2.4	Mitigation Methods . . . . .	23
2.5	Extensive Sources . . . . .	23
2.6	General ML biases . . . . .	24
2.7	Statistical biases . . . . .	42
2.8	Dermatology Bias . . . . .	42
<b>3</b>	<b>Ideas and Concepts</b>	<b>48</b>
3.1	PASSION Dataset . . . . .	48
3.2	Broad Methodology . . . . .	48
<b>4</b>	<b>Methods</b>	<b>50</b>
<b>5</b>	<b>Execution</b>	<b>51</b>
<b>6</b>	<b>Evaluation and Validation</b>	<b>52</b>
<b>7</b>	<b>Outlook</b>	<b>53</b>
<b>8</b>	<b>Glossary</b>	<b>54</b>
<b>9</b>	<b>Bibliography</b>	<b>55</b>

# Todo list

TODO: solve todos

TODO: also solve todos in the code ;)

TODO: also fix metadata entry!!!

TODO: Portfolio DB für Referenzarbeiten anschauen

TODO: remove all `\rawcitationstart \rawcitationend \baaCriteria`

Alle Fakten (fundiertes Wissen Dritter) sind korrekt zitiert. Es werden verschiedene Zitierweisen verwendet und teilweise mehrere Interpretationen gegenübergestellt. Der gemeinsam definierte Zitierstil im Text, in Abbildungen und Tabellen sowie im Literaturverzeichnis wird korrekt und durchgängig angewendet. Eigene Leistungen (sowie Bewertungen) und Fremdquellen sowie Recherchen sind klar unterscheidbar.

Die erstellten Artefakte sind von sehr hoher Qualität. Das trifft u.a. auf Diagramme, Skizzen sowie Notationen (z.B. BPMN/UML) zu. Darstellungen sind einwandfrei, alle statistisch notwendigen Qualitätskriterien sind erfüllt. Beschriftungen etc. sind vorhanden, keine Einwände, Text und Bild stimmen beschreibend gut überein. Es wurden angemessene Dokumentationsmethoden und -arten korrekt verwendet. Vereinbarte Interview Transkripte, Beobachtungsprotokolle bzw. Zusammenfassungen sind vorhanden. Daten, Ort, Kontext, Beschreibung, Zeilennummer, Verweise, Strukturen sind erkennbar, gut formatiert und korrekt mit dem Text/ der Analyse verknüpft. Alle Elemente und Themen sind im methodischen Teil/Text erklärt und verständlich, keine technischen oder strukturellen Einwände. Auch Zwischenanalysen, Zwischenschritte oder Gesamtauswertungen wurden durchgeführt, die Herkunft der Daten ist erkennbar und professionell aufbereitet.

Der Schreibstil aller Dokumente entspricht hohen Standards und enthält keine Übertreibungen oder unbegründete Beurteilungen. Die Sprache ist aussagekräftig, prägnant und präzise. Die Fachterminologie ist konsistent, d.h. für gleiche Gegenstände und Themen werden immer die gleichen Begriffe verwendet. Der Sprachgebrauch ist durchgängig geschlechtergerecht, einheitlich und sachlich.

# List of Figures

2.1	Bias definitions in a ML lifecycle (Mehrabi et al., 2021). . . . .	8
-----	--	---

# List of Tables

2.1	PASSION dataset - labels and descriptions (Gottfrois et al., 2024) .	4
2.2	PASSION dataset - existing analysis scripts (Gottfrois et al., 2024)	
	<b>TODO: decide on a table style</b> . . . . .	5
2.3	Biases - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness . . . . .	10
2.4	Features which often hold biases - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness . . . . .	22
2.5	Fairness Definitions - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness . . . . .	30
2.6	Mitigation Methods - Unbiasing Data - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness . . . . .	33
2.7	Mitigation Methods - Fair Classification - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness . . . . .	34
2.8	Mitigation Methods - Others - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness . . . . .	35



# Glossary

**Fitzpatrick skin type** A skin classifier based on the skins' reaction to light (Gotfrois et al., 2024). ix, 2,

**Jupyter Notebook** Executable files, often used in ML to write Python code and add explanations in text form. 4

**pediatric** A medical term for infants, children and adolescents. 1, 2

**proxy variable** "one or more variables that encode the protected attribute with a substantial degree of accuracy" according to <https://medium.com/bcggamma/practice-ai-responsibly-with-proxy-variable-detection-42c2156ad986>. 14

# Acronyms

**FST** Fitzpatrick skin type. *Glossary:* Fitzpatrick skin type, 2, 4, 5, 11, 48

**HSLU** long. 2

# 1 Problem Statement

Welche Ziele, Fragestellungen werden mit dem Projekt verfolgt? Die Bedeutung, Auswirkung und Relevanz dieses Projektes für die unterschiedlichen Beteiligten soll aufgeführt werden. Typischerweise wird hier ein Verweis auf die Aufgabenstellung im Anhang gemacht.

TODO: Formulate statement from those citations:

- AI systems can be used in many sensitive environments to make important and life-changing decisions; thus, it is crucial to ensure that these decisions do not reflect discriminatory behavior toward certain groups or populations (Mehrabi et al., 2021).
- There are clear benefits to algorithmic decision-making; unlike people, machines do not become tired or bored [45, 119], and can take into account orders of magnitude more factors than people can. However, like people, algorithms are vulnerable to biases that render their decisions “unfair” [6, 121]. In the context of decision-making, fairness is *the absence of any prejudice or favoritism toward an individual or group based on their inherent or acquired characteristics*. Thus, an unfair algorithm is one whose decisions are skewed toward a particular group of people. (Mehrabi et al., 2021).
- it is important for researchers and engineers to be concerned about the downstream applications and their potential harmful effects when modeling an algorithm or a system (Mehrabi et al., 2021).
- We should think responsibly, and recognize that the application of these tools, and their subsequent decisions affect peoples’ lives; therefore, considering fairness constraints is a crucial task while designing and engineering these types of sensitive tools (Mehrabi et al., 2021).

## 1.1 Context

This thesis is part of the PASSION project. The PASSION research team identified that in Africa, dermatology treatment is not accessible. There is less than one dermatologist per one million citizens. In contrast, there is high demand for dermatology treatment, especially among children and adolescents. 80% of the pediatric population is affected. The goal of PASSION is to make dermatology treatment more accessible by using AI supported telemedicine for triage (Gottfrois et al., 2024).

For AI supported triage, demographic biases in existing dermatology models is a problem since the corresponding datasets lack diversity, especially regarding skin tones (Gottfrois et al., 2024). This type of bias is important in dermatology, since different diseases present themselves differently depending on the skin-color (Diaz et al., 2022). Further, skin diseases are more advanced or severe at diagnosis in patients with lower socioeconomic status (British Association of Dermatologists (BAD), 2021).

PASSION tries to mitigate the demographic bias by providing a dataset of pigmented skin images of patients from Sub-Saharan Africa. The PASSION team focused on gathering data with Fitzpatrick skin type (FST) IV, V and VI. Further, the covered conditions represent up to 80% of the conditions in the pediatric population, the demographic group who is most affected by skin disease (Gottfrois et al., 2024). long (HSLU)

The PASSION dataset is complementary to the existing datasets and improves the diversity in a combined dataset. Within the dataset itself, there could potentially be further demographic biases, e.g. related to age or gender.

## 1.2 Objective

The goal of this research is to

1. Identify demographic biases in dermatology AI models, using established fairness metrics.
2. Identify mitigation strategies to minimize these biases.
3. Assess the effectiveness of the mitigation strategies.

It is important to identify the existent biases first, so that the mitigation strategies can be **TODO: proceed here to reason why you chose those objectives**

## 2 State of Research

Bezogen auf die eigenen Zielsetzungen und Fragestellungen soll aufgezeigt werden, wie andere dieses oder ähnliche Probleme gelöst haben. Worauf können Sie aufbauen, was müssen Sie neu angehen? Wodurch unterscheidet sich Ihre Lösung von anderen Lösungen? Für wissenschaftlich orientierte Arbeiten sei hier explizit auf (Balzert, S. 66 ff) verwiesen. Relevante, aktuelle und fundierte Fachliteratur wurde identifiziert, kritisch geprüft und verwendet. Die Begriffe der Fragestellung sind definiert und referenziert. Der gesamte Kontext ist verknüpft und eine Abgrenzung wurde vorgenommen. All dies ist in einer leicht verständlichen Struktur formuliert und überprüft.

### 2.1 PASSION for Dermatology

The PASSION research team provides a dataset including three analysis scripts and an AI model. For this thesis, it is important to understand which labels the dataset provides, so that the applicable bias mitigation methodologies can be chosen.

The provided analysis scripts show a first insight into the demographic distribution in the dataset, such as Fitzpatrick skin type and cases per country distribution. The results of those analyses reveal first biases.

There are also dermatology specific analysis scripts in regards of body localization by condition or impetigo cases. Those results

#### 2.1.1 PASSION Dataset

The PASSION dataset contains data from patients from four African countries. It contains 4901 images of dermatology cases and the corresponding demographic and clinical information, see Table 2.1. There is one record per patient and one or more corresponding images. The images are linked with the record by filename, which contains the `subject_id` of the row entry. Access to the dataset can be requested via <https://passionderm.github.io/> (Gottfrois et al., 2024).

Label	Data Type	Description
subject_id	string	Participant's unique identifier
country	string	Country of origin of the participant
age	integer	Age of the participant in years
sex	m/f/o	Gender of the participant
fitzpatrick	integer	FST
body_loc	string (list; null-able, semicolon-separated)	Specific affected body locations
impetig	0/1	Presence of impetigo (1=present), may occur alone or with other conditions, affects the treatment options for coexisting conditions
conditions_PASSION	Eczema, Scabies, Fungal, Others	Primary diagnosed skin condition

Table 2.1: PASSION dataset - labels and descriptions (Gottfrois et al., 2024)

### 2.1.2 PASSION Analysis Scripts

With the Dataset, the PASSION research team provides a Jupyter Notebook with code examples and analysis scripts. They are listed in Table 2.2 with a description and an indicator, how relevant the scripts are for this thesis.

Script Title	Description	Relevance - Reasoning
Linking CSV Data with Image Files	Creates mapping between the data records and images. It further counts the cases by country	<b>High</b> - Basis for other analysis's, potentially provides dermatological info
Extracting and Comparing Subject IDs	Checks the dataset complecity and accuracy in regards of linking records and images	<b>Low</b> - Checks loaded data for completeness, but is not providing more insight
Regrouping Malawi and Tanzania to EAS	data aggregation due to dataset size and geographical proximity	<b>Low</b> - Might be relevant to understand the dataset and for interpreting the results of the following scripts correctly
Conditions by Country	Relationship between clinical conditions and country	<b>Medium</b> - Currently unsure whether this information is relevant for this thesis <b>TODO: research relevance between country vs. clinical conditions in regards of demographic bias</b>
Body Localizations by Conditions	Shows correlation between the condition and primarily affected body parts; does not use all affected body parts listed in the data <b>TODO: check with Philippe why this was done</b>	<b>Low</b> - While the correlation can be interesting for other research, it is not relevant for demographic biases.
Impetigo Cases	Counts total number of impetigo cases as well as proportion to all cases	<b>Medium</b> - Currently unsure whether this information is relevant for this thesis <b>TODO: research relevance between impedigo and demographic bias</b>
Distribution of Fitzpatrick Skin Types	Counts and visualizes the skin type distribution	<b>High</b> - FST is a demographic information

Table 2.2: PASSION dataset - existing analysis scripts (Gottfrois et al., 2024)  
**TODO: decide on a table style**

### 2.1.3 PASSION Experiments

see <https://github.com/Digital-Dermatology/PASSION-Evaluation>

### 2.1.4 Telemedicine

- **Teledermatology.** Telemedicine may be one of the first fields to embrace AI, driven by demand for services, the necessity of collecting fit-for-purpose high-quality images, and the availability of existing technology (Xiong et al., 2019). Face-to-face diagnostic accuracy exceeds that of teledermatology (Finnane et al., 2017); however, inequalities surrounding access to dermatological care persist. Teledermatology has the potential to increase access by facilitating referrals and offering convenience and decreased wait times (Finnane et al., 2017), as well as providing diagnostic support at the time of case review. For teledermatology cases, the accuracy of a DL classifier (0.67) matched dermatologists' (0.63) and was higher than primary care physicians' (0.45) for 26 skin conditions (Liu et al., 2019b). AI may be integrated into smartphone apps to photograph skin lesions, collect relevant clinical information, and generate a referral if appropriate. Many smartphones already support on-device DL with Google's TensorFlow Lite (TensorFlow, 2020) or Apple's CoreML (Apple Inc, 2020), preserving privacy by keeping health information on the device. A systematic review found nine studies that evaluated six algorithm-based smartphone apps and concluded that evidence of diagnostic accuracy was poor and does not support current implementation, despite two apps having obtained the CE marking; no apps are Food and Drug Administration approved (Freeman et al., 2020). AI may also assist in automatic tracking and monitoring of skin lesions; although preliminary results are promising, existing studies used small datasets with little description, and there is no established standard metric of change (Navarro et al., 2019). Further study hinges on the prospective collection of large datasets. (Young et al., 2020)
- **(Tsetsi\_2017)** on smartphone / internet access divide between populations
- <https://www.tandfonline.com/doi/full/10.1080/08870446.2019.1579330> on how open people are to use AI
- **Image quality.** Several barriers to AI implementation in the clinic need to be overcome with regards to imaging (Figure 1). These include technical variations (e.g., camera hardware and software) and differences in image acquisition and quality (e.g., zoom level, focus, lighting, and presence of hair). For example, the presence of surgical ink markings is associated with decreased specificity (Winkler et al., 2019), field of view can significantly affect prediction quality (Mishra et al., 2019), and classification performance improves when hair and rulers are removed (Bisla et al., 2019). We have developed a method to measure how model predictions might be biased by the presence of a visual artifact (e.g., ink) and proposed methods to reduce such biases (Pfau et al., 2019). Poor quality images are often excluded from studies, but the problem of what makes an image adequate is not well studied. Ideally, models need to be able to express a level of confidence in a prediction as a function of image quality and appropriately direct a user to retake photos if needed. (Young et al., 2020)



## 2.2 Bias

The usage of AI with its algorithmic decisions affect peoples' lives. Due to that, tools in healthcare, like the PASSION project which should make triage and potentially diagnosis easier (Gottfrois et al., 2024), are sensitive tools. AI application decisions proved to be biased. Therefore, AI engineers should aim to address and mitigate those biases. To do so, it is crucial for AI engineers to know what bias is, what types of biases exist and where they are coming from. Further, AI engineers need to know what prevention methods are available (Mehrabi et al., 2021).

