



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS

ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΜΗΧΑΝΙΚΩΝ Η/Υ ΚΑΙ
ΠΛΗΡΟΦΟΡΙΚΗΣ

Ανάπτυξη συστήματος αναγνώρισης μεροληψίας
σε μεθόδους μηχανικής μάθησης

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

του Ορέστη Ι. Τσαγκέτα

Επιβλέπων:

Χρήστος Μακρής
Αναπληρωτής Καθηγητής
Τμήμα Μηχανικών Η/Υ και
Πληροφορικής
Πανεπιστήμιο Πατρών

Συνεπιβλέπον:

Ιωάννης Κανελλόπουλος

Πάτρα, Σεπτέμβριος 2024

Copyright © Ορέστης Ι. Τσαγκέτας, 2024

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, η αποθήκευση και η διανομή της παρούσας εργασίας, εξ' ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, η αποθήκευση και η διανομή για σκοπό μη-κερδοσκοπικό, εκπαίδευσης ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι αντιπροσωπεύουν τις επίσημες θέσεις του Πανεπιστημίου Πατρών.

Περίληψη

Η παρούσα διπλωματική εργασία ασχολείται με τη σχεδίαση και ανάπτυξη ενός εργαλείου που στοχεύει στον εντοπισμό και τη μείωση της μεροληψίας σε μεθόδους μηχανικής μάθησης. Το εργαλείο, το οποίο είναι μια διαδικτυακή εφαρμογή, αναπτύχθηκε σε Python χρησιμοποιώντας το Flask Framework και την βιβλιοθήκη Aif360 της IBM . Οι χρήστες της εφαρμογής καλούνται να ανεβάσουν ένα dataset με το οποίο εκπαιδεύουν ένα από τα διαθέσιμα μοντέλα μηχανικής μάθησης. Ο έλεγχος της μεροληψίας πραγματοποιείται βάσει του χαρακτηριστικού που επιλέγει ο χρήστης να ελεγχθεί, και αν επιβεβαιωθούν οι υποψίες του, να μειωθεί. Αυτή η διαδικασία διενεργείται μέσω της εφαρμογής, χρησιμοποιώντας διάφορες μετρικές μεροληψίας. Η μείωση της μεροληψίας επιτυγχάνεται με τη χρήση συγκεκριμένων αλγορίθμων, τους οποίους ο χρήστης μπορεί να επιλέξει ανάλογα με τις ανάγκες του. Για τη σωστή επιλογή των μετρικών και των αλγορίθμων, η εφαρμογή παρέχει καθοδήγηση στον χρήστη, λαμβάνοντας υπόψη τους περιορισμούς που προκύπτουν από τα χαρακτηριστικά των δεδομένων. Ο βασικός στόχος αυτής της διπλωματικής εργασίας είναι η εκπαίδευση και εξοικείωση των χρηστών που δεν διαθέτουν προγραμματιστικές γνώσεις ή βαθιά κατανόηση της μηχανικής μάθησης με την έννοια της δικαιοσύνης στους αλγορίθμους μηχανικής μάθησης. Επιπλέον, η εφαρμογή θα ελεγχθεί ώστε τα αποτελέσματά της να συμμορφώνονται με τον νόμο Local Law 144 of 2021 , που επιβάλλεται από το NYC Department of Consumer and Worker Protection (DCWP) . Ο νόμος αυτός απαιτεί διαφάνεια και δίκαιες πρακτικές στις αποφάσεις που λαμβάνονται μέσω αυτοματοποιημένων συστημάτων λήψης αποφάσεων, διασφαλίζοντας ότι δεν υπάρχει μεροληψία κατά συγκεκριμένων ομάδων πληθυσμού.

Λέξεις Κλειδιά: Αλγοριθμική Δικαιοσύνη, Μετρικές Δικαιοσύνης, Αλγόριθμοι Μείωσης Μεροληψίας, Μηχανική Μάθηση , Python, Aif360, Flask, Local Law 144 of 2021

Abstract

THE current Diploma Thesis focuses on the design and development of a tool aimed at detecting and reducing bias in machine learning methods. The tool, which is a web application, was developed in Python using the Flask Framework and the Aif360 toolkit from IBM. Users of the application are required to upload a dataset with which they train one of the available machine learning models. Bias detection is conducted based on the characteristic selected by the user to be checked, and if their suspicions are confirmed, to be reduced. This process is carried out through the application using various bias metrics. Bias reduction is achieved using specific algorithms that the user can choose based on their needs. For the correct selection of metrics and algorithms, the application provides guidance to the user, taking into account the constraints arising from the characteristics of the data. The primary goal of this thesis is to educate and familiarize users who do not have programming knowledge or a deep understanding of machine learning with the concept of fairness in machine learning algorithms. Additionally, the application will be tested to ensure its results comply with Local Law 144 of 2021, enforced by the NYC Department of Consumer and Worker Protection (DCWP). This law requires transparency and fair practices in decisions made through automated decision systems, ensuring that there is no bias against specific population groups.

Keywords: Algorithmic Fairness, Fairness Metrics, Bias Mitigation Algorithms, Machine Learning, Python, Aif360, Flask, Local Law 144 of 2021

*“Being good is easy, what is difficult is
being just.”*

— Victor Hugo

Ευχαριστίες

Θα ήθελα να ευχαριστήσω κ. Γ. Κανελλόπουλο και τον καθηγητή κ. Χρήστο Μακρή και για την επίβλεψη αλλά και για τη συμβολή τους στην εκπόνηση αυτής της διπλωματικής εργασίας.

Τέλος θα ήθελα να ευχαριστήσω τους γονείς μου για τη υποστήριξη και την ηθική συμπαράσταση που μου προσέφεραν όλα αυτά τα χρόνια.

Πάτρα, 1 Σεπτεμβρίου 2024

Περιεχόμενα

Κατάλογος Σχημάτων	11
Κατάλογος Πινάκων	13
1 Εισαγωγή	15
1.1 Πρόβλημα	16
1.2 Δομή Διπλωματικής Εργασίας	18
1.3 Συνεισφορά	19
2 Βιβλιογραφική Επισκόπηση	20
2.1 Η μηχανική μάθηση	20
2.2 Προεπεξεργασία δεδομένων	22
2.2.1 Καθαρισμός Δεδομένων	22
2.2.2 Μετασχηματισμός Δεδομένων	23
2.2.3 Μείωση Δεδομένων	24
2.2.4 Οφέλη της Προεπεξεργασίας Δεδομένων	24
2.3 Μόντελα μηχανικής Μάθησης	25
2.3.1 Λογιστική Παλινδρόμηση	25
2.3.1.1 Ο αλγόριθμος	25
2.3.1.2 Κατάλληλες Εφαρμογές	26
2.3.1.3 Πλεονεκτήματα	27
2.3.1.4 Τύποι Δεδομένων	28
2.3.1.5 Προεπεξεργασία	28
2.3.1.6 Περιορισμοί	28
2.3.2 Random Forest	28
2.3.2.1 Ο αλγόριθμος	29
2.3.2.2 Κατάλληλες Εφαρμογές	29

2.3.2.3	Πλεονεκτήματα	29
2.3.2.4	Τύποι Δεδομένων	30
2.3.2.5	Προεπεξεργασία	30
2.3.2.6	Περιορισμοί	30
2.3.3	Support Vector Machine	31
2.3.3.1	Ο αλγόριθμος	31
2.3.3.2	Πλεονεκτήματα	31
2.3.3.3	Κατάλληλες Εφαρμογές	32
2.3.3.4	Τύποι Δεδομένων	32
2.3.3.5	Προεπεξεργασία	32
2.3.3.6	Περιορισμοί	32
2.3.4	Naive Bayes	33
2.3.4.1	Ο αλγόριθμος	33
2.3.4.2	Πλεονεκτήματα	33
2.3.4.3	Κατάλληλες Εφαρμογές	34
2.3.4.4	Τύποι Δεδομένων	34
2.3.4.5	Προεπεξεργασία	34
2.3.4.6	Περιορισμοί	35
2.4	Μέθοδοι Αξιολόγησης Μοντέλου	35
2.4.1	Ακρίβεια (Accuracy)	36
2.4.2	Ανάκληση (Recall)	36
2.4.3	Ευστοχία (Precision)	36
2.4.4	F1-Score	37
2.5	Μεροληψία στην Τεχνητή Νοημοσύνη	37
2.6	Αλγοριθμική μεροληψία και δικαιοσύνη στη μηχανική μάθηση	39
2.6.1	Δικαιοσύνη στην τεχνητή νοημοσύνη	40
2.6.2	Σημασία της αλγοριθμικής και Δικαιοσύνης στη μηχανική μάθηση	41
2.7	Μετρικές Δικαιοσύνης	42