- Bias in facial recognition systems [128] and recommender systems [140] have also been largely studied and evaluated and in many cases shown to be discriminative towards certain populations and subgroups. In order to be able to address the bias issue in these applications, it is important for us to know where these biases are coming from and what we can do to prevent them.(Mehrabi et al., 2021).
- We should think responsibly, and recognize that the application of these tools, and their subsequent decisions affect peoples' lives; therefore, considering fairness constraints is a crucial task while designing and engineering these types of sensitive tools (Mehrabi et al., 2021).

This chapter provides an overview over biases and related features which were mentioned in ML and dermatology related research.

### 2.2.1 Introduction to Bias in AI

According to the Cambridge English dictionary, bias can be defined as "the fact of preferring a particular subject or thing" or even as "the action of supporting or opposing a particular person or thing in an unfair way, because of allowing personal opinions to influence your judgment" (**Cambridge\_2025**). While an AI does not hold personal opinions, the judgment of an AI algorithm is influenced by the underlying data and how the algorithm uses this data. The data and even the algorithm itself can hold biases, which affects the final outcome and potential lead to unfair decisions. In decision-making, fairness means "absence of any prejudice or favoritism toward an individual or group based on their inherent or acquired characteristics" - or in other terms, a fair algorithms decisions are not scewed toward a group of people (Mehrabi et al., 2021). **TODO: Do I need to say more about this? If so, add infos from section Discrimination vs. Biases**

- The Cambridge English dictionary defines bias as "the action of supporting or opposing a particular person or thing in an unfair way as a result of allowing personal opinions to influence your judgement."1 However, statistical bias is defined as any systematic error in the determination of the association between exposure and disease.2 (Chakraborty, 2024)
- These biased predictions stem from the hidden or neglected biases in data or algorithms (Mehrabi et al., 2021).

- There are clear benefits to algorithmic decision-making; unlike people, machines do not become tired or bored [45, 119], and can take into account orders of magnitude more factors than people can. However, like people, algorithms are vulnerable to biases that render their decisions “unfair” [6, 121]. In the context of decision-making, fairness is *the absence of any prejudice or favoritism toward an individual or group based on their inherent or acquired characteristics*. Thus, an unfair algorithm is one whose decisions are skewed toward a particular group of people. (Mehrabi et al., 2021).

### 2.2.2 Bias Sources

The general ML lifecycle consists of data gathering, training the algorithm and the user interaction with the trained model. Now, while data gathering, biases can arise either through the collection process or it is already inherited in the available data. Further, depending on the algorithm design, during training, the existing bias in the data can be amplified and new bias can be introduced. Lastly, the result of the algorithm can affect the user experience on inference which can lead to further bias amplification. This generates a feedback loop between the biases in each step of the ML lifecycle which can make it hard to identify the original bias source. The feedback loop is illustrated in 2.1, which also shows first bias definitions, which were categorized according to this feedback loop (Mehrabi et al., 2021).

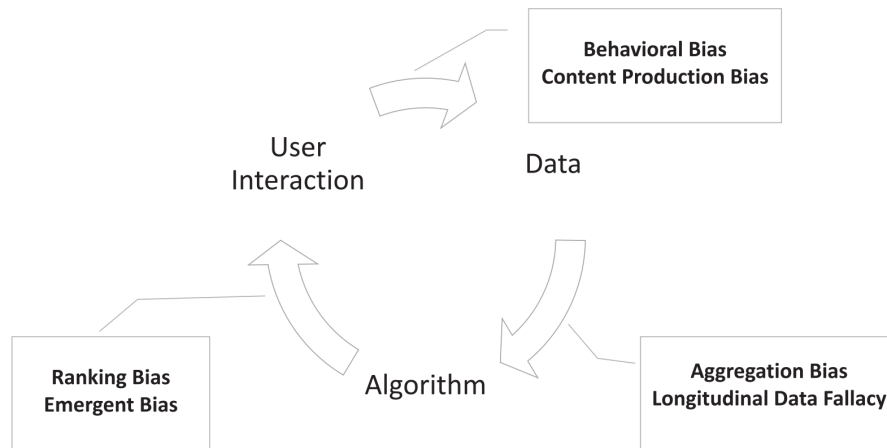


Figure 2.1: Bias definitions in a ML lifecycle (Mehrabi et al., 2021).

- two potential sources of unfairness in machine learning outcomes - those that arise from biases in the data and those that arise from the algorithms ... we observe that biased algorithmic outcomes might impact user experience, thus generating a feedback loop between data, algorithms and users that can perpetuate and even amplify existing sources of bias (Mehrabi et al., 2021).
- The loop capturing this feedback between biases in data, algorithms, and user interaction is illustrated in Figure 1. We use this loop to categorize definitions of bias in the section below (Mehrabi et al., 2021)

### 2.2.3 Bias Types

The Table 2.3 aims to provide an overview over what kind of biases exist according to research. The more detailed categories listed in the table try to capture similar kind of biases. This thesis follows roughly the categorization of Mehrabi et al., 2021. Some biases might actually fit in multiple categories. The definition of the categories including examples of specific biases follows.

TODO: check categorization of table ref18, statistical biases  
TODO: check the c citations pls

Bias	Mentioned in Context of	
	ML	Dermatology
<b>Data Biases</b>		
<i><b>Sampling Biases</b></i> <span style="color: red;">TODO: specify which one of note 18, was selection and survivorship, maybe also add HP_2022</span>	X <sup>1,16,18</sup>	X <sup>19</sup>
<i><b>Representation Biases</b></i>	X <sup>1,20</sup>	X <sup>21</sup>
<i><b>Measurement Biases</b></i> <span style="color: red;">TODO: only mester 2022 relevant</span>	X <sup>1,18</sup>	X <sup>19,21</sup>
<i><b>Research Biases</b></i> <span style="color: red;">TODO: specify which one of note 18, was funding bias</span>	X <sup>18</sup>	X <sup>19</sup>
<i><b>Feature Representation Biases</b></i> <span style="color: red;">TODO: only mester 2022 relevant</span>	X <sup>1,18</sup>	X <sup>19</sup>
<i><b>Imaging Biases</b></i>		X <sup>20</sup>
<i><b>Medical Biases</b></i>	X <sup>14</sup>	X <sup>19</sup>
<i><b>Temporal Data Biases</b></i>	X <sup>1</sup>	X <sup>19</sup>
<b>Algorithmic Biases</b>		
<i><b>Biases in Model Training and Learning</b></i>		
Algorithmic Bias	X <sup>1,4,5</sup>	
Emergent Bias	X <sup>1,9</sup>	
Evaluation Bias	X <sup>1,2,12</sup>	
Generalization Issues		X <sup>20</sup> , <span style="color: red;">TODO: add those from young</span>
<i><b>Biases in Predictions and User Interaction</b></i>		
User Interaction Bias	X <sup>1,4</sup>	
Presentation Bias	X <sup>1,4</sup>	
Ranking Bias	X <sup>1,4,6</sup>	
Popularity Bias	X <sup>1,10</sup>	X <sup>19,c9</sup>
Incorporation Bias		X <sup>19,c25,c26</sup>
<b>User Biases</b>		
<i><b>Cognitive and Decision-Making Biases</b></i>		
Cognitive Bias	X <sup>18</sup>	
Previous opinion Bias		X <sup>19, c32</sup>
Confirmation bias		X <sup>19,c15</sup>
Cause-Effect Bias	X <sup>18</sup>	
<i><b>Behavioral and Social Biases</b></i>		
Self-Selection Bias	X <sup>1,18,20</sup>	
Social Bias	X <sup>1,4,7</sup>	
Historical Bias	X <sup>1,2,17</sup>	
Behavioral Bias	X <sup>1,3</sup>	
Temporal Bias	X <sup>1,3</sup>	
<i><b>Medical and Publication Biases</b></i>		
Content Production Bias	X <sup>1,3</sup>	
Publication/All is Well Bias	<span style="color: red;">TODO: add citation</span>	X <sup>19,c10-c12</sup>
Non-Responder bias / Attrition bias		X <sup>19, c9</sup>
<i><b>Population and Perception Biases</b></i>		
Population Bias	X <sup>1,3,8</sup>	
<i><b>Other Biases</b></i>		
Apprehension bias		X <sup>19,c13</sup>
Rethoric bias		X <sup>19, c14</sup>
Centripetal bias		X <sup>19</sup>
Novelty bias		X <sup>19</sup>
Language Bias		X <sup>19</sup>

## Data Biases

**Sampling Biases** When gathering data, it's usually not possible to gather the data of a whole population. Instead, the data is gathered by sampling. A sample is a subgroup of individuals from the population. To get unbiased results, this sampling process should represent the true population, with a low sampling error (HP\_2022). This is often achieved with randomized samples. With non-random sampling processes, sampling bias arises. The consequence is, that the insights of one sampled population may not generalize with insights on another sampled population (Mehrabi et al., 2021).

Those biases can be introduced with a flawed sampling process:

- **Sampling bias**, due to nonrandom sampling of subgroups, leading to poor generalization (Mehrabi et al., 2021)
- **Selection bias**, working only on specific subset of the population which is not representative (Mestner\_2022Chakraborty, 2024)
- **Systematic selection bias**, chosen samples differ dramatically from the representative populations; e.g. in dermatology, when only the most severe patient data gets included (c5; c6; c33; Chakraborty, 2024)
- **Ascertainment bias**, tendency to exclude segments from the population due to e.g. cultural differences, such as which patient segment goes to government clinics vs. private clinics (usually influenced by socioeconomic status) (c5; Chakraborty, 2024) TODO: check this Mehrabi reference also M98 -> potentially remove ref from table
- **Availability bias**, focus on widely available data instead of most representative data (c9; c10; Chakraborty, 2024<empty citation>)
- **Survivorship bias**, focus only on pre-selected data, ignoring the initial data-points which got filtered out (Mestner\_2022).

**Potential Biases in PASSION** PASSION tries to reduce sampling bias in dermatology against high pigmented skin. PASSION might introduce (systematic) selection bias or Ascertainment bias, if in the dermatology centers only sickest / more severe patients are seen as indicated by Chakraborty, 2024 PASSION inherits availability bias as it is using FST scale. Survivorship bias could be relevant for PASSION, if dermatology diseases could be lethal. Further, all patients which are not able to go to one of the dermatology centers which were used in PASSION could be considered to left out by survivorship bias.

used

- **Sampling Bias**. Sampling bias is similar to representation bias, and it arises due to nonrandom sampling of subgroups. As a consequence of sampling bias, the trends estimated for one population may not generalize to data collected from a new population. (Mehrabi et al., 2021). This is what the PASSION dataset tries to improve

- Selection bias - wrong sampling method, working on a specific subset of audience; usually by working only with data that is easy to access (**Mester\_\_2022Mester\_\_2017**) - statistical bias
- Selection bias: Since it is not possible to work with large populations, for most dermatological studies, samples are chosen that are said to be representative of the original population. In selection bias, the selected subgroups are not representative of their original population. A variation of this is systematic selection bias, where samples chosen differ dramatically from their representative populations. Our experience suggests, such selection bias occurs more commonly in studies conducted in regional referral centers where only the sickest or more severe patients are usually seen. For example, a study compared the efficacy of thalidomide vs. prednisolone in hospitalised patients of erythema nodosum leprosum. It derived that thalidomide was more efficacious than steroids in erythema nodosum leprosum. Such findings cannot be generalised to all erythema nodosum leprosum since patients admitted to a regional referral center will likely have more severe disease.<sup>5,6,33</sup> (Chakraborty, 2024)
- Availability bias: More emphasis is placed on widely available data than scantily available data. A classic example is the use of antihistamines in pregnancy dermatoses, where nearly all standard books recommend first-generation antihistamine chlorpheniramine because more data is available.<sup>9,10</sup> (Chakraborty, 2024) - dermatology
- Survivorship bias (**Mester\_\_2022Mester\_\_2017**) - statistical bias
- Ascertainment Bias: This bias is commonly encountered in venereology practice. It is defined as a bias due to the tendency of some segments of the target population to get excluded due to cultural and other differences. For example, in most venereology clinics in government setups, studies show that venereal diseases are commoner in lower socioeconomic status. One reason might be that the higher socioeconomic status people tend to go to private practitioners and thereby get excluded from government-run clinics.<sup>9,10</sup> Allocation concealment and blinding are good ways to avoid this. <sup>5</sup>. (Chakraborty, 2024) - healthcare

even more extensive

- Selection bias is again divided into two types endogenous selection bias and exogenous selection bias. The best example of endogenous selection bias in dermatology is the inclusion of non-response. If a trial tests the efficacy of a particular biologic in psoriasis, the response is usually collected from trial participants via postal services. Certain participants will not respond, although they might have substantially improved. Their exclusion will result in significant differences in efficacy evaluation.<sup>33</sup> Exogenous selection bias results when both treatment and outcome result from dependency on

an external variable that is not controlled. For example, if sunlight exposure is not controlled, it will influence both the intervention and control groups since psoriasis is a photosensitive (and photoexacerbated) dermatosis. (Chakraborty, 2024) - dermatology

- survivorship bias - World War II planes (**Silfwer\_2017**) - <https://doctorspin.org/media-psychology/psychology/survivorship-bias/>

### Representation Biases TODO: still describe this category

Those biases can be introduced :

- **Representation bias**, non-representative sample lead to missing subgroups or other representation anomalies, which can be harmful to downstream applications. Popular ML datasets suffer from representation bias (Mehrabi et al., 2021**M142**\_\_)
- **Aggregation bias** occurs, when "false conclusions are drawn about individuals from observing the entire population". It doesn't matter, whether the subgroups are represented equally in the training set, any generalized assumptions can result in aggregation bias (Mehrabi et al., 2021). In medicine, diseases can present themselves differently across genders and ethnicities (Suresh & Guttag, 2021). Therefore, diagnostic models need to incorporate those differences to mitigate aggregation bias (Mehrabi et al., 2021).
- **Simpson's Paradox** is a type of aggregation bias, which arises in heterogeneous data analysis. Observed associations disappear or reverses in the subgroup data (Mehrabi et al., 2021).

**Potential Biases in PASSION** PASSION tries to mitigate representation bias, by including more FST skin types - however, it could introduce other representation biases Aggregation bias and Simpson's Paradox could potentially be an issue when the analyzed skin diseases present themselves differently in patients based on their genetics

used

- **Representation Bias.** Representation bias arises from how we sample from a population during data collection process (Suresh & Guttag, 2021). Non-representative samples lack the diversity of the population, with missing subgroups and other anomalies (Mehrabi et al., 2021).
- Popular machine-learning datasets that serve as a base for most of the developed algorithms and tools can also be biased—which can be harmful to the downstream applications that are based on these datasets. ... In [142], researchers showed that these datasets suffer from representation bias and advocate for the need to incorporate geographic diversity and inclusion while creating such datasets. (Mehrabi et al., 2021)

- **Aggregation Bias.** Aggregation bias (or ecological fallacy) arises when false conclusions are drawn about individuals from observing the entire population. An example of this type of bias can be seen in clinical aid tools. Consider diabetes patients who have apparent morbidity differences across ethnicities and genders. Specifically, HbA1c levels, that are widely used to diagnose and monitor diabetes, differ in complex ways across genders and ethnicities. Therefore, a model that ignores individual differences will likely not be well-suited for all ethnic and gender groups in the population (Suresh & Guttag, 2021). This is true even when they are represented equally in the training data. Any general assumptions about subgroups within the population can result in aggregation bias. (Mehrabi et al., 2021). -> could also be important for dermatology issues!!!
  - **Simpson’s Paradox.** Simpson’s paradox is a type of aggregation bias that arises in the analysis of heterogeneous data [18]. The paradox arises when an association observed in aggregated data disappears or reverses when the same data is disaggregated into its underlying subgroups (Fig. 2(a)). ... After analyzing graduate school admissions data, it seemed like there was bias toward women, a smaller fraction of whom were being admitted to graduate programs compared to their male counterparts. However, when admissions data was separated and analyzed over the departments, women applicants had equality and in some cases even a small advantage over men. The paradox happened as women tended to apply to departments with lower admission rates for both genders. Simpson’s paradox has been observed in a variety of domains, including biology [37], psychology [81], astronomy [109], and computational social science [91].(Mehrabi et al., 2021).

**Measurement Biases** How features are chosen, used and measured can lead to biases (Mehrabi et al., 2021; Suresh & Guttag, 2021).

Examples for such biases are:

- **Measurement bias** in general, e.g. using mismeasured proxy variables lead to misinterpretations of the outcome (Mehrabi et al., 2021)
- **Observer bias** is a subconscious bias which can occur in different forms. Either, researchers projects their own expectations on the research and influence the testers accordingly (**Mester\_2022**). In other cases, different observers report the same observation differently (**c29; c26; Chakraborty, 2024**)
- **Annotator bias** is a special form of observer bias. The labeling process of human annotators can be influenced by lots of factors (e.g. personal background, social context) and even minor design choices (e.g. scale order, image context). This can introduce inconsistencies when labeling the data (Montoya et al., 2025)



- **Recall bias.** This bias occurs when queried individuals do not remember things correctly, due to humans selective memory. This can cause misinterpretation, for example when analyzing causes and effects of behaviour on certain diseases in medicine (**Mester\_2022c3-6; c2**; Chakraborty, 2024).

**Potential Biases in PASSION** Measurement Bias (proxy var) - Country of Origin in PASSION depending on the interpretation - should not be used for ethnicity, as this is not linked directly to the genes, see example <https://medium.com/bcggamma/practical-ai-responsibly-with-proxy-variable-detection-42c2156ad986>

Annotator bias regarding skin tone labeling has been investigated in (Montoya et al., 2025). PASSION should evaluate its process.

used

- **Measurement Bias.** Measurement, or reporting, bias arises from how we choose, utilize, and measure particular features (Suresh & Guttag, 2021) (e.g. mismeasured proxy variables) (Mehrabi et al., 2021). (= e.g. someone who lives at that postal code probably has this ethnicity ); → could that be an issue with the country of origin feature?
- This study found that while using skin tone instead of race for fairness evaluations in computer vision seems objective, the annotation process remains biased by human annotators. Untested scales, unclear procedures, and a lack of awareness about annotator backgrounds and social context significantly influence skin tone labeling. This study exposes how even minor design choices in the annotation process, like scale order (dark to light instead of light to dark) or image context (face or no face, skin lesion presence), can sway agreement and introduce uncertainty in skin tone assessments. ... The researchers emphasize the need for greater transparency, standardized procedures, and careful consideration of annotator biases to mitigate these challenges and ensure fairer and more robust evaluations in computer vision. (Montoya et al., 2025) - demographic dermatology bias
- **Observer bias** - projecting expectations onto the research (**Mester\_2022Mester\_2017**) - statistical bias
- **Observer bias:** When different observers view the same observation, they report it differently e.g., different observers may give differing descriptions about subtle features in the histopathology report of a skin biopsy.<sup>29 26</sup>. (Chakraborty, 2024) - dermatology
- **Recall bias** - respondent doesn't remember things correctly; Recall bias is another common error of interview/survey situations. It happens when the respondent doesn't remember things correctly. It's not about bad or good memory – humans have selective memory by default. After a few years (or even a few days), certain things stay and others fade. It's normal, but it makes research much more difficult. **TODO: keep an eye on this when recalling evidences!!** (**Mester\_2022Mester\_2017**) - statistical bias

- **Memory or recall bias:** This is a type of bias where sufferers of a disease, often termed cases, have a greater tendency to recall a particular habit than non-sufferers, viz controls. This results in an uneven distribution of risk factors between the cases and controls. An example of this would be a case-control study to evaluate the association between dental amalgam use and the development of oral lichen planus. Those with lichen planus are more likely to recall a history of dental amalgam use than those who do not have the disease. This difference in recall between a diseased cohort and control has resulted in difficulties in assessing the association between diet and many dermatological diseases – like milk and chocolate consumption and acne, fatty meals and psoriasis, sugary meals and psoriasis, agricultural exposure to insecticides and pemphigus and so on.<sup>3–6</sup> 2. (Chakraborty, 2024) - dermatology

**Research Biases** TODO: consider to move at beginning / out of data biases  
 Researchers and their processes can also be biased in multiple ways:

- **Funding / Sponsorship bias**, when a study is deliberately supporting those findings, which the sponsor expects (c22; Chakraborty, 2024Mester\_2017)
- **Data dredging bias.** The statistical methods and model are chosen to provide a certain p-value, to improve the probability of the research hypothesis being true. TODO: consider to move this to an own reporting section (Chakraborty, 2024)
- **Hot stuff bias** is a subtype of publication bias, where Journals are less critical about trending topics, which lead to more frequent publishing of those topics. This in turn can lead to flawed meta-analyses regarding those topics (c22; c23; c19; Chakraborty, 2024).
- **Hypothetical bias.** Hypothetical questions lead to responses that do not reflect, what interviewees would do in real life. (c31; c28; Chakraborty, 2024) TODO: isn't this a user bias instead?

**Potential Biases in PASSION** Since the PASSION dataset is already published, the research biases might already be introduced. It is not feasible during the duration of this thesis to make an evaluation on those biases. Instead, I would recommend the PASSION team and researcher in general, to check the list above carefully and take measures against them. Maybe, an external evaluation could help to detect and prevent those biases even better.

used

- **Funding bias** (Mester\_2022Mester\_2017) - statistical bias
- **Industry sponsorship bias:** This has now been reclassified as conflict-of-interest bias. In short, the study deliberately supports the findings expected from it by its sponsors. 22.(Chakraborty, 2024) - dermatology

Reporting biases

- **Data dredging bias:** It is an entirely avoidable bias. This is subdivided into two types – Fishing type and “P-value hacking” type. It involves using multiple statistical methods to get the desired p-value and selecting the statistical model that gives the p-value the author wants. This is “lamentably common” in dermatological research.<sup>16</sup> To detect data dredging bias, always perform a “p-curve analysis” while performing a meta-analysis.<sup>17,18</sup> Much emphasis is nowadays given to the confidence interval instead of the p-value, which gives an approximate idea of the range in which one can be 95% (or 90%, depending on the confidence interval chosen) sure that the result is correct. The confidence interval remains unaffected by p-value dredging. This subject has been reviewed in depth in recent works.<sup>18,19</sup> 15.(Chakraborty, 2024)

publication bias

- **Hot stuff bias:** Editors of journals may be less critical about topics that are “fashionable” or currently in vogue and consequently end up publishing them more frequently, resulting in publication bias as well as hot stuff bias. It can result in flawed meta-analyses based on these studies. An example is how cutaneous manifestations of COVID-19 were published. Indian Journal of Dermatology Venereology and Leprosy stood out by choosing not to publish anything and everything related to COVID-19, thus reducing hot stuff bias.<sup>22,23</sup> 19. (Chakraborty, 2024)
- **Hypothetical bias:** Many dermatological researches (and some life quality questionnaires like vitiQoL) use hypothetical questions – like “What would you do when some stranger asks you about your lesion?”. The responses to these questions by the study participants often do not tally with what they would do in real life. This is called hypothetical bias and is avoided by adopting the ex-ante approach.<sup>31</sup> 28. (Chakraborty, 2024) - dermatology

**Feature Representation Biases** Some of those biases are:

- **Omitted Variable Bias** arises when variables are not included in the model, which leads to situations for which the model is not ready for (Mehrabi et al., 2021Mester\_\_2022)(Clarke, 2005; Riegg, 2008)(Mustard, 2003).
- **Collider Bias** Two variables can influence a common third variable, the collider variable. When sampling is restricted by this collider variable, it could lead to a distortion (c4; c8; c9; Chakraborty, 2024).

**Potential Biases in PASSION** The ethnicity is omitted in the PASSION dataset which could lead to issues See the medical section for more specific collider bias, maybe there could be others

used

- **Omitted Variable Bias.** Omitted variable bias<sup>4</sup> occurs when one or more important variables are left out of the model (Clarke, 2005; Riegg, 2008)(Mustard, 2003). Something that the model was not ready for(Mehrabi et al., 2021). did not take into account (Mehrabi et al., 2021)
- **Omitted variable bias (Mester\_2022Mester\_2017)** - statistical bias
- **Collider Bias:** This is an under-appreciated bias, and often confused with a confounder. This is especially seen in observational studies where it is defined as a distortion produced by the restriction of sampling by a collider variable. A collider variable is defined as one that has an independent effect on the outcome studied apart from the studied variable. In simpler terms, collider bias occurs when exposure and development influence a common third variable. That variable or collider is controlled by study design or in the analysis. An example is the observation that psoriasis patients tend to have more depression and anxiety disorders. Since severe psoriasis patients tend to get hospitalised and also get screened for mental health issues, a spurious association between them could have been obtained due to collider bias. The two variables viz psoriasis and depression converged, i.e., collided, into a single outcome – hospitalization.<sup>8,9</sup> 4. (Chakraborty, 2024) - dermatology

**Imaging Biases** Dealing with images can lead to a whole other set of challenges, which can lead to biases. The challenges are for example technical variations in hardware and software but also differences in how images are gathered or what is in it (Young et al., 2020).

Those biases can be introduced :

- **Image Quality Bias.** The quality of an image (zoom level, focus, lightning) could be associated with the classification (Young et al., 2020)
- **Visual Artifact Bias.** Other artifacts, such as presence of hair or surgical ink markings on dermatology images, can decrease classification performance (Winkler et al.; 2019 & Bisla et al.; 2019 (from Young\_2020))
- **Field of View Bias.** What view is captured in the image can interfere with prediction quality what is it, consequence (Mishra et al.; 2019 from Young\_2020)

**Potential Biases in PASSION** The PASSION model could learn to associate unrelated visual effects, hair, body parts or image quality with a disease.  
used

- **Image quality.** Several barriers to AI implementation in the clinic need to be overcome with regards to imaging (Figure 1). These include technical variations (e.g., camera hardware and software) and differences in image acquisition and quality (e.g., zoom level, focus, lighting, and presence of hair). For example, the presence of surgical ink markings is associated with

decreased specificity (Winkler et al., 2019), field of view can significantly affect prediction quality (Mishra et al., 2019), and classification performance improves when hair and rulers are removed (Bisla et al., 2019). We have developed a method to measure how model predictions might be biased by the presence of a visual artifact (e.g., ink) and proposed methods to reduce such biases (Pfau et al., 2019). Poor quality images are often excluded from studies, but the problem of what makes an image adequate is not well studied. Ideally, models need to be able to express a level of confidence in a prediction as a function of image quality and appropriately direct a user to retake photos if needed. (Young et al., 2020) - dermatology

**Medical Biases** In ML for health care, there are special medical versions of the mentioned biases as well as completely new biases. They require special attention, since they directly influence the diagnosis or treatment of a disease.

Those biases can be introduced:

- **Berksonian bias** occurs in hospital-based studies when two variables influence hospital or clinical attendance independently. This can lead to a distorted estimation of the relationship between those variables because the study population of hospitalized patients is not representative of the whole population (**c3**; **c7**; Chakraborty, 2024)
- **Informed presence bias**, the probability to get screened for other diseases is higher for people who seek medical care. Like Berksonian bias, this can lead to misleading interpretations of relationships between two diseases (**c27**; **c23**; Chakraborty, 2024)
- **Diagnostic access bias**, depending on the geographical location, individuals have better access to medical care. Therefore, their disease prevalence could appear to be higher and diseases could be diagnosed earlier. (**c19-c21**; Chakraborty, 2024)
- **Diagnostic reference test bias** is a **verification bias**, where not all individuals receive the same reference test for the diagnostic process, potentially leading to different diagnoses. (**c21**; Chakraborty, 2024)
- **Mimicry bias**, exposures to treatment options can cause a disease which presents itself similar to the study disease, which potentially creates misleading data (**c28**; **c25**; Chakraborty, 2024)
- **Unacceptable Disease bias**. When a disease is socially unacceptable, it can result in under-reporting of the same disease (**c30**; **c27**; Chakraborty, 2024)
- **Healthy volunteer selection bias**, is a type of self-selection bias where the volunteers are in general healthier than the population due to more interest in health (Delgado-Rodríguez\_2004)

**Potential Biases in PASSION** Berksonian bias depending on the chosen hospitals Informed presence bias regarding correlation between impetigo and the other diseases Diagnostic access bias can somewhat be addressed by PASSION, since its dataset includes samples of later states of diseases. However, in the PASSION context itself, this bias could still be relevant. Diagnostic reference test bias could be inherited in the PASSION dataset, depending on how the dermatologists work. Mimicry bias is not relevant regarding the exposures since PASSION does not hold any exposure data. However, diseases which mimicry others could lead to issues if they are not detected.