2.7.1	Μέση Διαφορά (Mean Difference)	43
2.7.2	Disparate Impact	43
2.7.3	Διαφορά Ίσων Ευκαιριών (Equal Opportunity Difference)	44
2.7.4	Διαφορά Μέσων Όρων (Average Odds Difference)	44
2.7.5	Δείκτης Theil (Theil Index)	45
2.8	Αλγόριθμοι Μείωσης Μεροληψίας	47
2.8.1	Αναπροσαρμογή Βαρών (Reweighting)	48
2.8.2	Adversarial Debiasing	48
2.8.3	Calibrated Equalized Odds	48
3	Τοπικός Νόμος 144 του 2021	49
3.1	Ανάλυση Τοπικού Νόμου 144 του 2021	50
3.2	Μεροληψία Διασταυρούμενων Χαρακτηριστικών	50
3.3	Αντιμετώπιση της Διασταυρούμενης Μεροληψίας	51
3.4	Αξιοποίηση του Εργαλείου AIF360	51
3.5	Πρακτική Εφαρμογή και Προκλήσεις	52
3.6	Ενίσχυση της Διαφάνειας με την Επεξηγήσιμη Τεχνητή Νοημοσύνη	52
3.7	Επιπτώσεις	53
4	Τεχνική Ανάλυση	54
4.1	Ανάλυση βάση του μοντέλου C4	54
4.2	Διάγραμμα Περιπτώσεων Χρήσης	63
4.3	Διάγραμματα Ευρωστίας	66
4.4	Διάγραμματα Ακολουθίας	68
4.5	Εργαλεία και Βιβλιοθήκες	69
4.5.1	Backend Τεχνολογίες	69
4.5.2	Frontend Τεχνολογίες	70
5	Παρουσίαση Εργαλείου και Εκπαιδευτικής Διαδικασίας	71
5.1	Αρχική Σελίδα	71

5.2	Επιλογή Συνόλου Δεδομένων	71
5.3	Επιλογή Προστατευόμενων Χαρακτηριστικών	72
5.4	Επιλογή Μετρικών Δικαιοσύνης	73
5.5	Παρουσίαση Αποτελεσμάτων Μετρικών Δικαιοσύνης	74
5.6	Επιλογή Αλγορίθμων Μείωσης Μεροληψίας	76
5.7	Τελικά Αποτελέσματα	76
6	Πειραματική Διαδικασία	79
6.1	Επιλογή και Περιγραφή Dataset	79
6.2	Επιλογή Μοντέλου Μηχανικής Μάθησης	79
6.3	Ανάλυση Dataset	81
6.3.1	Κατηγορηματικά χαρακτηριστικά:	81
6.3.2	Αριθμητικά χαρακτηριστικά:	81
6.3.3	Σημαντικότητα Χαρακτηριστικών	81
6.3.4	Συσχέτιση Χαρακτηριστικών	82
6.4	Επιλογή Προστατευμένων Χαρακτηριστικών και Ομάδων	84
6.5	Επιλογή Μετρικών Δικαιοσύνης	86
6.6	Επιλογή Αλγορίθμου Μείωσης Μεροληψίας	87
7	Συμπεράσματα και Μελλοντικές Επεκτάσεις	89
7.1	Συνοψη Αποτελεσμάτων	89
7.1.1	Συμμόρφωση με τον τοπικό Νόμο 144 του 2021	89
7.2	Συμπεράσματα	89
7.3	Μελλοντικές Κατευθύνσεις	90
	Βιβλιογραφία	91

Κατάλογος Σχημάτων

1	Είδη Μηχανικής Μάθησης	21
2	Διαδικασία εκπαίδευσης των Μοντέλων	22
3	Οπτική αναπαράσταση πίνακα σύγχυσης	37
4	Σταδία πιθανής εμφάνισης Μεροληψίας	38
5	Στάδια Εφαρμογής των αλγόριθμων μείωσης μεροληψίας	47
6	Διάγραμμα Περιβάλλοντος συστήματος	55
7	Container Diagram του συστήματος	55
8	Component Diagram του συστήματος	56
9	Διάγραμμα Περιπτώσεων Χρήσης	63
10	Διάγραμμα Ευρωστίας Περίπτωσης Χρήσης 1	66
11	Διάγραμμα Ευρωστίας Περίπτωσης Χρήσης 2	67
12	Διάγραμμα Ακολουθίας Περίπτωσης Χρήσης 1	68
13	Διάγραμμα Ακολουθίας Περίπτωσης Χρήσης 2	69
14	Σελίδα υποδοχής της εφαρμογής	71
15	Σελίδα επιλογής δεδομένων	72
16	Demo Σελίδα επιλογής δεδομένων	73
17	User metric guide	73
18	Demo Επιλογή Μετρικών Δικαιοσύνης	74
19	Demo Επιλογή Μετρικών Δικαιοσύνης	75
20	Demo Επιλογή Αλγορίθμων Μείωσης Μεροληψίας	76
21	Demo Επιλογή Αλγορίθμων Μείωσης Μεροληψίας	77
22	Demo Επιλογή Αλγορίθμων Μείωσης Μεροληψίας	78
23	Demo Επιλογή Μοντέλου Κατά τη Χρήση του Εργαλείου	80
24	Σημαντικότητα Χαρακτηριστικών	82
25	Συσχέτιση Χαρακτηριστικών	83
26	Συνδυασμοί φύλου και εθνικότητας με μέσους όρους	84

27	Επιλογή Προστατευμένων Χαρακτηριστικών και Ομάδων	86
28	Επιλογή Μετρικών Δικαιοσύνης	87
29	Επιλογή Αλγορίθμου Μείωσης Μεροληψίας	88

Κατάλογος Πινάκων

1	Βήματα του Αλγορίθμου Λογιστικής Παλινδρόμησης	27
2	Βήματα του Αλγορίθμου Δάσους Τυχαίων Δέντρων	29
3	Βήματα του Αλγορίθμου SVM	31
4	Βήματα του Αλγορίθμου Naive Bayes	33
5	Mean Difference	43
6	Disparate Impact	44
7	Equal Opportunity Difference	44
8	Average Odds Difference	45
9	Βήματα Υπολογισμού Δείκτη Theil (Theil Index Calculation Steps) .	46
10	Reweighting	48
11	Adversarial Debiasing	49
12	Calibrated Equalized Odds	49
13	Αρχεία Python και οι λειτουργίες τους	57
14	Συναρτήσεις του αρχείου proccess.py (Part 1)	58
15	Συναρτήσεις του αρχείου proccess.py (Part 2)	59
16	Συναρτήσεις του αρχείου proccess.py (Part 3)	60
17	Συναρτήσεις του αρχείου examples.py	61
18	HTML αρχεία και οι λειτουργίες τους	62
19	Επιλεγμένες Μετρικές Δικαιοσύνης Εκπαιδευτικής Διαδικασίας	74
20	Συνδυασμοί φύλου και εθνικότητας με μέσους όρους προβλεπόμενων και πραγματικών αποτελεσμάτων	84

Κεφάλαιο 1: Εισαγωγή

1 Εισαγωγή

Η ραγδαία εξάπλωση της μηχανικής μάθησης (Machine Learning) σε διάφορους τομείς, από εγκρίσεις δανείων και συστήματα αναγνώρισης προσώπου μέχρι προβλέψεις στην ποινική δικαιοσύνη, έχει φέρει σημαντικά οφέλη στην κοινωνία, αυξάνοντας την παραγωγικότητα και την ακρίβεια για τη λήψη αποφάσεων και ενεργειών. Ωστόσο, υπάρχει μια αυξανόμενη ανησυχία σχετικά με τη δυνατότητα εμφάνισης προκαταλήψεων σε αυτούς τους ισχυρούς αλγόριθμους. Οι προκαταλήψεις στα μοντέλα ΜΛ μπορούν να οδηγήσουν ενίσχυση των κοινωνικών ανισοτήτων και διακρίσεων, καθώς και την έλλειψη διαφάνειας και λογοδοσίας, που δυσχεραίνει τον εντοπισμό και τη διόρθωση αυτών των προκαταλήψεων.

Αυτή η διπλωματική εργασία ασχολείται με το κρίσιμο ζήτημα της ανίχνευσης και μείωσης των προκαταλήψεων στις μεθόδους ML. Παρουσιάζουμε το σχεδιασμό και την ανάπτυξη μιας φιλικής προς τον χρήστη διαδικτυακής εφαρμογής που δίνει τη δυνατότητα σε άτομα, ακόμη και χωρίς εκτεταμένη γνώση προγραμματισμού, να εντοπίζουν και να μειώνουν τις πιθανές προκαταλήψεις στα μοντέλα μηχανικής μάθησης τους.