used

- **Berksonian Bias:** Named after Dr. Joseph Berkson, this bias reflects the variation in rates of hospital admission or clinic attendance for different diseases. For example, if a study is conducted to know the effect of pregnancy on syphilis in an antenatal clinic, we are likely to get biased data since the two conditions, viz pregnancy and syphilis, are both likely to affect clinic attendance and all observations related to the relationship between pregnancy and syphilis.<sup>7 3</sup>. (Chakraborty, 2024) - dermatology
- **Informed presence bias:** Simply, a person attending a health center is more likely to get screened for other unrelated comorbidities than those not attending a health center e.g., the finding psoriasis is associated with depression has now been criticised because those having psoriasis also have a greater chance to be screened for depression since they are already attending a health center.<sup>27 23</sup>. (Chakraborty, 2024) - dermatology
- **Diagnostic Access Bias:** Individuals in certain geographical localities have better access to medical care and, hence, may appear to have higher disease prevalence. For example, atopic dermatitis is believed to be commoner in the West – this could be due to better and earlier diagnostic facilities available than in India.<sup>19,20 17</sup>.(Chakraborty, 2024)
- **Diagnostic reference test bias:** These bias results when all individuals do not receive the same reference test. e.g., direct immunofluorescence studies may not be done for all patients with pemphigus vulgaris some patients may receive only a skin biopsy-based diagnosis. It is a subtype of verification bias. Another variation of this type of bias is partial reference bias, where only some of the study participants receive the index and the reference tests.<sup>21</sup>(Chakraborty, 2024)
- **Mimicry bias:** When an exposure causes a disease that resembles the study disease, mimicry bias can result. For example, certain drugs are known to cause a pityriasis rosea-like reaction, which, although looks like pityriasis rosea, differs from it.<sup>28 25</sup>.(Chakraborty, 2024) - dermatology
- **Unacceptable disease bias:** This occurs in socially unacceptable diseases like leprosy and STDs, which result in under-reporting.<sup>30 27</sup>. (Chakraborty, 2024) - dermatology



- **TODO:** Other such studies were conducted in [54] which states that UK Biobank, a large and widely used genetic dataset, may not represent the sampling population. Researchers found evidence of a “healthy volunteer” selection bias. [150] has other examples of studies on existing biases in the data used in the medical domain. [157] also looks at machine-learning algorithms and data utilized in medical fields, and writes about how artificial intelligence in health care has not impacted all patients equally. (Mehrabi et al., 2021) → [150] also provides an overview over the impact of social determinants on health, such as Economic stability, neighborhood and physical environment, education, food, community and social context, access to healthcare and quality
- The healthy volunteer effect is a particular case: when the participants are healthier than the general population. (**Delgado-Rodríguez\_2004**)

**Temporal Data Biases** Certain studies require to track temporal data, to learn about behaviour or disease progression over time (Mehrabi et al., 2021). For PASSION, temporal biases are currently irrelevant, since PASSION contains images independently of time and is not tracking the disease progression. Therefore, the listed biases in this chapter are not explained in detail, refer to the sources for further information.

Examples for temporal data biases are:

- **Longitudinal Data Fallacy** (Mehrabi et al., 2021)
- **Chronological bias** (c9; c13; Chakraborty, 2024)
- **Immortal time bias** (c24; c20; Chakraborty, 2024)

## 2.2.4 Sensitive Features

Research showed sensitive features which are prone to bias. Those led to issues in existing AI applications and their usage should therefore be investigated carefully.

2.4 summarizes mentioned features. Since PASSION tries to improve the classification of skin diseases with photographs alone, those features which directly influence the appearance of a disease on the skin are most important. Skin type including undertone are affecting the diseases appearance, but also genetic factors, gender and the age **TODO: doublecheck impact of age** of a patient can influence how the diseases manifest themselves. Therefore, they are directly relevant for skin detection and must be considered in the dataset.

Other demographic factors can have an impact on the progression of skin diseases or increase the prevalence of skin conditions (e.g. tropical vs dry climates) **TODO: add source**. While this information could support the diagnostic process, it is not visible in the images. Therefore, they could potentially be relevant and enhance the diagnostic process when included into the PASSION dataset. However, the linkage is weak and could introduce further biases.

For completeness, 2.4 shows further sensitive demographic features which seem not to be linked to dermatology according to research.

**TODO: check what to do with those additional features:** Other important features according to ((Montoya et al., 2025) 13): lesion type, anatomical location of lesion, img characteristics such as source, imaging techniques, resolution, real vs. artificially generated

Bias sensitive Features	Mentioned in Context of	
	ML	Dermatology
<b><i>Dermatology Related Features</i></b>		
Skin Type	X <sup>1,2,7,12</sup>	X <sup>2,13</sup>
Skin Undertones		X <sup>13</sup>
<b><i>Demographic Features with Direct Relevance for Skin Disease Detection</i></b>		
Age	X <sup>7,11</sup>	X <sup>13</sup>
Ethnicity/Race <b>TODO: check definitions</b>	X <sup>1,2,4,5,6,7,11,12</sup>	X <sup>13</sup>
Gender/Sex	X <sup>1,2,7,8,9,10,11</sup>	X <sup>13</sup>
Gender and Skin Type Subgroups	X <sup>1,2</sup>	
<b><i>Demographic Features Potential Relevance for Skin Disease Detection</i></b>		
Geographic Location	X <sup>1,3</sup>	
Socio-Economic Status	X <sup>6,12</sup>	
Disabilities	X <sup>7,11</sup>	
<b><i>Demographic Features no Relevance for Skin Disease Detection</i></b>		
Familial status	X <sup>7</sup>	
Marital status	X <sup>7,11</sup>	
Nationality/National origin	X <sup>7,11</sup>	
Recipient of public assistance	X <sup>7</sup>	
Religion	X <sup>7,11</sup>	
<sup>1</sup> (Mehrabi et al., 2021)	<sup>5</sup> (M54__)	<sup>9</sup> (M20__)
<sup>2</sup> (Buolamwini & Gebru, 2018)	<sup>6</sup> (M150__)	<sup>10</sup> (M168__)
<sup>3</sup> (M142__)	<sup>7</sup> (M30__)	<sup>11</sup> (M62__)
<sup>4</sup> (M98__)	<sup>8</sup> (M167__)	<sup>12</sup> (Young et al., 2020)
		<sup>13</sup> (Montoya et al., 2025)

Table 2.4: Features which often hold biases - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness



## 2.3 Fairness

## 2.4 Mitigation Methods

## 2.5 Extensive Sources

### 2.5.1 Discrimination vs. Biases

- bias and discrimination = source of unfairness. Discrimination can be considered as a source for unfairness that is due to human prejudice and stereotyping based on the sensitive attributes, which may happen intentionally or unintentionally, while bias can be considered as a source for unfairness that is due to the data collection, sampling, and measurement. Although bias can also be seen as a source of unfairness that is due to human prejudice and stereotyping, in the algorithmic fairness literature it is more intuitive to categorize them as such according to the existing research in these areas. In this survey, we mainly focus on concepts that are relevant to algorithmic fairness issues. (Mehrabi et al., 2021)
- Explainable Discrimination. Differences in treatment and outcomes amongst different groups can be justified and explained via some attributes in some cases. In situations where these differences are justified and explained, it is not considered to be illegal discrimination and hence called explainable [77]. In [77], authors present a methodology to quantify the explainable and illegal discrimination in data. They argue that methods that do not take the explainable part of the discrimination into account may result in non-desirable outcomes, so they introduce a reverse discrimination which is equally harmful and undesirable. They explain how to quantify and measure discrimination in data or a classifier’s decisions which directly considers illegal and explainable discrimination.(Mehrabi et al., 2021)
- Unexplainable Discrimination. In contrast to explainable discrimination, there is unexplainable discrimination in which the discrimination toward a group is unjustified and therefore considered illegal. Authors in [77] also present local techniques for removing only the illegal or unexplainable discrimination, allowing only for explainable differences in decisions. These are preprocessing techniques that change the training data such that it contains no unexplainable discrimination. We expect classifiers trained on this preprocessed data to not capture illegal or unexplainable discrimination. Unexplainable discrimination consists of direct and indirect discrimination.(Mehrabi et al., 2021)
  - Direct Discrimination. Direct discrimination happens when protected attributes of individuals explicitly result in non-favorable outcomes toward them [164]. ... these traits that are considered to be “protected”

- or “sensitive” attributes in computer science literature (Mehrabi et al., 2021)
- Indirect Discrimination. In indirect discrimination, individuals appear to be treated based on seemingly neutral and non-protected attributes; however, protected groups, or individuals still get to be treated unjustly as a result of implicit effects from their protected attributes (Mehrabi et al., 2021)
- Discrimination can be either direct or indirect. Direct discrimination occurs when decisions are made based on sensitive attributes. Indirect discrimination occurs when decisions are made based on nonsensitive attributes which are strongly correlated with biased sensitive ones (**M62\_\_**)
- ML algorithms reflect cognitive biases in humans, which can result in predictions and decisions that are “unfair” [5]. Unfairness in ML decision-making equates to prejudice or favoritism based on inherent or acquired characteristics of an individual or group [6]. Within medical research, the term bias refers to “a feature of the design of a study, or the execution of a study, or the analysis of the data from a study, that makes evidence misleading” [7,8]. (**Montonaya\_2025**)

## 2.6 General ML biases

### 2.6.1 Bias Introduction

- compared SAVRY, a tool used in risk assessment frameworks that includes human intervention in its process, with automatic machine learning methods in order to see which one is more accurate and more fair. Conducting these types of studies should be done more frequently, but prior to releasing the tools in order to avoid doing harm (Mehrabi et al., 2021).
- **Assessment Tools** An interesting direction that researchers have taken is introducing tools that can assess the amount of fairness in a tool or system. For example, Aequitas [136] is a toolkit that lets users to test models with regards to several bias and fairness metrics for different population subgroups. Aequitas produces reports from the obtained data that helps data scientists, machine learning researchers, and policymakers to make conscious decisions and avoid harm and damage toward certain populations. AI Fairness 360 (AIF360) is another toolkit developed by IBM in order to help moving fairness research algorithms into an industrial setting and to create a benchmark for fairness algorithms to get evaluated and an environment for fairness researchers to share their ideas [11]. These types of toolkits can be helpful for learners, researchers, and people working in the industry to move towards developing fair machine learning application away from discriminatory behavior (Mehrabi et al., 2021).

- At first sight, automating decisions may give a sense of fairness: classification rules do not guide themselves by personal preferences. However, at a closer look, one realizes that classification rules are actually learned by the system (e.g., loan granting) from the training data. If the training data are inherently biased for or against a particular community (e.g., foreigners), the learned model may show a discriminatory prejudiced behavior. In other words, the system may infer that just being foreign is a legitimate reason for loan denial. (M62\_\_)
- One might think of a straightforward preprocessing approach consisting of just removing the discriminatory attributes from the data set. Although this would solve the direct discrimination problem, it would cause much information loss and in general it would not solve indirect discrimination. (M62\_\_)
- Hence, there are two important challenges regarding discrimination prevention: one challenge is to consider both direct and indirect discrimination instead of only direct discrimination; the other challenge is to find a good trade-off between discrimination removal and the quality of the resulting training data sets and data mining models. (M62\_\_)
- Most AI systems and algorithms are data driven and require data upon which to be trained. Thus, data is tightly coupled to the functionality of these algorithms and systems. In the cases where the underlying training data contains biases, the algorithms trained on them will learn these biases and reflect them into their predictions. As a result, existing biases in data can affect the algorithms using the data, producing biased outcomes. Algorithms can even amplify and perpetuate existing biases in the data. (Mehrabi et al., 2021).
- In addition, algorithms themselves can display biased behavior due to certain design choices, even if the data itself is not biased. The outcomes of these biased algorithms can then be fed into real-world systems and affect users' decisions, which will result in more biased data for training future algorithms. (Mehrabi et al., 2021).
- Bias can exist in many shapes and forms, some of which can lead to unfairness in different downstream learning tasks. In (Suresh & Gutttag, 2021), authors talk about sources of bias in machine learning with their categorizations and descriptions in order to motivate future solutions to each of the sources of bias introduced in the paper. In (Olteanu et al., 2019), the authors prepare a complete list of different types of biases with their corresponding definitions that exist in different cycles from data origins to its collection and its processing. (Mehrabi et al., 2021).
- The list of biases that can occur in any research is considerably long, and certainly not all of them can be avoided. However, dermatologists should be well aware of them because: 1. While conducting systematic reviews and

metaanalysis, PRISMA guidelines need to be followed, and the PRISMA checklist requires a very exhaustive list of declarations to be made by the authors as to how biases in the individual studies were detected, their types and whether they were included in the systematic review or not. This usually requires more than two authors working independently. 2. Biases can result in dramatically opposite inferences, which may not be biologically plausible; the knowledge of the biases can help detect them and thereby negate such findings. 3. The knowledge of biases is a vital part of the postgraduate dermatology curriculum and is a mustknow area for thesis/dissertation purposes. 4. Since there is no fool-proof way to avoid all biases, the help of a well-qualified biostatistician might help detect and prevent many of these biases in research. (Chakraborty, 2024)

Already rewritten: The following categorization was modeled with the intent to show that the different biases are intertwined and one should consider the effects between each other in the cycle to address them correctly (Mehrabi et al., 2021)

## 2.6.2 Biases Extensive Sources

**Data Biases** Data biases (data to algorithm (biases in data which might have an impact on biased algorithmic outcomes (Mehrabi et al., 2021)))

•

**Algorithmic Biases** Algorithmic biases (Algorithm to user (A modulates U behaviour, biases in algorithm might lead to introduce biases in user behaviour and affect it as a consequence)) (Mehrabi et al., 2021)

- **Algorithmic Bias.** Algorithmic bias is when the bias is not present in the input data and is added purely by the algorithm (Baeza-Yates, 2018). The algorithmic design choices, such as use of certain optimization functions, regularizations, choices in applying regression models on the data as a whole or considering subgroups, and the general use of statistically biased estimators in algorithms (Danks & London, 2017), can all contribute to biased algorithmic decisions that can bias the outcome of the algorithms.(Mehrabi et al., 2021).
- **User Interaction Bias.** User Interaction bias is a type of bias that can not only be observant on the Web but also get triggered from two sources—the user interface and through the user itself by imposing his/her self-selected biased behavior and interaction (Baeza-Yates, 2018). This type of bias can be influenced by other types and subtypes, such as presentation and ranking biases. (Mehrabi et al., 2021). – more relevant for later, when the application would become bigger
  - **Presentation Bias.** Presentation bias is a result of how information is presented (Baeza-Yates, 2018) (can only click on content they see,

could be the case that user does not see all info on web) (Mehrabi et al., 2021).

- Ranking Bias. The idea that top-ranked results are the most relevant and important will result in attraction of more clicks than others. This bias affects search engines (Baeza-Yates, 2018) and crowdsourcing applications (Lerman & Hogg, 2014). (Mehrabi et al., 2021).
- Popularity Bias. Items that are more popular tend to be exposed more. However, popularity metrics are subject to manipulation—for example, by fake reviews or social bots (Ciampaglia et al., 2018). ... this presentation may not be a result of good quality; instead, it may be due to other biased factors. (Mehrabi et al., 2021).
- Emergent Bias. Emergent bias occurs as a result of use and interaction with real users. This bias arises as a result of change in population, cultural values, or societal knowledge usually some time after the completion of design (Friedman & Nissenbaum, 1996). This type of bias is more likely to be observed in user interfaces, ... This type of bias can itself be divided into more subtypes, as discussed in detail in (Friedman & Nissenbaum, 1996). (Mehrabi et al., 2021). probably less relevant at the first stage
- Evaluation Bias. Evaluation bias happens during model evaluation (Suresh & Guttag, 2021). This includes the use of inappropriate and disproportionate benchmarks for evaluation of applications such as Adience and IJB-A benchmarks. These benchmarks are used in the evaluation of facial recognition systems that were biased toward skin color and gender (Buolamwini & Gebru, 2018), and can serve as examples for this type of bias (Suresh & Guttag, 2021). (Mehrabi et al., 2021). – important for this thesis

**User Biases** User to Data (user-generated data, inherent biases in users could be reflected in the data they generate; biases in last section might introduce further bias in this process) (Mehrabi et al., 2021)

- Historical Bias. Historical bias is the already existing bias and socio-technical issues in the world and can seep into from the data generation process even given a perfect sampling and feature selection (Suresh & Guttag, 2021). ... search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering (Mehrabi et al., 2021) - maybe relevant
- Population Bias. Population bias arises when statistics, demographics, representatives, and user characteristics are different in the user population of the platform from the original target population (Olteanu et al., 2019). Population bias creates non-representative data. ... More such examples and statistics related to social media use among young adults according to gender, race, ethnicity, and parental educational background can be found in (Hargittai, 2007). (Mehrabi et al., 2021)

- **Self-Selection Bias.** Self-selection bias<sup>4</sup> is a subtype of the selection or sampling bias in which subjects of the research select themselves. (Mehrabi et al., 2021)
- **Social Bias.** Social bias happens when others' actions affect our judgment (Baeza-Yates, 2018). (case where we want to rate or review an item with a low score, but when influenced by other high ratings, we change our scoring thinking that perhaps we are being too harsh [(Baeza-Yates, 2018), (Wang & Wang, 2014).]) (Mehrabi et al., 2021)
- **Behavioral Bias.** Behavioral bias arises from different user behavior across platforms, contexts, or different datasets (Olteanu et al., 2019). (Mehrabi et al., 2021) maybe, people from different countries go to the dermatologist for different diseases, based on cultural differences?
- **Temporal Bias.** Temporal bias arises from differences in populations and behaviors over time (Olteanu et al., 2019). (Mehrabi et al., 2021) – could this also be differences in the year, when people go to dermatologists? over which timeline has the PASSION data been captured?
- **Content Production Bias.** Content Production bias arises from structural, lexical, semantic, and syntactic differences in the contents generated by users (Olteanu et al., 2019). (Mehrabi et al., 2021) – could the quality of the pictures been related to this as well?

### **Dermatology Biases**

- **Equity.** AI has the potential to worsen health-care disparities, as recognized by the popular media (Khullar, 2019), particularly in dermatology (Adamson and Smith, 2018). The first concern is adequate representation of underserved populations in training data. Existing DL models have been trained on mainly European or East Asian populations, and the relative lack of training on darker skin pigmentation may limit overall diagnostic accuracy. This possibility is demonstrated by the increased error rates in commercial systems, trained on predominantly white datasets, for facial analysis in identifying black individuals (Buolamwini and Gebru, 2018). Second, AI may entrench existing social and economic biases and perpetuate inadvertent discriminatory practices, for example, in recommending less follow-up for black patients than for whites, when health costs are used as a proxy for health needs (Obermeyer et al., 2019). Third, disproportionate adoption by different groups may exacerbate existing inequities. Access to and use of technology differs based on sociodemographics (Tsetsi and Rains, 2017), and more techsavvy users may be more likely to embrace AI for skin screening (Tong and Sopory, 2019). The issue of equity in AI diagnosis needs to be carefully addressed to avoid inadvertent exacerbation of health-care disparities. (Young et al., 2020) - dermatology

- Model generalizability. Generalizability is a major concern for AI models; studies of computer-assisted diagnosis of melanoma report lower sensitivity for melanoma on independent test sets than on nonindependent test sets (Dick et al., 2019). It is difficult to study generalizability because published DL models are not publicly available, making it impossible to compare performance, unless each study uses a standardized benchmark database, such as the Melanoma Classification Benchmark (Brinker et al., 2019d). Han et al. (2018a) reported excellent metrics of performance and made their model available for image submission; however, the model prediction was not robust when images from an outside clinic were submitted, image magnification or contrast was altered, or images were rotated (NavarreteDechent et al., 2018). On ImageNet, a nonmedical dataset of 1,000 object categories, training on a dataset of 300 ... (Young et al., 2020)

### **Papers I would need access to**

- <https://jamanetwork.com/journals/jamadermatology/article-abstract/2688587>  
(linked to (Young et al., 2020))
- <https://academic.oup.com/bjd/article-abstract/190/6/789/7603706?redirectedFrom=fulltext>  
Ethical considerations for artificial intelligence in dermatology: a scoping review



### 2.6.3 Fairness Overview

Fairness Definitions	Mentioned in Context of	
	ML	Dermatology
<b>Group Fairness</b>		
Conditional Statistical Parity	X <sup>1,3,10</sup>	
Demographic/Statistical Parity	X <sup>1,3,4,5</sup>	
Equalized Odds	X <sup>1,2,3</sup>	
Equal Opportunity	X <sup>1,2,3</sup>	
Treatment Equality	X <sup>1,7</sup>	
Test Fairness	X <sup>1,3,8</sup>	
<b>Subgroup Fairness</b>		
Subgroup Fairness	X <sup>1,11,12</sup>	
<b>Individual Fairness</b>		
Counterfactual Fairness	X <sup>1,5</sup>	
Fairness Through Awareness	X <sup>1,4,5</sup>	
Fairness Through Unawareness	X <sup>1,5,6</sup>	
<b>Not Categorized</b>		
Fairness in Relational Domains	X <sup>1,9</sup>	

- <sup>1</sup> (Mehrabi et al., 2021)      <sup>5</sup> (Kusner et al., 2017)      <sup>9</sup> (Farnadi et al., 2018)
- <sup>2</sup> (Hardt et al., 2016)      <sup>6</sup> (Grgic-Hlača et al., 2016)      <sup>10</sup> (Corbett-Davies et al., 2017)
- <sup>3</sup> (Verma & Rubin, 2018)      <sup>7</sup> (Berk et al., 2017)      <sup>11</sup> (Kearns et al., 2018)
- <sup>4</sup> (Dwork et al., 2012)      <sup>8</sup> (Chouldechova, 2017)      <sup>12</sup> (Kearns et al., 2019) **TODO:**  
potential bias in this bc same  
author of the algorithm tested  
it

Table 2.5: Fairness Definitions - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness

### 2.6.4 Fairness Extensive Sources

#### Algorithmic Fairness

- in order to be able to fight against discrimination and achieve fairness, one should first define fairness. (Mehrabi et al., 2021)
- The fact that no universal definition of fairness exists shows the difficulty of solving this problem [138]. Different preferences and outlooks in different cultures lend a preference to different ways of looking at fairness, which makes it harder to come up with just a single definition that is acceptable to everyone in a situation. there is still no clear agreement on which constraints are the most appropriate for those problems. (Mehrabi et al., 2021)
- Broadly, fairness is the absence of any prejudice or favoritism towards an individual or a group based on their intrinsic or acquired traits in the context of



decision-making [139]. Even though fairness is an incredibly desirable quality in society, it can be surprisingly difficult to achieve in practice. (Mehrabi et al., 2021)

- Here we will reiterate and provide some of the most widely used definitions, along with their explanations inspired from (Verma & Rubin, 2018).(Mehrabi et al., 2021)
- Definitions on page 12, 13, 14 (Mehrabi et al., 2021)
  - Equalized Odds (TP and FP rate should be the same for individuals in different sub groups) (Mehrabi et al., 2021)
  - Equal Opportunity (TP rate should be the same) (Mehrabi et al., 2021)
  - Demographic Parity / Statistical Parity (likelihood of positive outcome the same regardless of protected group) (Mehrabi et al., 2021)
  - Fairness Through Awareness (similar predictions to similar individuals (similarity = inverse distance)) (Mehrabi et al., 2021)
  - Fairness Through Unawareness (no protected attributes explicitly used in decision-making process) (Mehrabi et al., 2021)
  - Treatment Equality (Ration FN and FP same for both protected group categories) (Mehrabi et al., 2021)
  - Test Fairness (for predicted probability scores, people in both groups must have equal probability of TP) (Mehrabi et al., 2021)
  - Counterfactual Fairness (same outcome in actual world and counterfactual world where the individual belonged to a different demographic group) (Mehrabi et al., 2021)
  - Fairness in Relational Domains (“A notion of fairness that is able to capture the relational structure in a domain—not only by taking attributes of individuals into consideration but by taking into account the social, organizational, and other connections between individuals” (Farnadi et al., 2018)) (Mehrabi et al., 2021) probably not relevant since not relational
  - Conditional Statistical Parity (people in both groups have equal possibilities of being assigned to a positive outcome given a set of legitimate factors) (Mehrabi et al., 2021)
  - My text: The survey categorizes those fairness notions in three different groups: Individual Fairness, Group Fairness and Subgroup fairness. (Mehrabi et al., 2021)
  - Subgroup fairness: Subgroup fairness intends to obtain the best properties of the group and individual notions of fairness. It is different than these notions but uses them in order to obtain better outcomes. It picks a group fairness constraint like equalizing false positive and asks whether this constraint holds over a large collection of subgroups (Kearns et al., 2018)(Kearns et al., 2019)(Mehrabi et al., 2021)

- it is impossible to satisfy some of the fairness constraints at once except in highly constrained special cases. In [83], the authors show the inherent incompatibility of two conditions: calibration and balancing the positive and negative classes. These cannot be satisfied simultaneously with each other unless under certain constraints; therefore, it is important to take the context and application in which fairness definitions need to be used into consideration and use them accordingly [141](Mehrabi et al., 2021)
- Another important aspect to consider is time and temporal analysis of the impacts that these definitions may have on individuals or groups. In [95] authors show that current fairness definitions are not always helpful and do not promote improvement for sensitive groups—and can actually be harmful when analyzed over time in some cases. They also show that measurement errors can also act in favor of these fairness definitions; therefore, they show how temporal modeling and measurement are important in evaluation of fairness criteria and introduce a new range of trade-offs and challenges toward this direction. It is also important to pay attention to the sources of bias and their types when trying to solve fairness-related questions. (Mehrabi et al., 2021)

### 2.6.5 Mitigation Methods Overview

TODO: write definitions of pre-in and post-processing, see Methods for fair machine learning below [43, 11, 14]

TODO: double check and further improve groups

Mitigation Methods - Unbiasing Data (Pre-Processing)	Mentioned in Context of	
	ML	Dermatology
<b>Documentation and Transparency</b>		
Good Practices while using Data	X <sup>1,2,3</sup>	
Datasheets as supporting document for dataset creation method, characteristics, motivations and skews	X <sup>1,2,3</sup>	
Datasheets as supporting document for model method, characteristics, motivations and skews	X <sup>1,4</sup>	
Dataset (Nutrition) Labels	X <sup>1,5,6</sup>	X <sup>18</sup> , <b>TODO: add spec source</b>
<b>Communication and Reporting</b>		
Messaging	X <sup>1,12</sup>	
<b>Bias Detection and Evaluation</b>		
Test for Simpson's Paradox <b>TODO: Discribe Simpson's Paradox</b>	X <sup>1,7,8,9</sup>	
Detect Direct Discrimination with Causal Models and Graphs	X <sup>1,10</sup>	
Out-of-Distribution Detection in Dermatology Using Input Perturbation and Subset Scanning		X <sup>19</sup>
Check confidence interval and p-curve analysis instead of p-value		X <sup>17</sup>
<b>Study Design</b>		
Allocation concealment and blinding		X <sup>17</sup>
Preventing Direct and Indirect Discrimination	X <sup>1,11</sup>	
<b>Data Gathering</b>		
Data Collection from diverse sources (incl. primary care clinics)	X <sup>18</sup>	
Robust standards for external validation	X <sup>18</sup>	
Preferential Sampling	X <sup>1,13,14</sup>	
Geographical Diversity and Inclusion for Dataset creation	X <sup>16</sup>	
Balanced Representation accross skin tones and genders		X <sup>19</sup>
Disparate Impact Removal	X <sup>1,15</sup>	
<b>Labeling</b>		
Multidimensional Scale for Skin Tones		X <sup>19</sup>
<b>Data Availability and Open Science</b>		
Publish Datasets accessible for the public		X <sup>18</sup> , <b>TODO: add source</b>

<sup>1</sup> (Mehrabi et al., 2021)	<sup>7</sup> (M81__)	<sup>13</sup> (M75__)
<sup>2</sup> (M13__)	<sup>8</sup> (M3__)	<sup>14</sup> (M76__)
<sup>3</sup> (M55__)	<sup>9</sup> (M4__)	<sup>15</sup> (M51__)
<sup>4</sup> (M110__)	<sup>10</sup> (M163__)	<sup>16</sup> (M142__)
<sup>5</sup> (M66__)	<sup>11</sup> (M62__)	<sup>17</sup> (Chakraborty, 2024)
<sup>6</sup> (M66Successor__)	<sup>12</sup> (M74__)	<sup>18</sup> (Young et al., 2020)
		<sup>19</sup> (Montoya et al., 2025)

Table 2.6: Mitigation Methods - Unbiasing Data - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness

Mitigation Methods - Fair Classification	Mentioned in Context of	
	ML	Dermatology
<b>Satisfy Fairness Definitions</b>		
Satisfy Subgroup Fairness <b>TODO: unclear if</b> <b>* in <sup>3</sup> as well, or if <sup>2</sup> also handles *</b>	X <sup>1,2</sup>	
Satisfy Equality of Opportunity <sup>*</sup>	X <sup>1,3,6</sup>	
Satisfy Equalized Odds <sup>*</sup>	X <sup>1,3</sup>	
Disparate Treatment <sup>**</sup>	X <sup>1,4,5</sup>	
Disparate Impact <sup>**</sup>	X <sup>1,4,5</sup>	
<b>TODO: find out exact method</b>	X <sup>1,7</sup>	
<b>TODO: find out exact method</b>	X <sup>1,8</sup>	
<b>TODO: find out exact method</b>	X <sup>1,9</sup>	
<b>TODO: find out exact method</b>	X <sup>1,10</sup>	
<b>Satisfy Fairness and Stability Under Distribution Shifts</b>		
<b>TODO: find out exact method</b>	X <sup>1,11</sup>	
<b>Fair Representation Learning (Pre/In-processing)</b>		
Representation Learning by	X <sup>1,2</sup>	
Disentanglement		
Variational Fair Autoencoder	X <sup>1,3,15</sup>	
VAE without adversarial training	X <sup>1,4</sup>	
Adversarial Learning with FairGAN	X <sup>1,16</sup>	
Removing correlation between protected and unprotected features with a geometric solution	X <sup>1,17</sup>	
<b>Algorithmic Adaptions for Fairness</b>		
Modified Discrimination-Free Naive Bayes Classifier	X <sup>1,12</sup>	
<b>Fairness-Aware ML Frameworks</b>		
Fairness-Aware Classification Framework	X <sup>1,13</sup>	
Fairness Constraints in Multitask Learning (MTL) Framework	X <sup>1,14</sup>	
Decoupled Classification System with Transfer Learning	X <sup>1,15</sup>	
<b>Preferential Data Selection and Representation</b>		
Wasserstein Distance Measure for	X <sup>1,16</sup>	
Dependence Mitigation		
Preferential Sampling (PS) for	X <sup>1,17</sup>	
Discrimination-Free Training Data		
<b>Model Interpretability</b>		
Post-Processing with Attention Mechanism	X <sup>1,18</sup>	
Use Brier Score and Response Rate Accuracy		X <sup>19</sup> , <b>TODO: add clear source</b>
some more methods <b>TODO: describe</b>		X <sup>19</sup>

<sup>*</sup> possible to satisfy together	<sup>6</sup> (M154__)	<sup>13</sup> (M155__)
<sup>**</sup> possible to satisfy together	<sup>7</sup> (M57__)	<sup>14</sup> (M12__)
<sup>1</sup> (Mehrabi et al., 2021)	<sup>8</sup> (M78__)	<sup>15</sup> (M49__)
<sup>2</sup> (M147__)	<sup>9</sup> (M85__)	<sup>16</sup> (M73__)
<sup>3</sup> (Hardt et al., 2016)	<sup>10</sup> (M106__)	<sup>17</sup> (M75__)
<sup>4</sup> (M2__)	<sup>11</sup> (M69__)	<sup>18</sup> (M102__)
<sup>5</sup> (M159__)	<sup>12</sup> (M25__)	<sup>19</sup> (Young et al., 2020)

Table 2.7: Mitigation Methods - Fair Classification - Mentioned in Contextual Research. grouped like in Mehrabi et al., 2021. the author cannot guarantee for

TODO: check categorization

Mitigation Methods - not so relevant for us	Mentioned in Context of	
	ML	Dermatology
<b>Fair NLP</b>		
Fair Word-Embedding	X <sup>1,5,6,7</sup>	
Train-Time Data Augmentation	X <sup>1,8</sup>	
Test-Time Neutralization	X <sup>1,8</sup>	
<b>Fair Regression (In-processing)</b>		
Price of Fairness (POF)	X <sup>1,10</sup>	
XY TODO: check this and bounded group loss	X <sup>1,11</sup>	
Decision Tree for Disparate Impact and Treatment	X <sup>1,12</sup>	
<b>Structured Prediction (In-processing)</b>		
Reducing Bias Amplification (RBA) as calibration algorithm	X <sup>1,13</sup>	
<b>Principal Component Analysis (PCA) (In-processing)</b>		
Fair PCA	X <sup>1,14</sup>	
<b>Graph-Based Fairness Methods</b>		
Community Detection / Graph Embedding	X	
TODO: how to proceed with this		
<b>Causal Fairness and Disparate Learning</b>		
Disparate Learning Processes (DLP)	X <sup>1,9</sup>	
Causal Approach to Fairness TODO: how to proceed with this	X <sup>TODO: add clear source</sup>	
Disregard path in causal graph which result in sensitive attributes affecting decision outcome	X <sup>1</sup>	
<b>Removing Sensitive Attributes</b>		
Disregard sensitive attributes in effect on decision making	X <sup>1</sup>	

<sup>1</sup> (Mehrabi et al., 2021)	<sup>7</sup> (M169__)	<sup>13</sup> (M167__)
<sup>2</sup> (M42__)	<sup>8</sup> (M166__)	<sup>14</sup> (M137__)
<sup>3</sup> (M97__)	<sup>9</sup> (M94__)	<sup>15</sup> (M5__)
<sup>4</sup> (M112__)	<sup>10</sup> (M14__)	<sup>16</sup> (M90__)
<sup>5</sup> (M20__)	<sup>11</sup> (M1__)	<sup>17</sup> (M65__)
<sup>6</sup> (M58__)	<sup>12</sup> (M2__)	

Table 2.8: Mitigation Methods - Others - Mentioned in Contextual Research, grouped like in Mehrabi et al., 2021, the author cannot guarantee for completeness

TODO: mention also the IBM AI Fairness 360 toolkit [11] and that authors evaluated their work in benchmark datasets [65], [72], [158], [159]

### 2.6.6 Mitigation Methods Extensive Sources

**Bias Examples and Mitigation Ideas** Data bias examples and mitigation ideas

- Bias in ML Data - (Buolamwini & Gebru, 2018) IJB-A / Adience imbalanced (mainly light-skinned subjects) - Bias towards dark-skinned groups (under-represented). Other instance - when we do not consider different subgroups in the data. Considering only male-female groups not enough, use race to further subdivide gender groups. Only then, clear biases in sub groups can be found, since otherwise part of the groups would compromise the other group and hide the underlying bias towards that subgroup (Mehrabi et al., 2021)
- Popular machine-learning datasets that serve as a base for most of the developed algorithms and tools can also be biased—which can be harmful to the downstream applications that are based on these datasets. ... In [142], researchers showed that these datasets suffer from representation bias and advocate for the need to incorporate geographic diversity and inclusion while creating such datasets. (Mehrabi et al., 2021)
- Examples of Data Bias in Medical Applications. These data biases can be more dangerous in other sensitive applications. For example, in medical domains there are many instances in which the data studied and used are skewed toward certain populations—which can have dangerous consequences for the underrepresented communities. [98] showed how exclusion of African-Americans resulted in their misclassification in clinical studies, so they became advocates for sequencing the genomes of diverse populations in the data to prevent harm to underrepresented populations (Mehrabi et al., 2021) **TODO: What does sequencing data mean?, is it relevant**

### Methods for Fair Machine Learning

- While this section is largely domain-specific, it can be useful to take a cross-domain view. Generally, methods that target biases in the algorithms fall under three categories (Mehrabi et al., 2021)
- Pre-processing. Pre-processing techniques try to transform the data so that the underlying discrimination is removed [43]. If the algorithm is allowed to modify the training data, then pre-processing can be used [11].(Mehrabi et al., 2021)
- In-processing. In-processing techniques try to modify and change state-of-the-art learning algorithms in order to remove discrimination during the model training process [43]. If it is allowed to change the learning procedure for a machine learning model, then in-processing can be used during the training of a model— either by incorporating changes into the objective function or imposing a constraint [11, 14].(Mehrabi et al., 2021)

- **Post-processing.** Post-processing is performed after training by accessing a holdout set which was not involved during the training of the model [43]. If the algorithm can only treat the learned model as a black box without any ability to modify the training data or learning algorithm, then only post-processing can be used in which the labels assigned by the black-box model initially get reassigned based on a function during the post-processing phase [11, 14].(Mehrabi et al., 2021)
- we concentrate on discrimination prevention based on preprocessing, because the preprocessing approach seems the most flexible one: it does not require changing the standard data mining algorithms, unlike the inprocessing approach, and it allows data publishing (rather than just knowledge publishing), unlike the postprocessing approach. (M62\_\_) -> TODO: this is an important point which we should consider for PASSION, also, some more insight in regards of the different phases can be found in this paper
- From learning fair representations [42, 97, 112] to learning fair word embeddings [20, 58, 169], debiasing methods have been proposed in different AI applications and domains. (Mehrabi et al., 2021) -> seems to refer mostly to NLP domains
- Most of these methods try to avoid unethical interference of sensitive or protected attributes into the decision-making process, while others target exclusion bias by trying to include users from sensitive groups. (Mehrabi et al., 2021)
- However, a recent paper [58] argues against these debiasing techniques and states that many recent works on debiasing word embeddings have been superficial, that those techniques just hide the bias and don't actually remove it. (Mehrabi et al., 2021)
- some works try to satisfy one or more of the fairness notions in their methods, such as disparate learning processes (DLPs) which try to satisfy notions of treatment disparity and impact disparity by allowing the protected attributes during the training phase but avoiding them during prediction time [94].(Mehrabi et al., 2021)
- Some of the existing work tries to treat sensitive attributes as noise to disregard their effect on decision-making, while some causal methods use causal graphs, and disregard some paths in the causal graph that result in sensitive attributes affecting the outcome of the decision.(Mehrabi et al., 2021)
- Different bias-mitigating methods and techniques are discussed below for different domains—each targeting a different problem in different areas of machine learning in detail. (Mehrabi et al., 2021)

## Unbiasing Data

- Every dataset is the result of several design decisions made by the data curator. Those decisions have consequences for the fairness of the resulting dataset, which in turn affects the resulting algorithms. In order to mitigate the effects of bias in data, some general methods have been proposed that advocate having good practices while using data, such as having datasheets that would act like a supporting document for the data reporting the dataset creation method, its characteristics, motivations, and its skews [13, 55]. A similar suggestion has been proposed for models in [110].(Mehrabi et al., 2021)
- Authors in [66] also propose having labels, just like nutrition labels on food, in order to better categorize each data for each task. (Mehrabi et al., 2021)
- some work has targeted more specific types of biases. For example, [81] has proposed methods to test for cases of Simpson’s paradox in the data, and [3, 4] proposed methods to discover Simpson’s paradoxes in data automatically. (Mehrabi et al., 2021)
- Causal models and graphs were also used in some work to detect direct discrimination in the data along with its prevention technique that modifies the data such that the predictions would be absent from direct discrimination [163].(Mehrabi et al., 2021)
- in [62] also worked on preventing discrimination in data mining, targeting direct, indirect, and simultaneous effects.(Mehrabi et al., 2021)
- Other pre-processing approaches, such as messaging [74], preferential sampling [75, 76], disparate impact removal [51], also aim to remove biases from the data. (Mehrabi et al., 2021)
- Image quality. Several barriers to AI implementation in the clinic need to be overcome with regards to imaging (Figure 1). These include technical variations (e.g., camera hardware and software) and differences in image acquisition and quality (e.g., zoom level, focus, lighting, and presence of hair). For example, the presence of surgical ink markings is associated with decreased specificity (Winkler et al., 2019), field of view can significantly affect prediction quality (Mishra et al., 2019), and classification performance improves when hair and rulers are removed (Bisla et al., 2019). We have developed a method to measure how model predictions might be biased by the presence of a visual artifact (e.g., ink) and proposed methods to reduce such biases (Pfau et al., 2019). Poor quality images are often excluded from studies, but the problem of what makes an image adequate is not well studied. Ideally, models need to be able to express a level of confidence in a prediction as a function of image quality and appropriately direct a user to retake photos if needed. (Young et al., 2020) - dermatology

## Fair Classification



- certain methods have been proposed [57, 78, 85, 106] that satisfy certain definitions of fairness in classification. For instance, in [147] authors try to satisfy subgroup fairness in classification, equality of opportunity and equalized odds in [63], both disparate treatment and disparate impact in [2, 159], and equalized odds in [154]. (Mehrabi et al., 2021)
- Other methods try to not only satisfy some fairness constraints but to also be stable toward change in the test set [69] (Mehrabi et al., 2021)
- The authors in [155], propose a general framework for learning fair classifiers. This framework can be used for formulating fairness-aware classification with fairness guarantees. In another work [25], authors propose three different modifications to the existing Naive Bayes classifier for discrimination-free classification.(Mehrabi et al., 2021)
- paper [122] takes a new approach into fair classification by imposing fairness constraints into a Multitask learning (MTL) framework. In addition to imposing fairness during training, this approach can benefit the minority groups by focusing on maximizing the average accuracy of each group as opposed to maximizing the accuracy as a whole without attention to accuracy across different groups. In a similar work [49], authors propose a decoupled classification system where a separate classifier is learned for each group. They use transfer learning to reduce the issue of having less data for minority groups.(Mehrabi et al., 2021)
- In [73] authors propose to achieve fair classification by mitigating the dependence of the classification outcome on the sensitive attributes by utilizing the Wasserstein distance measure.(Mehrabi et al., 2021)
- In [75] authors propose the Preferential Sampling (PS) method to create a discrimination free train data set. They then learn a classifier on this discrimination free dataset to have a classifier with no discrimination.(Mehrabi et al., 2021)
- In [102], authors propose a post-processing bias mitigation strategy that utilizes attention mechanism for classification and that can provide interpretability. (Mehrabi et al., 2021)

**Fair Regression** TODO: only summarize briefly, as PASSION is a classification and not a regression task

- “price of fairness” (POF) to measure accuracy-fairness trade-offs, 3 penalites: Individual fairness, group fairness and hybrid fairness [14] (Mehrabi et al., 2021)
- In addition to the previous work, [1] considers the fair regression problem formulation with regards to two notions of fairness statistical (demographic) parity and bounded group loss. [2] uses decision trees to satisfy disparate impact and treatment in regression tasks in addition to classification. (Mehrabi et al., 2021)

**Structured Prediction** TODO: only summarize briefly, as PASSION is a classification task

- RBA (reducing bias amplification) as calibration algorithm to prevent risk of leveraging social bias, distributions in training data are followed in the predictions. multi-label object and visual semantic role labeling classification amplify existing bias in data [167] (Mehrabi et al., 2021) → TODO: be careful with this if the approach would be to generate new images for training!!

**Fair PCA** TODO: only summarize briefly, as PASSION is a classification task with only like 10 features

- Principal Component Analysis (PCA) <https://www.geeksforgeeks.org/principal-component-analysis-pca/> → dimensionality reduction, statistical technique, high-dimensional data into lower-dimensional space while maximising variance in new space → most important patterns and relationships is preserved
- vanilla PCA exaggerate error in reconstruction for one group of people [137] (Mehrabi et al., 2021)
- And their proposed algorithm is a two-step process listed below: (1) Relax the Fair PCA objective to a semidefinite program (SDP) and solve it. (2) Solve a linear program that would reduce the rank of the solution. [137] (Mehrabi et al., 2021)

**Community Detection** TODO: use this as an example for out of scope text, - Ludovic approved Community detection algorithms are specifically tailored to analyze network data and find connections in such datasets. For example, they can be used to detect groups of people with similar interest in social networks (Jayawickrama, 2021). This kind of data is not found in the context of PASSION, which is a classification task. Please refer to Mehrabi et al., 2021 for more information on bias mitigation in community detection algorithms.