Αυτή η εργασία συμβάλλει στον τομέα της αλγοριθμικής δικαιοσύνης παρέχοντας ένα πρακτικό εργαλείο που ενισχύει τη διαφάνεια και προάγει τις ανησυχίες δικαιοσύνης καθ' όλη τη διάρκεια ανάπτυξης των μοντέλων μηχανικής μάθησης. Η εφαρμογή αξιοποιεί τη βιβλιοθήκη Aif360 από την IBM [1] για την ανάλυση των δεδομένων που παρέχουν οι χρήστες και την ανίχνευση πιθανών προκαταλήψεων βάσει καθορισμένων από τον χρήστη χαρακτηριστικών, όπως η φυλή, το φύλο ή η ηλικία. Αυτό επιτρέπει στους χρήστες να εντοπίζουν περιοχές όπου τα μοντέλα τους μπορεί να παρουσιάζουν

άδικες προτιμήσεις προς συγκεκριμένες δημογραφικές ομάδες.

Επιπλέον, η εφαρμογή προχωρά πέρα από την απλή ανίχνευση προκαταλήψεων, προτείνοντας κατάλληλες τεχνικές μείωσης προκαταλήψεων. Συνιστά κατάλληλους αλγόριθμους για τη μείωση του ανιχνευόμενου τύπου προκατάληψης, λαμβάνοντας υπόψη τα συγκεκριμένα χαρακτηριστικά των δεδομένων και τις πιθανές περιορισμούς. Αυτό δίνει τη δυνατότητα στους χρήστες να αντιμετωπίζουν ενεργά τις ανησυχίες δικαιοσύνης και να διασφαλίζουν ότι τα μοντέλα τους λειτουργούν με υπεύθυνο και ηθικό τρόπο.

Επιπλέον, προωθεί τη συμμόρφωση με κανονισμούς όπως ο Τοπικός Νόμος 144 του 2021 που επιβάλλεται από το Τμήμα Προστασίας Καταναλωτών και Εργαζομένων της Νέας Υόρκης (DCWP) [2]. Αυτός ο νόμος απαιτεί διαφάνεια και δικαιοσύνη στα αυτοματοποιημένα συστήματα αποφάσεων, ευθυγραμμιζόμενος απόλυτα με τον στόχο μας για την προώθηση υπεύθυνων και ηθικών πρακτικών της τεχνητής νοημοσύνης. Με την ενεργή μείωση των προκαταλήψεων και τη διασφάλιση της διαφάνειας στη διαδικασία ανάπτυξης μοντέλων, η εφαρμογή δίνει τη δυνατότητα στους χρήστες να συμμορφώνονται με τέτοιους κανονισμούς χωρίς να διακυβεύεται η λειτουργικότητα των MLμοντέλων τους.

1.1 Πρόβλημα

Υπάρχουν πολλές αξιοσημείωτες περιπτώσεις που υπογραμμίζουν τη σημασία της δικαιοσύνης στα AI συστήματα. Ένα χαρακτηριστικό παράδειγμα είναι το αυτοματοποιημένο εργαλείο πρόσληψης της Amazon [3]. Ξεκίνησε το 2014, αυτό το εργαλείο χρησιμοποιήθηκε για την αξιολόγηση βιογραφικών και την βαθμολόγηση υποψηφίων. Ωστόσο, ένα χρόνο μετά το 2015, διαπιστώθηκε ότι το σύστημα πρόσληψης δεν βαθμολογούσε δίκαια τους υποψηφίους, καθώς ευνοούσε τους άντρες υποψηφίους έναντι των γυναικών. Αυτή η μεροληψία προέκυψε επειδή το εργαλείο είχε εκπαιδευτεί με βιογραφικά που είχαν υποβληθεί στην Amazon κατά τη διάρκεια μιας δεκαετίας, τα περισσότερα από τα οποία προέρχονταν από άντρες [4].

Ένα άλλο χαρακτηριστικό παράδειγμα είναι το λογισμικό COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)[5] που χρησιμοποιήθηκε από την κυβέρνηση των Η.Π.Α., το οποίο υπολόγιζε βάσει των δεδομένων των κατηγορουμένων ένα σκορ (1 έως 10). Αυτά τα σκορ βοηθούσαν τους δικαστές να αποφασίσουν για ποινές, αναστολές και άλλες δικαστικές αποφάσεις. Ωστόσο, μελέτες αποκάλυψαν ότι ο αλγόριθμος είχε μεγαλύτερη πιθανότητα να προβλέψει λανθασμένα ότι οι μαύροι κατηγορούμενοι θα επαναλάμβαναν το έγκλημα σε σύγκριση με τους λευκούς κατηγορούμενους, οδηγώντας σε δυσανάλογες επιπτώσεις στις αποφάσεις για την ποινή και την αναστολή [6] .

Αυτά τα παραδείγματα υπογραμμίζουν την κρίσιμη ανάγκη αντιμετώπισης των μεροληψιών στα συστήματα AI για να διασφαλιστεί η δικαιοσύνη και η ισότητα. Η έρευνα εστιάζει σε μεθόδους για τον μετριασμό αυτών των μεροληψιών, όπως η ανάπτυξη πιο διαφανών αλγορίθμων και η ενσωμάτωση μετρήσεων δικαιοσύνης στον σχεδιασμό και την αξιολόγηση των συστημάτων. Μία προσέγγιση είναι η χρήση τεχνικών μηχανικής μάθησης που είναι ευαίσθητες στη δικαιοσύνη και προσαρμόζουν τη διαδικασία μάθησης για να ελαχιστοποιήσουν τη μεροληψία. Μια άλλη μέθοδος είναι η διενέργεια ελέγχων μεροληψίας για τον εντοπισμό και τη διόρθωση των μεροληψιών πριν από την ανάπτυξη των AI συστημάτων σε κρίσιμες εφαρμογές.

Για παράδειγμα, η ενσωμάτωση περιορισμών για την επίτευξη δίκαιων συστημάτων κατά τη διαδικασία εκπαίδευσης μπορεί να βοηθήσει στη διασφάλιση ότι τα παραγόμενα μοντέλα δεν επηρεάζουν δυσανάλογα καμία συγκεκριμένη ομάδα [7]. Επιπλέον, η χρήση τεχνικών εξηγήσιμης AI (XAI) μπορεί να προσφέρει πληροφορίες για το πώς λαμβάνονται οι αποφάσεις από τα AI συστήματα, καθιστώντας ευκολότερο τον εντοπισμό και τη διόρθωση της μεροληπτικής συμπεριφοράς.

Οι πρόσφατες εξελίξεις στη μείωση της μεροληψίας περιλαμβάνουν την ανάπτυξη τεχνικών αντίπαλης αποπροκατάληψης (adversarial debiasing), οι οποίες περιλαμβάνουν την εκπαίδευση των AI μοντέλων με τρόπο που οι αντίπαλοι προσπαθούν να εισαγάγουν μεροληψία και το κύριο μοντέλο μαθαίνει να την εξουδετερώνει [8]. Αυτή η

μέθοδος έχει δείξει υποσχέσεις για τη μείωση της μεροληψίας σε διάφορες εφαρμογές, από τις προσλήψεις μέχρι τη δικαιοσύνη. Επιπλέον, χρησιμοποιούνται αλγόριθμοι προεπεξεργασίας που τροποποιούν τα δεδομένα πριν από την εκπαίδευση των μοντέλων, καθώς και αλγόριθμοι μετά την επεξεργασία, οι οποίοι προσαρμόζουν τις εξόδους των μοντέλων για να εξασφαλίσουν δίκαια αποτελέσματα.

Η αντιμετώπιση της μεροληψίας στην τεχνητή νοημοσύνη απαιτεί λύσεις που συνδυάζουν γνώσεις από την πληροφορική, το δίκαιο, την ηθική και τις κοινωνικές επιστήμες. Γι' αυτό, η έμφαση στρέφεται ολοένα και περισσότερο στις διεπιστημονικές προσεγγίσεις. Μέσα από συνεργατικές προσπάθειες και ερευνητικές πρακτικές, μπορούμε να αναπτύξουμε λύσεις που διασφαλίζουν ότι δεν αδικούνται και θα αντιμετωπίζονται ισότιμα όλα τα μέλη της κοινωνίας.

1.2 Δομή Διπλωματικής Εργασίας

Στα πλαίσια της διπλωματικής εργασίας μελετήθηκαν διάφορες διεθνείς δημοσιεύσεις που αφορούν την έννοια της αλγοριθμικής δικαιοσύνης, τις μετρικές εκτίμησης της αλγοριθμικής μεροληψίας και τους αλγορίθμους μείωσης της. Το εργαλείο που παρουσιάζουμε, κατασκευάστηκε με χρήση του Flask [9], ενώ οι λειτουργίες της που αφορούν τη μέτρηση και τη μείωση της αλγοριθμικής μεροληψίας υλοποιήθηκαν με τη βιβλιοθήκη AI Fairness 360 της IBM [10].

Στο κεφάλαιο 2 γίνεται βιβλιογραφική ανασκόπηση σχετικά με τη μηχανική μάθηση. Το κεφάλαιο περιλαμβάνει μια εισαγωγή στη μηχανική μάθηση, την προεπεξεργασία δεδομένων όπως καθαρισμός, ενοποίηση, μετασχηματισμός και μείωση δεδομένων, καθώς και ανάλυση διαφόρων μοντέλων μηχανικής μάθησης, όπως Logistic Regression, Random Forest, Support Vector Machine και Naive Bayes. Επιπλέον, παρουσιάζονται οι μέθοδοι αξιολόγησης μοντέλων, όπως η ακρίβεια, η ανάκληση και άλλες μετρικές. Συζητείται η μεροληψία στη μηχανική μάθηση και οι επιπτώσεις της αλγοριθμικής δικαιοσύνης, ενώ παρουσιάζονται παρουσιάζονται οι μετρικές και οι αλγόριθμοι μείωσης της μεροληψίας που συμπεριλήφθηκαν στην παρούσα εργασία. Επιπλέον, γίνεται μελέτη

των νομικών και ηθικών ζητημάτων, όπως ο νόμος Local Law 144 of 2021 .

Στο κεφάλαιο 3 παρουσιάζεται αναλυτικά η αρχιτεκτονική του συστήματος, ο σχεδιασμός και οι τεχνολογίες στις οποίες βασίζεται.

Στο κεφάλαιο 4 περιγράφεται η διαδικασία αξιολόγησης της εφαρμογής και τα αποτελέσματα που προέκυψαν από αυτή.

Στο κεφάλαιο 5 παρατίθενται τα συμπεράσματα από την αξιολόγηση και τα αποτελέσματα της εργασίας. Τέλος, σκιαγραφούνται οι μελλοντικές ερευνητικές κατευθύνσεις και τα ζητήματα που προκύπτουν από την εργασία, καθώς και οι περιορισμοί που εντοπίστηκαν κατά την υλοποίηση και την αξιολόγηση. Συγκεκριμένα, επισημαίνονται τα προβλήματα και οι προκλήσεις που συνδέονται με την εξασφάλιση της αλγοριθμικής δικαιοσύνης σε διαφορετικά πλαίσια εφαρμογής και προτείνονται λύσεις και κατευθύνσεις για περαιτέρω έρευνα. Οι περιορισμοί που εντοπίστηκαν περιλαμβάνουν την ανάγκη για μεγαλύτερα και πιο ποικίλα δεδομένα εκπαίδευσης, την αυξημένη πολυπλοκότητα των αλγορίθμων μείωσης της μεροληψίας, καθώς και την ανάγκη για συνεχή ενημέρωση και προσαρμογή στις τρέχουσες νομικές και ηθικές απαιτήσεις.

1.3 Συνεισφορά

Η παρούσα διπλωματική εργασία εστιάζει στην ανάπτυξη ενός εργαλείου αξιολόγησης της αλγοριθμικής δικαιοσύνης, το οποίο ευθυγραμμίζεται με τις προδιαγραφές του Νόμου 144 του 2021 ("Local Law 144 of 2021") της Νέας Υόρκης. Στόχος είναι να καταστήσουμε το εργαλείο εύχρηστο και προσβάσιμο σε χρήστες με ή χωρίς εξειδικευμένες γνώσεις, ώστε να μπορούν να αξιολογούν τη λειτουργία των μοντέλων τεχνητής νοημοσύνης που υλοποιούν, να εντοπίζουν τυχόν προκαταλήψεις και να υιοθετούν στρατηγικές για την μείωση ή την εξάλειψή τους.

Πέρα από την τήρηση του νομικού πλαισίου, το εργαλείο μας φιλοδοξεί να προσφέρει ουσιαστική αξία στον πραγματικό κόσμο, συμβάλλοντας στην υιοθέτηση ηθικών και δίκαιων εφαρμογών της τεχνητής νοημοσύνης σε διάφορους τομείς.

Κεφάλαιο 2

2 Βιβλιογραφική Επισκόπηση

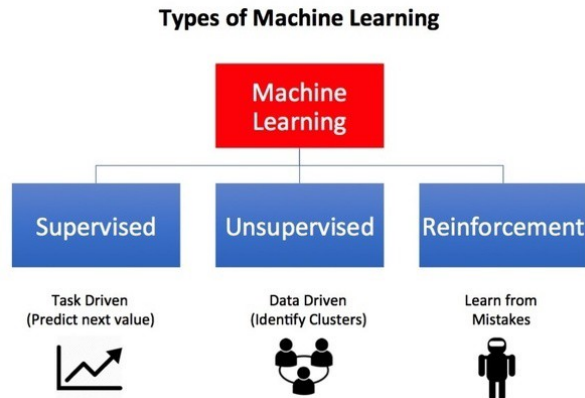
Η βιβλιογραφική επισκόπηση εστιάζει στην τρέχουσα κατάσταση της έρευνας στον τομέα της αλγοριθμικής προκατάληψης και της δικαιοσύνης στη μηχανική μάθηση. Η ενότητα αυτή καλύπτει διεξοδικά θεμελιώδεις έννοιες, μετρικές και τεχνικές μετριάσμού που σχετίζονται με την παρούσα μελέτη. Επιπλέον, φέρνει στο προσκήνιο το νομικό πλαίσιο του Νόμου 144 του 2021 (Local Law 144 of 2021),που θέτει το πλαίσιο για την τήρηση των κανονισμών στην παρούσα έρευνα.

2.1 Η μηχανική μάθηση

Η μηχανική μάθηση (ML) είναι ένας υποτομέας της τεχνητής νοημοσύνης (AI) που περιλαμβάνει την ανάπτυξη αλγορίθμων που επιτρέπουν στους υπολογιστές να μαθαίνουν από δεδομένα και να κάνουν προβλέψεις ή να λαμβάνουν αποφάσεις με βάση αυτά. Σε αντίθεση με τον παραδοσιακό προγραμματισμό, όπου ο υπολογιστής ακολουθεί ρητές οδηγίες, τα μοντέλα μηχανικής μάθησης εκπαιδεύονται σε δεδομένα για να αναγνωρίζουν πρότυπα και να λαμβάνουν αποφάσεις με ελάχιστη ανθρώπινη παρέμβαση.

Υπάρχουν τρεις κύριοι τύποι μηχανικής μάθησης: η επιβλεπόμενη μάθηση (supervised learning), η μη επιβλεπόμενη μάθηση (unsupervised learning) και η ενισχυτική μάθηση (reinforcement learning). Στην επιβλεπόμενη μάθηση, το μοντέλο εκπαιδεύεται σε ένα δεδομένο σύνολο με ετικέτες, πράγμα που σημαίνει ότι κάθε παράδειγμα εκπαίδευσης συνοδεύεται από μια ετικέτα εξόδου.

Παραδείγματα αλγορίθμων επιβλεπόμενης μάθησης περιλαμβάνουν τους αλγόριθμους logistic regression, support vector machines και τα νευρωνικά δίκτυα. Η μη



Σχήμα 1: Είδη Μηχανικής Μάθησης

επιβλεπόμενη μάθηση περιλαμβάνει την εκπαίδευση ενός μοντέλου σε δεδομένα χωρίς ετικέτες, αναζητώντας κρυφά πρότυπα ή εσωτερικές δομές. Παραδείγματα περιλαμβάνουν αλγόριθμους ομαδοποίησης όπως ο k-means. Η ενισχυτική μάθηση είναι ένας τύπος μάθησης όπου ένας πράκτορας μαθαίνει να λαμβάνει αποφάσεις εκτελώντας ενέργειες σε ένα περιβάλλον για να μεγιστοποιήσει κάποια έννοια συνολικής ανταμοιβής.

Η μηχανική μάθηση έχει εφαρμοστεί επιτυχώς σε διάφορους τομείς όπως η υγειονομική περίθαλψη, η χρηματοοικονομική, το μάρκετινγκ και τα αυτόνομα συστήματα, δείχνοντας τη δυνατότητά της να μετασχηματίζει βιομηχανίες αυτοματοποιώντας πολύπλοκες εργασίες και παρέχοντας πολύτιμες πληροφορίες από δεδομένα [11, 12].

Η συνηθέστερη διαδικασία κατασκευής ενός μοντέλου μηχανικής μάθησης περιλαμβάνει τα εξής βασικά βήματα:

1. Συλλογή Δεδομένων

2. Προεπεξεργασία Δεδομένων

- Κωδικοποίηση κατηγορηματικών χαρακτηριστικών (*Encoding Categorical Features*)
 - One-Hot Encoding [13]
 - Label Encoding [14, 15]

3. Επιλογή Αλγορίθμου

4. Εκπαίδευση Μοντέλου

5. Αξιολόγηση Μοντέλου

Κάθε βήμα είναι σημαντικό για τη δημιουργία ενός αποτελεσματικού μοντέλου που μπορεί να χρησιμοποιηθεί για προβλέψεις ή λήψη αποφάσεων.