**Causal Approach to Fairness** TODO: only relevant, if our variables have a dependency on the variables, e.g. age / gender determines how the disease is presenting itself in the images; check (Mehrabi et al., 2021) page 18 if relevant

**Fair Representation Learning** <https://medium.com/superlinear-eu-blog/representation-learning-breakthroughs-what-is-representation-learning-5dda2e2fed2e>

- Variational Auto encoders → Variational Fair Autoencoder introduced in [97]. Here, they treat the sensitive variable as the nuisance variable, so that by removing the information about this variable they will get a fair representation. They use a maximum mean discrepancy regularizer to obtain invariance in the posterior distribution over latent variables. Adding this maximum mean discrepancy (MMD) penalty into the lower bound of their

VAE architecture satisfies their proposed model for having the Variational Fair Autoencoder.

In [5] authors also propose a debiased VAE architecture called DB-VAE which learns sensitive latent variables that can bias the model (e.g., skin tone, gender, etc.) and propose an algorithm on top of this DB-VAE using these latent variables to debias systems like facial detection systems.

In [112] authors model their representation-learning task as an optimization objective that would minimize the loss of the mutual information between the encoding and the sensitive variable. The relaxed version of this assumption is shown in Equation 1. They use this in order to learn fair representation and show that adversarial training is unnecessary and in some cases even counter-productive.

In [42], authors introduce flexibly fair representation learning by disentanglement that disentangles information from multiple sensitive attributes. Their flexible and fair variational autoencoder is not only flexible with respect to downstream task labels but also flexible with respect to sensitive attributes. They address the demographic parity notion of fairness, which can target multiple sensitive attributes or any subset combination of them. (Mehrabi et al., 2021)

- **Adversarial Learning** - In [90] authors present a framework to mitigate bias in models learned from data with stereotypical associations. using adversarial networks by introducing FairGAN which generates synthetic data that is free from discrimination and is similar to the real data. They use their newly generated synthetic data from FairGAN, which is now debiased, instead of the real data for training and testing. They do not try to remove discrimination from the dataset, unlike many of the existing approaches, but instead generate new datasets similar to the real one which is debiased and preserves good data utility. (Mehrabi et al., 2021) **TODO: address challenges in creating synthetic data in dermatology?**

**Fair NLP** **TODO: for PASSION irrelevant, if it wants to stick to ResNet50 Architecture (Gottfrois et al., 2024) and not use Visual Encoders, which would make sense bc of the small dataset**

- **Word Embedding** **TODO: potentially relevant, if the labels are used in training, e.g. age / gender determines how the disease is presenting itself in the images; check (Mehrabi et al., 2021) page 21 if relevant**
- **Coreference Resolution** "Coreference resolution involves identifying when two or more expressions in a text refer to the same entity, be it a person, place, or thing." <https://medium.com/@datailm/the-key-to-unlocking-true-language-understanding-coreference-resolution-c01d569e2e87> **TODO: irrelevant for the PASSION Context**

**comparison of different mitigation algorithms**

- The field of algorithmic fairness is a relatively new area of research and work still needs to be done for its improvement. With that being said, there are already papers that propose fair AI algorithms and bias mitigation techniques and compare different mitigation algorithms using different benchmark datasets in the fairness domain. For instance, authors in [65] propose a geometric solution to learn fair representations that removes correlation between protected and unprotected features. The proposed approach can control the trade-off between fairness and accuracy via an adjustable parameter. In this work, authors evaluate the performance of their approach on different benchmark datasets, such as COMPAS, Adult and German, and compare them against various different approaches for fair learning algorithms considering fairness and accuracy measures [65, 72, 158, 159]. In addition, IBM's AI Fairness 360 (AIF360) toolkit [11] has implemented many of the current fair learning algorithms and has demonstrated some of the results as demos which can be utilized by interested users to compare different methods with regards to different fairness measures. (Mehrabi et al., 2021)

## 2.7 Statistical biases

<https://data36.com/statistical-bias-types-explained/>

- Self-selection bias - when you let the subjects of the analyses select themselves, less proactive people will be excluded TODO: could be an issue as well for PASSION, couldn't it? since the doctors probably ask the clients. One way to go is to default should be to provide access to the data. but is it ethical? (Mester\_2022Mester\_2017)- statistical bias
- Cause-effect bias (Mester\_2022Mester\_2017) - statistical bias
- Cognitive bias (Mester\_2022Mester\_2017) - statistical bias

## 2.8 Dermatology Bias

- <https://ijdvl.com/biases-in-dermatology-a-primer/> 29 biases, 4 reasons to know about it, 7 mitigation methods (Chakraborty, 2024)
- Popularity Bias: This bias arises when a particular disease is more popular (i.e. either more well-known or more stigmatised) among the participants than the disease with which it is compared. For example, if a study compares clinic attendance rates among various dermatological disorders, one would see vitiligo patients are over-represented over melasma. While melasma is commoner in the normal population, vitiligo, due to its popularity because of media publicity and other factors, tends to present earlier.<sup>9 6</sup>. (Chakraborty, 2024)

- All is well bias: It is a subjective bias where theories supported by the majority tend to get more easily published than the opposing view supported by the minority. For example, ideas on the origin of endemic pemphigus supporting autoimmunity are more likely to be published than theories exploring an infectious trigger. According to some authors, this bias is very difficult to eliminate and is a variant of publication bias.<sup>10-12</sup> 7.(Chakraborty, 2024)
- Apprehension bias: This results from fear and apprehensions related to an impending procedure. The classic example is the false elevation of blood pressure because the person is apprehensive of his or her blood pressure being measured.<sup>13</sup> A variant of this is the Hawthorne bias, where subjects modify their behavior, such as regularly taking a prescribed drug or exercising, simply because they know they are being watched, but not due to any apprehensions. Hawthorne bias is practically utilised in many leprosy clinics since regular follow-up has been shown to improve adherence to therapy based on Hawthorne bias. 8. (Chakraborty, 2024)
- Attrition bias: This occurs due to lack of follow-up. This is a common problem in studies evaluating the efficacy of biologics in psoriasis – where many patients are lost to follow-up. A remedy is performing intention to treat analysis.<sup>9</sup> A variation of this is non-responder bias, where non-responders to a questionnaire differ significantly from responders.<sup>9</sup> 9. (Chakraborty, 2024)
- Rhetoric bias: A more charismatic piece of writing has a greater influence on the study participants than other available literature. An example is the wider use of sunscreen for polymorphous light eruption over photoprotective strategies like umbrellas, broadbrimmed hats, etc, because the lay press is more vocal about sunscreens.<sup>14</sup> 11. (Chakraborty, 2024)
- Centripetal bias: Patients tend to go to more reputed physicians and hospitals than others. For example, a famous or better-known cosmetologist with a good reputation tends to see more cases than other cosmetologists.<sup>12</sup>.(Chakraborty, 2024)
- Confirmation bias: This bias occurs when study participants have a preconceived notion of their disease that may not be based on facts. For example, we have observed that in North India many tinea patients report an increase in their disease due to taking meat, fish, and other so-called “hot foods”. They may also present information they have collected from the internet which reinforces their beliefs.<sup>15</sup> 14.(Chakraborty, 2024)
- Novelty bias: The newer an intervention, the better it appears, and with time, its efficacy seems to decrease. When ligelizumab, an IgE antagonist was first discovered, ligelizumab was believed to be better than omalizumab; however, evidence soon pointed to the contrary. 16.(Chakraborty, 2024)
- Incorporation bias: This is principally relevant for diagnostic accuracy studies when the index test forms a part of the reference test, resulting in elevated

sensitivity e.g., if one wants to compare the grattage test vs. dermoscopy in psoriasis and does dermoscopy only from areas of grattage positivity, one would get a very high sensitivity for the grattage test because it was incorporated into the reference test, i.e., dermoscopy.<sup>25,26</sup> 21.(Chakraborty, 2024)

- Language bias: Articles with significant findings tend to get published more often in English (since that is the most common language in medical research) than in other languages. It is crucial in many studies involving dermatological quality of life measurements. <sup>24</sup>.(Chakraborty, 2024) - dermatology
- Previous opinion bias: In performing a second diagnostic test, if the result of a previous test is known, it is likely to influence the result. An extension of this is the Greenwald’s law of lupus: the Sontheimer amendment – anything and everything that happens to a lupus erythematosus patient is correctly or incorrectly attributed to lupus.<sup>32</sup> 29. (Chakraborty, 2024) - dermatology
- A recent study reported mean top-1 and top-5 model accuracy of 44.8% and 78.1%, respectively, for the classification of 134 diseases (Han et al., 2019b). Most datasets are proprietary, often with minimal description, and datasets collected in dermatology clinics may be skewed toward more complex cases, to those patients with better access to care, or by the choice of camera used in one clinic versus another. Data should be collected from as many diverse sources as possible, including primary care clinics, and robust standards for external validation are needed. (Young et al., 2020)
- There have been successful efforts to support reproducibility and open access. For example, the study by Han et al. (2018a) details the number and characteristics of images from each data source and makes thumbnails of the images publicly available. Additionally, several studies classifying dermoscopic images use the publicly available International Skin Imaging Collaboration archive (Gutman et al., 2016). By making datasets public, it becomes possible to examine them for bias (Bissoto et al., 2019). Alternatively, reporting a model training database’s patient demographics and disease classes would be helpful in predicting model performance on external populations. (Young et al., 2020)
- Metrics of model performance. Standard metrics are needed to assess the performance of different models (Figure 1). Currently, standard performance metrics such as accuracy and area under the receiver operating characteristic and precision recall curves are routinely reported. However, for use in the clinic, studies should additionally describe how well their models deal with uncertainty by reporting (i) the Brier Score, or mean-squared calibration error (Rufibach, 2010), which measures how reliably a model can forecast its accuracy, and (ii) area under the response rate accuracy curve, which measures how capably a model can identify examples it is likely to predict falsely and thus abstain from predicting (Hendrycks et al., 2019) (Young et al., 2020)



- **Model interpretability.** Acceptance of AI in clinical decision making hinges on being able to understand the decisionmaking process fundamental to its predictions. DL models are inherently difficult to interpret because they are complex, routinely containing millions of learned parameters; interpretation of DL models' output is an active field of research (Murdoch et al., 2019). One approach for interpreting model diagnoses is contentbased image retrieval, a method for retrieving training images that are visually similar to a test image (Tschandl et al., 2019a). This method may reassure the physician if all the retrieved training images have the same diagnosis as the predicted diagnosis but is less helpful if the test image looks similar to two or more training images with conflicting diagnoses. A second approach is to highlight pixels in an image most relevant for a model's prediction, using methods such as saliency mapping (Figure 1). However, it is often the case that highlighted pixels correspond to the entire lesion or visually distinctive features that are already obvious to clinicians without indication as to why these pixels are important to the diagnosis. A third approach is to see through the eyes of a model by plotting an activation atlas (Carter et al., 2019), which shows how subtle changes, in particular visual features, may tip the model over into choosing one diagnosis over another. These activation atlases are experimental and have yet to be applied in dermatology. Understanding a model's predictions and how the prediction is applicable to the patient at hand is necessary to build trust. As AI exceeds human performance in various tasks, interpreting models may help to advance scientific knowledge by understanding what the machine sees that is relevant to its predications (Young et al., 2020)

### 2.8.1 Demographic Bias in Dermatology

#### fairness melanoma detection

- Some biases can be easily detected and countered, such as through appropriate data curation; for example, having a balanced representation across skin tones and genders in training sets. However, in other cases, biases are hidden and untraceable [9]. (Montoya et al., 2025)
- whether information on demographic diversity (age, gender, race, or ethnicity of patients), clinical diversity (skin type, lesion type, anatomical location of lesion), or image characteristics (source, imaging techniques, resolution, and whether the images were real or artificially generated) (Montoya et al., 2025)
- The most popular skin color scale currently being used for data annotation for image recognition techniques is the Fitzpatrick Skin Tone Scale (FST) [10] which has six skin tones. Dating from the 1970s, it originally featured just 4 light tones and was designed for detecting photo sensitivity for white skin, with two darker tones added later [11]. The Monk Skin Scale was recently developed and still needs testing, but promisingly has 10 tones, 5 light and 5 dark [12].(Montoya et al., 2025) **TODO: highlight this (FST alternatives)**

- Fig. 4. Comparison of skin tone scales that can be used for skin cancer detection utilizing AI. Recreation of fitzpatrick skin type scale, monk skin tone scale, and sampling of L'Oreal color chart map for reference. (Montoya et al., 2025) **TODO: include this figure**
- While this systemic review provides a comprehensive review of the literature on fairness in AI for melanoma detection, it is primarily based on existing research. To validate the proposed recommendations or frameworks, continuing work is necessary to complete empirical analysis and experiments. Additionally, the suggested adoption of new skin tone scales, while beneficial, may face practical challenges in implementation. Furthermore, while the paper strongly advocates for specific skin tone scales, it's important to note that other methods or tools might also effectively address fairness issues in AI for melanoma detection. Finally, while the study addresses fairness in AI, it could benefit from further exploration of the practical implementation of these recommendations in real-world clinical settings. Potential obstacles and the feasibility of widespread adoption should be considered to ensure that the proposed solutions are not only theoretically sound but also practically viable. (Montoya et al., 2025) **TODO: also mention the limitations regarding FST alternatives**
- Recent research [13] adds another axis, skin hue, which is described as ranging from red to yellow. This offers a more complete representation of variations of skin color by providing a multidimensional scale [13]. (Montoya et al., 2025)
- The effect of hue (blue, red, yellow, green) on skin tones adds depth to each face producing a range of undertones (cold, neutral, warm, and olive). In the realm of color theory, the concept of 'contrast of hue' emphasizes the distinctiveness among fundamental colors, with primary hues like yellow, red, and blue exhibiting the most pronounced differences [14]. Because skin cancer appears differently on different colored skin, it is important to acknowledge a full range of colors present in both healthy skin and suspicious lesions within datasets used to train skin cancer detection ML tools. (Montoya et al., 2025)
- These findings should correlate to AI for melanoma detection since the contrast between skin color and skin lesions is a preliminary marker during feature extraction. Although the Fitzpatrick Skin Tone (FST) measurement scale is not diverse enough and leads to biased AI tools, it is continually used and has even been used to test a recently FDA-approved AI device for detecting melanoma. (Montoya et al., 2025)
- We advocate for the adoption of improved scales like the Monk and L'Oreal maps. Future studies should ensure equitable representation and testing across skin tones to guarantee AI's effectiveness for all. Please refer to Tables 2 through 7 in the discussion section for further recommendations for curating a diverse dataset, including purpose, ownership, funding, and data



annotation, as well as recommendations for each stage of the data life cycle. (Montoya et al., 2025) **TODO: Link for further mitigation methods**

- This study found that while using skin tone instead of race for fairness evaluations in computer vision seems objective, the annotation process remains biased by human annotators. Untested scales, unclear procedures, and a lack of awareness about annotator backgrounds and social context significantly influence skin tone labeling. This study exposes how even minor design choices in the annotation process, like scale order (dark to light instead of light to dark) or image context (face or no face, skin lesion presence), can sway agreement and introduce uncertainty in skin tone assessments. ... The researchers emphasize the need for greater transparency, standardized procedures, and careful consideration of annotator biases to mitigate these challenges and ensure fairer and more robust evaluations in computer vision. (Montoya et al., 2025) - demographic dermatology bias

## 3 Ideas and Concepts

Hier geht es um die Fragestellung, wie Sie die formulierten Ziele der Arbeit erreichen wollen. Sie halten z.B. erste, grobe Ideen, skizzenhafte Lösungsansätze fest. Gibt es mehrere Wege, Ansätze um dieses Ziel zu erreichen, begründen Sie hier, warum Sie einen bestimmten Weg einschlagen. Beispiel für ein Softwareprojekt: Erste Gedanken über eine grobe Systemarchitektur. Ist z.B. eine Microservice-Architektur angebracht? Welche Alternativen bestehen, wo gibt es Problempunkte? Die Umsetzung, die Beurteilung der Machbarkeit und die detaillierte Beschreibung der umgesetzten Architektur sind dann Teil der Realisierung.