Σχήμα 2: Διαδικασία εκπαίδευσης των Μοντέλων

2.2 Προεπεξεργασία δεδομένων

Η προεπεξεργασία δεδομένων αποτελεί θεμέλιο λίθο στη διαδικασία της μηχανικής μάθησης. Σαν ουσιαστικό βήμα, φροντίζει για την κατάλληλη προετοιμασία των δεδομένων, ώστε να τροφοδοτήσουν με ακρίβεια και αποτελεσματικότητα τα μοντέλα μηχανικής μάθησης, τόσο κατά την εκπαίδευση όσο και κατά την αξιολόγησή τους.

Σαν μια απαραίτητη διαδικασία, η προεπεξεργασία δεδομένων περιλαμβάνει ένα

Στη συνέχεια αναλύται το εύρος εργασιών που εμπλέκονται στην προεπεξεργασία δεδομένων.

2.2.1 Καθαρισμός Δεδομένων

Στόχος του καθαρισμού είναι ο εντοπισμός και η αντιμετώπιση τυχόν σφαλμάτων, ελλείψεων ή ασυνεπειών που δύναται να επηρεάσουν αρνητικά την εκπαίδευση και την απόδοση ενός μοντέλου.

- **Διαχείριση Ελλιπών Τιμών:** Οι ελλιπείς τιμές μπορούν να συμπληρωθούν χρησιμοποιώντας διάφορες τεχνικές όπως μέση τιμή, διάμεσος ή και να διαγραφούν αν η χρήση των παραπάνω τεχνικών δεν είναι ορθολογικά σωστή.
- **Κανονικοποίηση Τιμών:** Οι ακραίες τιμές μπορούν να αφαιρεθούν ή να αναπροσαρμοστούν, δηλαδή να περιοριστούν σε ένα συγκεκριμένο όριο για να μειωθεί η επίδρασή τους στο μοντέλο.
- **Εξαληψή Ασυνεπειών:** Οι ασυνέπειες στα δεδομένα, δηλαδή δεδομένα που δεν ακολουθούν την ίδια μορφοποίηση ή περιστάσεις όπου διαφορετικές πηγές δεδομένων δίνουν διαφορετικές τιμές για το ίδιο χαρακτηριστικό. Μπορούν να διορθωθούν ή να αφαιρεθούν για να εξασφαλιστεί η ομοιομορφία και η ακρίβεια.

2.2.2 Μετασχηματισμός Δεδομένων

Αυτό το βήμα περιλαμβάνει τον μετασχηματισμό των δεδομένων σε μορφή κατάλληλη για τον επιλεγμένο αλγόριθμο μηχανικής μάθησης.

- **Κλιμάκωση Αριθμητικών Χαρακτηριστικών:** Η κλιμάκωση διασφαλίζει ότι τα αριθμητικά χαρακτηριστικά έχουν συνεπή εμβέλεια, κάτι που είναι κρίσιμο για αλγόριθμους που είναι ευαίσθητοι στην κλίμακα των εισαγωγών δεδομένων.
- **Κωδικοποίηση Κατηγοριματικών Χαρακτηριστικών:** Τα κατηγορικά χαρακτηριστικά μετατρέπονται σε αριθμητικές τιμές χρησιμοποιώντας τεχνικές όπως οι one-hot encoding ή label encoding[15].
- **Κατασκευή Χαρακτηριστικών:** Νέα χαρακτηριστικά δημιουργούνται από τα υπάρχοντα για να παρέχουν πρόσθετες πληροφορίες στο μοντέλο.

2.2.3 Μείωση Δεδομένων

Σε ορισμένες περιπτώσεις, τα σύνολα δεδομένων μπορεί να είναι πολύ μεγάλα και υπολογιστικά ακριβά για να δουλεутούν. Τεχνικές μείωσης δεδομένων όπως η μείωση διαστάσεων μπορούν να χρησιμοποιηθούν για να μειώσουν τον αριθμό των χαρακτηριστικών χωρίς να χανθεί σημαντική πληροφορία.

- **Μείωση Διαστάσεων:** Τεχνικές όπως η Ανάλυση Κύριων Συνιστωσών (Principal Component Analysis ή PCA) ή η t-Distributed Stochastic Neighbor Embedding (t-SNE) βοηθούν στη μείωση του αριθμού των χαρακτηριστικών διατηρώντας την ουσιώδη πληροφορία.

2.2.4 Οφέλη της Προεπεξεργασίας Δεδομένων

- **Βελτιωμένη Απόδοση Μοντέλου:** Τα προεπεξεργασμένα δεδομένα οδηγούν σε πιο ακριβή και αποδοτικά μοντέλα μηχανικής μάθησης.
- **Μειωμένος Χρόνος Εκπαίδευσης:** Τα καθαρά και οργανωμένα δεδομένα επιτρέπουν στα μοντέλα να εκπαιδεύονται πιο γρήγορα, βελτιστοποιώντας τους υπολογιστικούς πόρους.
- **Αυξημένη Ερμηνευσιμότητα Μοντέλου:** Η προεπεξεργασία βοηθά στον εντοπισμό σημαντικών χαρακτηριστικών και σχέσεων μέσα στα δεδομένα, κάνοντας το μοντέλο πιο εύκολο να ερμηνευτεί και να κατανοηθεί.

Εφαρμόζοντας με επιμέλεια αυτά τα βήματα προεπεξεργασίας δεδομένων, μπορούμε να βελτιώσουμε σημαντικά την ποιότητα του συνόλου δεδομένων μας, οδηγώντας σε καλύτερη απόδοση και αξιοπιστία των μοντέλων μηχανικής μάθησης. Κάθε βήμα παίζει ζωτικό ρόλο στην εξασφάλιση ότι τα δεδομένα είναι καθαρά, ενσωματωμένα, μετασχηματισμένα και μειωμένα κατάλληλα, θέτοντας μια στέρεη βάση για οποιοδήποτε έργο μηχανικής μάθησης.

2.3 Μόντελα μηχανικής Μάθησης

Αυτή η ενότητα τα διαθέσιμα μοντέλα εποπτευόμενης μάθησης που μπορεί ο χρήστης να χρησιμοποιήσει μέσω του εργαλείου. Η εποπτευόμενη μάθηση, ένας ακρογωνιαίος λίθος της μηχανικής μάθησης, δίνει τη δυνατότητα στους αλγόριθμους να μαθαίνουν από δεδομένα που έχουν προεπισημανθεί με επιθυμητά αποτελέσματα. Αυτό επιτρέπει στα μοντέλα να εκτελούν ενέργειες όπως η ταξινόμηση και η παλινδρόμηση. Θα εξετάσουμε τέσσερα σημαντικά μοντέλα εποπτευόμενης μάθησης: Logistic Regression, Random Forest, Support Vector Machine (SVM), and Naive Bayes. Κάθε μοντέλο διαθέτει διακριτά χαρακτηριστικά και δείχνει αποτελεσματικότητα στην αντιμετώπιση διαφόρων τύπων προβλημάτων στο πλαίσιο αυτής της έρευνας.

2.3.1 Λογιστική Παλινδρόμηση

Αποτελεί θεμελιώδη αλγόριθμο εποπτευόμενης μάθησης που χρησιμοποιείται ευρέως για προβλήματα δυαδικής ταξινόμησης [16]. Εκτιμά την πιθανότητα ενός σημείου δεδομένων να ανήκει σε μια συγκεκριμένη κατηγορία (π.χ., η πιθανότητα έγκρισης ενός δανείου) [17].

2.3.1.1 Ο αλγόριθμος Ο αλγόριθμος χρησιμοποιεί ένα γραμμικό μοντέλο παλινδρόμησης για την πρόβλεψη μιας συνεχούς τιμής μεταξύ αρνητικού άπειρου και θετικού άπειρου. Ωστόσο, για εργασίες ταξινόμησης, χρειαζόμαστε μια πιθανότητα μεταξύ 0 και 1 [;]. Για να το επιτύχουμε αυτό, η λογιστική παλινδρόμηση εφαρμόζει μια συνάρτηση sigmoid (επίσης γνωστή ως λογιστική συνάρτηση) στην έξοδο του γραμμικού μοντέλου [18]. Η συνάρτηση sigmoid μετατρέπει την συνεχή έξοδο σε τιμή πιθανότητας.

Ακολουθεί η μαθηματική αναπαράσταση του μοντέλου λογιστικής παλινδρόμησης:

$$p(y = 1|x) = \sigma(w^T x + b)$$

όπου:

- $p(y = 1|x)$ είναι η πιθανότητα της μεταβλητής στόχου y να είναι 1 δεδομένων των εισαγόμενων χαρακτηριστικών x .
- σ είναι η συνάρτηση sigmoid.
- w είναι το διάνυσμα βαρών που αντιπροσωπεύει τους συντελεστές του γραμμικού μοντέλου.
- b είναι ο όρος μετατόπισης.
- x είναι το διάνυσμα των εισαγόμενων χαρακτηριστικών.
- w^T δηλώνει το μετασχηματισμένο (transposed) διάνυσμα βαρών.

Η συνάρτηση sigmoid ορίζεται ως:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

όπου $z = w^T x + b$.

Τα βήματα για την εκπαίδευση ενός μοντέλου λογιστικής παλινδρόμησης είναι τα εξής:

2.3.1.2 Κατάλληλες Εφαρμογές Ο αλγόριθμος είναι θεωρείται κατάλληλος για εφαρμογές δυαδικής ταξινόμησης όπου η μεταβλητή εξόδου μπορεί να κατηγοριοποιηθεί σε δύο κατηγορίες. Εδώ είναι μερικά παραδείγματα:

- Ανίχνευση ανεπιθύμητων μηνυμάτων: Κατηγοριοποίηση των emails ως ανεπιθύμητα ή μη.
- Πρόβλεψη εγκατάλειψης πελατών: Πρόβλεψη αν ένας πελάτης είναι πιθανό να σταματήσει να χρησιμοποιεί μια υπηρεσία.
- Πρόβλεψη έγκρισης δανείου: Πρόβλεψη αν μια αίτηση δανείου θα εγκριθεί.

Πίνακας 1: Βήματα του Αλγορίθμου Λογιστικής Παλινδρόμησης

Βήμα	Περιγραφή
Αρχικοποίηση	Αρχικοποιούμε τα βάρη w και την μετατόπιση b σε μικρές τυχαίες τιμές.
Προώθηση	Υπολογίζουμε την προβλεπόμενη πιθανότητα $\hat{y} = \sigma(w^T x + b)$ για κάθε παράδειγμα εκπαίδευσης.
Υπολογισμός Απώλειας	Υπολογίζουμε τη συνάρτηση απώλειας, συνήθως τη δυαδική διασταύρωση.
Οπισθοπροώθηση	Υπολογίζουμε τις κλίσεις της απώλειας σε σχέση με τα βάρη και τη μετατόπιση.
Ενημέρωση Παραμέτρων	Ενημερώνουμε τα βάρη και τη μετατόπιση χρησιμοποιώντας την καθοδική κλίση.
Επανάληψη	Επαναλαμβάνουμε τα βήματα 2-5 μέχρι να συγκλίνει το μοντέλο ή να φτάσει τον μέγιστο αριθμό επαναλήψεων.

2.3.1.3 Πλεονεκτήματα

- **Ερμηνευσιμότητα:** Οι συντελεστές του γραμμικού μοντέλου αντανακλούν άμεσα τη σημασία κάθε χαρακτηριστικού στην πρόβλεψη της μεταβλητής στόχου. Αυτό μας επιτρέπει να κατανοήσουμε πώς τα διάφορα χαρακτηριστικά συμβάλλουν στις προβλέψεις του μοντέλου.
- **Απλότητα:** Η λογιστική παλινδρόμηση είναι ένας σχετικά απλός αλγόριθμος που είναι εύκολο να υλοποιηθεί και να κατανοηθεί. Απαιτεί επίσης λιγότερη υπολογιστική ισχύ σε σύγκριση με ορισμένα άλλα μοντέλα.

- **Αποτελεσματικότητα:** Η λογιστική παλινδρόμηση μπορεί να είναι πολύ αποτελεσματική για τη διαχείριση μεγάλων συνόλων δεδομένων.

2.3.1.4 Τύποι Δεδομένων Ο αλγόριθμος μπορεί να διαχειριστεί τόσο αριθμητικά όσο και κατηγορηματικά δεδομένα. Ωστόσο, τα κατηγορηματικά δεδομένα πρέπει να προεπεξεργαστούν σε αριθμητικά χαρακτηριστικά πριν την εισαγωγή τους στο μοντέλο. Αυτό μπορεί να γίνει με τεχνικές όπως η κωδικοποίηση one-hot.

2.3.1.5 Προεπεξεργασία Εδώ είναι μερικά σημαντικά βήματα προεπεξεργασίας για τη λογιστική παλινδρόμηση:

- **Διαχείριση ελλειπόντων τιμών:** Οι ελλείπουσες τιμές μπορούν να συμπληρωθούν με τεχνικές όπως η συμπλήρωση με μέσο/διάμεσο ή να αφαιρεθούν αν το ποσοστό των ελλειπόντων δεδομένων είναι μικρό.
- **Κλιμάκωση χαρακτηριστικών:** Χαρακτηριστικά με διαφορετικές κλίμακες μπορούν να επηρεάσουν την απόδοση του μοντέλου. Η κλιμάκωση των χαρακτηριστικών σε ένα παρόμοιο εύρος μπορεί να βελτιώσει τη σύγκλιση του μοντέλου.

2.3.1.6 Περιορισμοί

- **Περιορισμένη σε δυαδική ταξινόμηση:** Η λογιστική παλινδρόμηση μπορεί να διαχειριστεί μόνο εργασίες δυαδικής ταξινόμησης. Για προβλήματα πολλαπλών κατηγοριών, χρειάζονται άλλα μοντέλα όπως η πολυωνυμική λογιστική παλινδρόμηση ή η στρατηγική one-vs-rest.

2.3.2 Random Forest

Το Δάσος Τυχαίων Δέντρων είναι ένας ευέλικτος αλγόριθμος εποπτευόμενης μάθησης που χρησιμοποιείται για εργασίες ταξινόμησης και παλινδρόμησης [19]. Συνδυάζει πολλά δέντρα απόφασης για να βελτιώσει την ακρίβεια των προβλέψεων και να μειώσει τον κίνδυνο υπερεκπαίδευσης.

2.3.2.1 Ο αλγόριθμος Ο αλγόριθμος Δάσους Τυχαίων Δέντρων δημιουργεί πολλαπλά δέντρα απόφασης από διαφορετικά υποσύνολα των δεδομένων εκπαίδευσης και στη συνέχεια συνδυάζει τις προβλέψεις τους [19]. Η τελική πρόβλεψη γίνεται με ψηφοφορία για προβλήματα ταξινόμησης ή με τον μέσο όρο των προβλέψεων για προβλήματα παλινδρόμησης [20].

Πίνακας 2: Βήματα του Αλγορίθμου Δάσους Τυχαίων Δέντρων

Βήμα	Περιγραφή
Επιλογή Υποσυνόλων	Δημιουργήστε πολλαπλά υποσύνολα δεδομένων από το αρχικό σύνολο με επαναδειγματοληψία (bootstrapping).
Δημιουργία Δέντρων Απόφασης	Εκπαιδεύστε ένα δέντρο απόφασης σε κάθε υποσύνολο δεδομένων.
Συνδυασμός Προβλέψεων	Συνδυάστε τις προβλέψεις από όλα τα δέντρα για να δώσετε την τελική πρόβλεψη.

2.3.2.2 Κατάλληλες Εφαρμογές

- Αναγνώριση Μοτίβων: Κατηγοριοποίηση εικόνων και ήχων.
- Ανίχνευση Απάτης: Πρόβλεψη απάτης σε συναλλαγές.
- Πρόβλεψη Ασθενειών: Πρόβλεψη της πιθανότητας εμφάνισης ασθενειών από ιατρικά δεδομένα.

2.3.2.3 Πλεονεκτήματα

- **Ανθεκτικότητα στην Υπερεκπαίδευση:** Η χρήση πολλαπλών δέντρων μειώνει την πιθανότητα υπερεκπαίδευσης.

- **Διαχείριση Ελλιπών Δεδομένων:** Τα δέντρα απόφασης μπορούν να διαχειριστούν ελλιπή δεδομένα, καθιστώντας το Δάσος Τυχαίων Δέντρων ανθεκτικό [20].
- **Υψηλή Ακρίβεια:** Ο συνδυασμός πολλαπλών δέντρων συνήθως οδηγεί σε καλύτερη απόδοση σε σχέση με μεμονωμένα δέντρα [21].

2.3.2.4 Τύποι Δεδομένων Ο αλγόριθμος μπορεί να διαχειριστεί τόσο αριθμητικά όσο και κατηγορηματικά δεδομένα. Τα κατηγορηματικά δεδομένα μπορούν να κωδικοποιηθούν χρησιμοποιώντας τεχνικές όπως η κωδικοποιήσεις που έχουμε αναφέρει π.χ. one-hot πριν την εισαγωγή τους στο μοντέλο [22].