### 3.1 PASSION Dataset

TODO: write things to consider more precisely:

- Include more details in gender attribute - transgender have probably different genes / hormones, and should be indicated for more accuracy
- include profession / at least an adapted version to indicate high risk patients for certain diseases? -> might lead to other biases?
- change country of origin to ethnicity (less of a proxy variable)
- are the data collectors specialized in some fields? That could lead to bias towards the center's country and the diagnosed diseases
- include images of healthy skin

### 3.2 Broad Methodology

TODO: write things to consider more precisely:

- Divide and Conquer vs. All-In-One-Model
  - An algorithm per ethnicity / subgroup running at the same time
  - Running 1 Algorithm chosen based on Fitzpatrick skin type
  - Running 1 Algorithm which detects first the demographic subgroup (FST, gender, age, ...) and runs the specific subgroup algorithm afterwards

- Hint Ludovic: Still not of data, maybe also others; often limited because the data is missing, you are missing data from others
- BLIND performance vs. Including the demographic data
  - Idea to try if the labels are not relevant for the diagnosis and should only be used for evaluating fairness purposes as some papers suggest
  - Might be obsolete after demographic biases in dermatology research, since melanin response and melanoma risk is different in male and female according to research <https://pmc.ncbi.nlm.nih.gov/articles/PMC4797181/>
- Hint Ludovic: Maybe Focal Loss more relevant → emphasis on data vs model
- Divide and Conquer vs. All-In-One-Model (either by ethnicity x algorithms at a time or one which separates the imgs first by demographic subgroup (incl. Fitzpatrick skin type))
- BLIND performance vs. Including the demographic data

## 4 Methods

Hier halten Sie fest und begründen, welches Vorgehensmodell Sie für Ihr Projekt wählen. Sie verweisen allenfalls auf die daraus entstandenen, konkreten Terminpläne mit Meilensteinen, welche z.B. unter Realisierung (Kapitel 5) oder im Anhang versorgt sind. Bei Projekten mit einer verlangten wissenschaftlichen Tiefe werden hier die geplanten Forschungsmethoden wie quantitative/qualitative Interviews, Befragungen, Beobachtungen, Feldexperiment etc. beschrieben und begründet. Warum ist in Ihrer Situation ein Interview besser als eine Umfrage? Wer soll interview werden? Die gewählten Methoden sind nachvollziehbar und begründet. Eine methodische Übersicht (Methodisches BigPicture) wurde aufgezeigt und Abgrenzungen erläutert.

## 5 Execution

Dies ist das Hauptkapitel Ihrer Arbeit! Hier wird die Umsetzung der eigenen Ideen und Konzepte (Kapitel 3) anhand der gewählten Methoden (Kapitel 4) beschrieben, inkl. der dabei aufgetretenen Schwierigkeiten und Einschränkungen. Die gewählten Methoden werden systematisch, konsistent und korrekt auf den Kontext der Arbeit angewendet. Die Bearbeitungs- bzw. Forschungsobjekte sind einheitlich benannt, im Kontext dargestellt und sinnvoll in die Arbeit integriert. Praxis- und Erfahrungswissen (z.B. aus Interviews) wird zur Validierung und Ergänzung der erarbeiteten Ergebnisse herangezogen.

## 6 Evaluation and Validation

Auswertung und Interpretation der Ergebnisse. Nachweis, dass die Ziele erreicht wurden, oder warum welche nicht erreicht wurden. Die Ziele / Forschungsfragen sind dem Umfang der Arbeit entsprechend sehr klar abgegrenzt; sie sind präzise, überprüfbar und nach den Standards der Zielformulierung definiert. Die Zielerreichung wurde systematisch und korrekt validiert. Die Herleitung und Bedeutung der Ergebnisse, mögliche Varianten, Gütekriterien und eine Validierung allgemein werden nachvollziehbar diskutiert

## 7 Outlook

Reflexion der eigenen Arbeit, ungelöste Probleme, weitere Ideen. Die Ergebnisse und Empfehlungen schaffen einen konkreten Mehrwert für die Auftraggebenden. Einschränkungen und Grenzen werden kritisch diskutiert und die nächsten Schritte im Ausblick festgehalten, so dass die Ergebnisse direkt in der Praxis weiterverwendet und/oder angewendet werden können.

## 8 Glossary

TODO: Add List of Formulas if necessary TODO: add AI declarations somewhere



## 9 Bibliography

- Baeza-Yates, R. (2018). Bias on the web. *Commun. ACM*, 61(6), 54–61. <https://doi.org/10.1145/3209581>  
Mehrabi 9.
- Berk, R., Heidari, H., Jabbari, S., Kearns, M., & Roth, A. (2017). Fairness in criminal justice risk assessments: The state of the art [Publisher: SAGE Publications Inc]. *Sociological Methods & Research*, 50(1), 3–44. <https://doi.org/10.1177/0049124118782533>  
Mehrabi 15.
- British Association of Dermatologists (BAD). (2021, July 7). *Lower socioeconomic status linked with more severe skin disease, including melanoma* [Bad patient hub] [Research was presented at the BAD’s Annual Meeting.]. Retrieved February 17, 2025, from <https://www.skinhealthinfo.org.uk/lower-socioeconomic-status-linked-with-more-severe-skin-disease-including-melanoma/>
- Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification [ISSN: 2640-3498]. *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, 77–91. Retrieved March 16, 2025, from <https://proceedings.mlr.press/v81/buolamwini18a.html>  
Mehrabi 24, demographic (skin type and gender).
- Chakraborty, A. (2024). Biases in dermatology: A primer [Publisher: Scientific Scholar]. *Indian J Dermatol Venereol Leprol*, 90(2), 250–254. [https://doi.org/10.25259/IJDVL\\_126\\_2023](https://doi.org/10.25259/IJDVL_126_2023)  
0 citations (but from 2024), list of lots of biases.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments [Publisher: Mary Ann Liebert, Inc., publishers]. *Big Data*, 5(2), 153–163. <https://doi.org/10.1089/big.2016.0047>  
Mehrabi 34.
- Ciampaglia, G. L., Nematzadeh, A., Menczer, F., & Flammini, A. (2018). How algorithmic popularity bias hinders or promotes quality [Publisher: Nature Publishing Group]. *Sci Rep*, 8(1), 15951. <https://doi.org/10.1038/s41598-018-34203-2>  
Mehrabi 117.
- Clarke, K. A. (2005). The phantom menace: Omitted variable bias in econometric research [Publisher: SAGE Publications Ltd]. *Conflict Management and Peace Science*, 22(4), 341–352. <https://doi.org/10.1080/07388940500339183>

- Mehrabi 38, difficultis regarding omitted variable and overcoming methods.
- Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., & Huq, A. (2017). Algorithmic decision making and the cost of fairness. *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 797–806. <https://doi.org/10.1145/3097983.3098095>  
Mehrabi 41.
- Danks, D., & London, A. J. (2017). Algorithmic bias in autonomous systems. *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, 4691–4697. <https://doi.org/10.24963/ijcai.2017/654>  
Mehrabi 44.
- Diaz, M., Lucke-Wold, B., Batchu, S., & Kleinberg, G. (2022). Racial underrepresentation in dermatological datasets leads to biased machine learning models and inequitable healthcare. *3*, 42–47.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, 214–226. <https://doi.org/10.1145/2090236.2090255>  
Mehrabi 48.
- Farnadi, G., Babaki, B., & Getoor, L. (2018). Fairness in relational domains. *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 108–114. <https://doi.org/10.1145/3278721.3278733>  
Mehrabi 50.
- Friedman, B., & Nissenbaum, H. (1996). Bias in computer systems. *ACM Trans. Inf. Syst.*, *14*(3), 330–347. <https://doi.org/10.1145/230538.230561>  
Mehrabi 53.
- Gottfrois, P., Gröger, F., Andriambololoniaina, F. H., Amruthalingam, L., Gonzalez-Jimenez, A., Hsu, C., Kessy, A., Lionetti, S., Mavura, D., Ng’ambi, W., Ngongonda, D. F., Pouly, M., Rakotoarisaona, M. F., Rapelanoro Rabenja, F., Traoré, I., & Navarini, A. A. (2024). Passion for dermatology: Bridging the diversity gap with pigmented skin images from sub-saharan africa. *Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, 703–712.
- Grgic-Hlača, N., Zafar, M. B., Gummadi, K. P., & Weller, A. (2016). The case for process fairness in learning: Feature selection for fair decision making  
Mehrabi 61.
- Hardt, M., Price, E., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, *29*. Retrieved March 16, 2025, from <https://proceedings.neurips.cc/paper/2016/hash/9d2682367c3935defcb1f9e247a97c0d-Abstract.html>  
Mehrabi 63.
- Hargittai, E. (2007). Whose space? differences among users and non-users of social network sites. *Journal of Computer-Mediated Communication*, *13*(1), 276–297. <https://doi.org/10.1111/j.1083-6101.2007.00396.x>  
Mehrabi 64.

- Jayawickrama, T. D. (2021, February 1). *Community detection algorithms* [Medium]. Retrieved March 24, 2025, from <https://medium.com/data-science/community-detection-algorithms-9bd8951e7dae>
- Kearns, M., Neel, S., Roth, A., & Wu, Z. S. (2018). Preventing fairness gerrymandering: Auditing and learning for subgroup fairness [ISSN: 2640-3498]. *Proceedings of the 35th International Conference on Machine Learning*, 2564–2572. Retrieved March 16, 2025, from <https://proceedings.mlr.press/v80/kearns18a.html>  
Mehrabani 79.
- Kearns, M., Neel, S., Roth, A., & Wu, Z. S. (2019). An empirical study of rich subgroup fairness for machine learning. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 100–109. <https://doi.org/10.1145/3287560.3287592>  
Mehrabani 80.
- Kusner, M. J., Loftus, J., Russell, C., & Silva, R. (2017). Counterfactual fairness. *Advances in Neural Information Processing Systems*, 30. Retrieved March 16, 2025, from [https://proceedings.neurips.cc/paper\\_files/paper/2017/hash/a486cd07e4ac3d270571622f4f316ec5-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2017/hash/a486cd07e4ac3d270571622f4f316ec5-Abstract.html)  
Mehrabani 87.
- Lerman, K., & Hogg, T. (2014). Leveraging position bias to improve peer recommendation [Publisher: Public Library of Science]. *PLOS ONE*, 9(6), e98914. <https://doi.org/10.1371/journal.pone.0098914>  
Mehrabani 93.
- Mehrabani, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning [Publisher: ACM PUB27New York, NY, USA]. *ACM Computing Surveys (CSUR)*. <https://doi.org/10.1145/3457607>
- Montoya, L. N., Roberts, J. S., & Hidalgo, B. S. (2025). Towards fairness in AI for melanoma detection: Systemic review and recommendations. In K. Arai (Ed.), *Advances in information and communication* (pp. 320–341). Springer Nature Switzerland. [https://doi.org/10.1007/978-3-031-84460-7\\_21](https://doi.org/10.1007/978-3-031-84460-7_21)  
2025.
- Mustard, D. B. (2003). Reexamining criminal behavior: The importance of omitted variable bias. *The Review of Economics and Statistics*, 85(1), 205–211. <https://doi.org/10.1162/rest.2003.85.1.205>  
Mehrabani 114.
- Olteanu, A., Castillo, C., Diaz, F., & Kiciman, E. (2019). Social data: Biases, methodological pitfalls, and ethical boundaries [Publisher: Frontiers]. *Front. Big Data*, 2. <https://doi.org/10.3389/fdata.2019.00013>  
Mehrabani 120.
- Riegg, S. K. (2008). Causal inference and omitted variable bias in financial aid research: Assessing solutions [Publisher: Johns Hopkins University Press]. *The Review of Higher Education*, 31(3), 329–354. Retrieved March 16, 2025, from <https://muse.jhu.edu/pub/1/article/232773>  
Mehrabani 131.

- Suresh, H., & Gutttag, J. (2021). A framework for understanding sources of harm throughout the machine learning life cycle. *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, 1–9. <https://doi.org/10.1145/3465416.3483305>  
Mehrabi 144.
- Verma, S., & Rubin, J. (2018). Fairness definitions explained. *Proceedings of the International Workshop on Software Fairness*, 1–7. <https://doi.org/10.1145/3194770.3194776>  
Mehrabi 149.
- Wang, T., & Wang, D. (2014). Why amazon’s ratings might mislead you: The story of herding effects [Publisher: Mary Ann Liebert, Inc., publishers]. *Big Data*, 2(4), 196–204. <https://doi.org/10.1089/big.2014.0063>  
Mehrabi 151.
- Young, A. T., Xiong, M., Pfau, J., Keiser, M. J., & Wei, M. L. (2020). Artificial intelligence in dermatology: A primer. *Journal of Investigative Dermatology*, 140(8), 1504–1512. <https://doi.org/10.1016/j.jid.2020.02.026>  
209 citations.

Projektspezifisch können weitere Dokumentationsteile angefügt werden wie: Aufgabenstellung, Projektmanagement-Plan/Bericht, Testplan/Testbericht, Bedienungsanleitungen, Details zu Umfragen, detaillierte Anforderungslisten, Referenzen auf projektspezifische Daten in externen Entwicklungs- und Datenverwaltungstools etc.

TODO: check the gls all unused.