2.3.2.5 Προεπεξεργασία

- **Διαχείριση Ελλειπόντων Τιμών:** Οι ελλείπουσες τιμές μπορούν να συμπληρωθούν με τεχνικές όπως η συμπλήρωση με μέσο/διάμεσο ή να αφαιρεθούν εάν το ποσοστό των ελλειπόντων δεδομένων είναι μικρό [22].
- **Κλιμάκωση Χαρακτηριστικών:** Παρόλο που το Δάσος Τυχαίων Δέντρων είναι λιγότερο ευαίσθητο σε χαρακτηριστικά με διαφορετικές κλίμακες, η κλιμάκωση των χαρακτηριστικών μπορεί να βελτιώσει την απόδοση του μοντέλου [22].

2.3.2.6 Περιορισμοί

- **Πολυπλοκότητα:** Το Δάσος Τυχαίων Δέντρων είναι πιο περίπλοκο και απαιτεί περισσότερους υπολογιστικούς πόρους από ένα μεμονωμένο δέντρο απόφασης.
- **Μειωμένη Ερμηνευσιμότητα:** Η ερμηνεία του μοντέλου μπορεί να είναι πιο δύσκολη λόγω του μεγάλου αριθμού δέντρων που συνδυάζονται [;]

2.3.3 Support Vector Machine

Η Υποστηρικτική Μηχανή Διανυσμάτων (Support Vector Machine, SVM) είναι ένας ισχυρός αλγόριθμος εποπτευόμενης μάθησης που χρησιμοποιείται για εργασίες ταξινόμησης και παλινδρόμησης. Ο SVM προσπαθεί να βρει το βέλτιστο υπερεπίπεδο που διαχωρίζει τις κατηγορίες των δεδομένων με τον μέγιστο περιθώριο.

2.3.3.1 Ο αλγόριθμος Ο αλγόριθμος SVM δημιουργήθηκε από τους Vapnik και Chervonenkis και λειτουργεί βρίσκοντας το υπερεπίπεδο που μεγιστοποιεί τον περιθώριο μεταξύ των διαφορετικών κατηγοριών [23]. Το υπερεπίπεδο αυτό καθορίζεται από ένα μικρό υποσύνολο των δεδομένων εκπαίδευσης, τα οποία ονομάζονται διανύσματα υποστήριξης [24].

Πίνακας 3: Βήματα του Αλγορίθμου SVM

Βήμα	Περιγραφή
Επιλογή Χαρακτηριστικών	Επιλέγουμε τα χαρακτηριστικά που θα χρησιμοποιηθούν για την εκπαίδευση του μοντέλου.
Επιλογή Υποκειμένων	Επιλέγουμε τα υποκείμενα δεδομένα που θα χρησιμοποιηθούν ως διανύσματα υποστήριξης.
Βελτιστοποίηση Υπερεπιπέδου	Υπολογίζουμε το υπερεπίπεδο που μεγιστοποιεί τον περιθώριο μεταξύ των κατηγοριών.

2.3.3.2 Πλεονεκτήματα

- **Υψηλή Ακρίβεια:** Προσφέρει υψηλή ακρίβεια στις προβλέψεις μεγιστοποιώντας το περιθώριο μεταξύ των δεδομένων. [23].
- **Αντοχή σε Υπερεκπαίδευση:** Ο SVM έχει καλές ιδιότητες γενίκευσης, ιδιαίτερα σε μικρά σύνολα δεδομένων [25].

- **Ευελιξία:** Μπορεί να χρησιμοποιηθεί με διαφορετικούς πυρήνες (kernels) για να προσαρμοστεί σε μη γραμμικά δεδομένα.

2.3.3.3 Κατάλληλες Εφαρμογές

- Αναγνώριση Προτύπων: Κατηγοριοποίηση εικόνων και αναγνώριση χειρογράφων.
- Βιοπληροφορική: Ταξινόμηση βιολογικών δεδομένων και ανάλυση γονιδίων.
- Ανίχνευση Απάτης: Πρόβλεψη και ανίχνευση απάτης σε συναλλαγές.

2.3.3.4 Τύποι Δεδομένων Ο αλγόριθμος μπορεί να διαχειριστεί τόσο αριθμητικά όσο και κατηγορηματικά δεδομένα. Τα κατηγορηματικά δεδομένα πρέπει να μετατραπούν σε αριθμητικά πριν την εισαγωγή τους στο μοντέλο.

2.3.3.5 Προεπεξεργασία Ακολουθούν μερικά σημαντικά βήματα προεπεξεργασίας για τον SVM:

- **Διαχείριση Ελλειπόντων Τιμών:** Οι ελλείπουσες τιμές μπορούν να συμπληρωθούν ή να αφαιρεθούν από το σύνολο δεδομένων.
- **Κλιμάκωση Χαρακτηριστικών:** Η κλιμάκωση των χαρακτηριστικών είναι σημαντική για την απόδοση του SVM, καθώς επηρεάζει την απόδοση του μοντέλου.

2.3.3.6 Περιορισμοί Αν και ο SVM αποτελεί έναν ισχυρό αλγόριθμο μηχανικής μάθησης με ευρεία εφαρμογή σε διάφορους τομείς. Η ικανότητά του να χειρίζεται τόσο γραμμικά όσο και μη γραμμικά δεδομένα, η υψηλή ακρίβεια και η ανθεκτικότητά του στην υπερεκπαίδευση το καθιστούν ένα πολύτιμο εργαλείο για την επίλυση προβλημάτων ταξινόμησης και παλινδρόμησης. Ωστόσο:

- **Πολυπλοκότητα:** Μπορεί να είναι υπολογιστικά απαιτητικός, ιδιαίτερα με μεγάλα σύνολα δεδομένων.
- **Δυσκολία Ερμηνείας:** Η ερμηνεία του μοντέλου μπορεί να είναι δύσκολη, ειδικά με μη γραμμικούς πυρήνες .

2.3.4 Naive Bayes

Ο Naive Bayes είναι ένας αλγόριθμος εποπτευόμενης μάθησης που χρησιμοποιείται για εργασίες ταξινόμησης. Βασίζεται στον θεώρημα του Bayes με την απλοποιητική υπόθεση της ανεξαρτησίας μεταξύ των χαρακτηριστικών [26].

2.3.4.1 Ο αλγόριθμος Ο αλγόριθμος Naive Bayes, βασισμένος στο θεώρημα Bayes, αποτελεί μια απλή αλλά ισχυρή μέθοδο ταξινόμησης δεδομένων. Η λειτουργία του βασίζεται στην υποθετική ανεξαρτησία των χαρακτηριστικών που περιγράφουν κάθε κατηγορία. Με απλά λόγια, υποθέτει πως η ύπαρξη μίας ιδιότητας σε μια κατηγορία δεν επηρεάζεται από την ύπαρξη άλλων ιδιοτήτων στην ίδια κατηγορία.

Πίνακας 4: Βήματα του Αλγορίθμου Naive Bayes

Βήμα	Περιγραφή
Υπολογισμός Πιθανοτήτων	Υπολογίζουμε τις πιθανότητες των διαφορετικών τάξεων.
Υπολογισμός Συνθήκης Πιθανοτήτων	Υπολογίζουμε τις συνθήκες πιθανότητες των χαρακτηριστικών δεδομένης της τάξης.
Πρόβλεψη	Χρησιμοποιούμε τις πιθανότητες για να προβλέψετε την τάξη ενός νέου δείγματος.

2.3.4.2 Πλεονεκτήματα

- **Απλότητα:** Είναι εύκολος στην κατανόηση και την υλοποίηση [27].

- **Ταχύτητα:** Είναι γρήγορος και αποδοτικός, ιδιαίτερα για μεγάλα σύνολα δεδομένων [28].
- **Χρήση με Μικρά Δεδομένα:** Μπορεί να έχει καλή απόδοση ακόμα και με μικρά σύνολα δεδομένων.

2.3.4.3 Κατάλληλες Εφαρμογές Ο Naive Bayes είναι κατάλληλος για πολλές εφαρμογές όπως:

- **Ανίχνευση Ανεπιθύμητων Μηνυμάτων:** Ταξινόμηση email ως ανεπιθύμητα ή μη.
- **Ανάλυση Συναισθήματος:** Ταξινόμηση κειμένων με βάση το συναίσθημα που εκφράζουν.
- **Ιατρική Διάγνωση:** Πρόβλεψη πιθανών ασθενειών με βάση τα συμπτώματα.

2.3.4.4 Τύποι Δεδομένων Ο αλγόριθμος μπορεί να διαχειριστεί τόσο αριθμητικά όσο και κατηγορηματικά δεδομένα. Τα κατηγορηματικά δεδομένα μπορούν να μετατραπούν σε αριθμητικά πριν την εισαγωγή τους στο μοντέλο.

2.3.4.5 Προεπεξεργασία Ακολουθούν μερικά σημαντικά βήματα προεπεξεργασίας για τον Naive Bayes:

- **Διαχείριση Ελλειπόντων Τιμών:** Οι ελλείπουσες τιμές μπορούν να συμπληρωθούν ή να αφαιρεθούν από το σύνολο δεδομένων.
- **Κλιμάκωση Χαρακτηριστικών:** Η κλιμάκωση των χαρακτηριστικών μπορεί να βελτιώσει την απόδοση του μοντέλου, αν και ο Naive Bayes είναι λιγότερο ευαίσθητος στις διαφορές κλίμακας [29].

2.3.4.6 Περιορισμοί

- **Απλοποιητική Υπόθεση Ανεξαρτησίας:** Ο αλγόριθμος υποθέτει ανεξαρτησία μεταξύ των χαρακτηριστικών, κάτι που σπάνια ισχύει στην πράξη [26].
- **Ευαισθησία στις Σπάνιες Κατηγορίες:** Ο Naive Bayes μπορεί να μην αποδίδει καλά όταν οι κατηγορίες έχουν πολύ λίγα δείγματα .

2.4 Μέθοδοι Αξιολόγησης Μοντέλου

Η αξιολόγηση αποτελεί θεμελιώδες κατά τη διαδικασία δημιουργία ενός μοντέλου μηχανικής μάθησης, καθώς μας επιτρέπει να ερμηνεύσουμε την απόδοση των μοντέλων και να επιλέξουμε το πλέον κατάλληλο για κάθε περίπτωση.

Πλήθος μετρικών έρχονται να φωτίσουν την αποτελεσματικότητα των μοντέλων, όπως:

- **Ακρίβεια (Accuracy):** Η συχνότητα με την οποία το μοντέλο προβλέπει σωστά.
- **Ανάκληση (Recall):** Το ποσοστό των πραγματικά θετικών περιπτώσεων που ταυτοποιούνται σωστά.
- **Ευστοχία (Precision):** Το ποσοστό των αρνητικών περιπτώσεων που ταυτοποιούνται σωστά.
- **F1-score (F1 Score):** Μέτρηση που συνδυάζει ακρίβεια και ανάκληση, προσφέροντας ισορροπημένη εικόνα.

Αναλύοντας το μοντέλο βάσει των παραπάνω μετρικών οδηγούμαστε στην επιλογή του ιδανικού εργαλείου για το πρόβλημα που καλούμαστε να αντιμετωπίσουμε. Στη συνέχεια, αναλύονται αναλυτικότερα οι προαναφερθέντες μετρικές απόδοσης.

2.4.1 Ακρίβεια (Accuracy)

Η ακρίβεια είναι το ποσοστό των σωστών προβλέψεων προς τον συνολικό αριθμό των προβλέψεων. Είναι μια κοινή μετρική που χρησιμοποιείται όταν τα δεδομένα είναι ισορροπημένα, δηλαδή δεν υπάρχει σημαντική διαφορά στον αριθμό των παρατηρήσεων μεταξύ των διαφορετικών κατηγοριών [30].

2.4.2 Ανάκληση (Recall)

Η ανάκληση, επίσης γνωστή ως ευαισθησία, μετρά την ικανότητα του μοντέλου να εντοπίζει με ακρίβεια τις θετικές περιπτώσεις. Είναι ιδιαίτερα χρήσιμη όταν είναι σημαντικό να μην χάνονται σημαντικά θετικά αποτελέσματα [31].

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

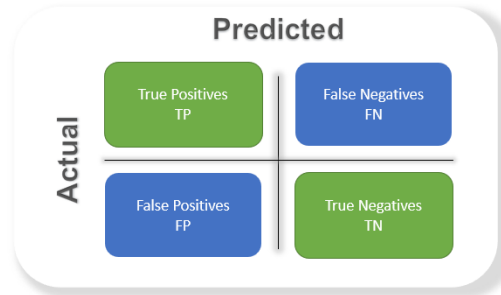
2.4.3 Ευστοχία (Precision)

Η ευστοχία (Precision) είναι ένα από τα βασικά μέτρα αξιολόγησης της απόδοσης ενός ταξινομητή, ειδικά σε προβλήματα δυαδικής ταξινόμησης. Η ευστοχία ορίζεται ως το ποσοστό των σωστά προβλεπόμενων θετικών περιπτώσεων προς όλες τις περιπτώσεις που προβλέφθηκαν ως θετικές [32].

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad (1)$$

όπου:

- Αληθώς Θετικά (True Positives, TP): Ο αριθμός των περιπτώσεων που ταξινομήθηκαν σωστά ως θετικές.
- Ψευδώς Θετικά (False Positives, FP): Ο αριθμός των περιπτώσεων που ταξινομήθηκαν λανθασμένα ως θετικές.



Σχήμα 3: Οπτική αναπαράσταση πίνακα σύγχυσης

2.4.4 F1-Score

Το F1-score είναι ο αρμονικός μέσος όρος της ακρίβειας και της ανάκλησης. Είναι χρήσιμο όταν υπάρχει ανισορροπία μεταξύ των κατηγοριών δεδομένων και επιθυμούμε ισορροπία μεταξύ ακρίβειας και ανάκλησης [33].

$$\text{F1-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

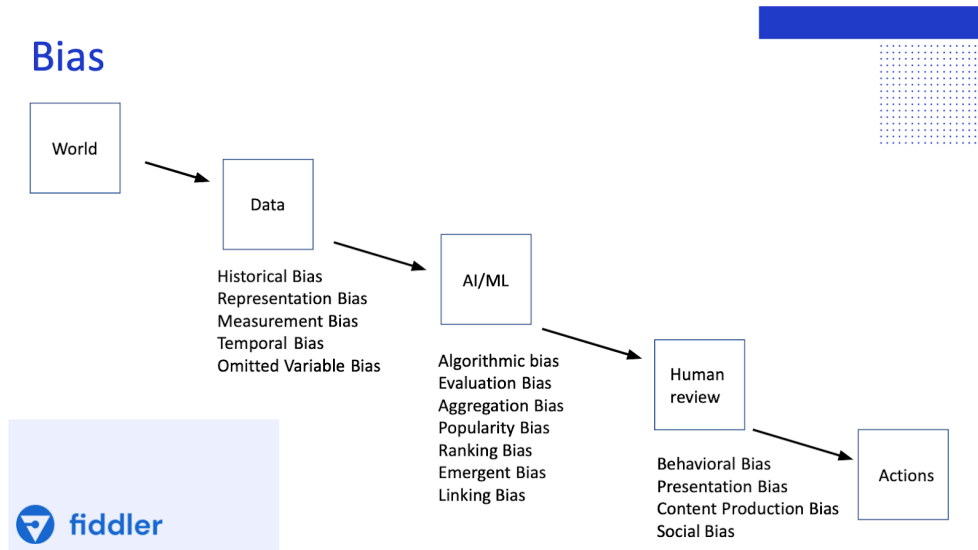
2.5 Μεροληψία στην Τεχνητή Νοημοσύνη

Η μεροληψία (bias) στην τεχνητή νοημοσύνη (AI) αποτελεί ένα από τα πλέον κρίσιμα ζητήματα, καθώς επηρεάζει άμεσα την αξιοπιστία και την ηθική χρήση των συστημάτων AI. Με την αυξανόμενη χρήση των συστημάτων AI σε τομείς όπως η υγειονομική περίθαλψη, η απασχόληση, η ποινική δικαιοσύνη και η πιστοληπτική αξιολόγηση, η ανησυχία για τη μεροληψία και την αμεροληψία των συστημάτων αυτών έχει ενταθεί.

Οι πηγές της μεροληψίας είναι ποικίλες και περιλαμβάνουν τη συλλογή δεδομένων, το σχεδιασμό αλγορίθμων και τις ανθρώπινες αποφάσεις [34]. Η μεροληψία μπορεί να οδηγήσει σε άδικο αποτελέσματα και να διακωλύσει υπάρχουσες ανισότητες, καθώς τα μοντέλα AI μπορεί να αναπαράγουν και να ενισχύουν κοινωνικά στερεότυπα.

Οι αρνητικές επιπτώσεις της μεροληψίας επηρεάζουν τόσο τα άτομα όσο και την κοινωνία συνολικά. Για την αντιμετώπιση της μεροληψίας, έχουν προταθεί διάφορες στρατηγικές μετριασμού, όπως η βελτίωση της ποιότητας των δεδομένων, η επιλογή

κατάλληλων μοντέλων και η επεξεργασία μετά την εκπαίδευση των μοντέλων.



Σχήμα 4: Σταδία πιθανής εμφάνισης Μεροληψίας

Υπάρχουν διάφοροι τύποι μεροληψίας που μπορούν να επηρεάσουν τα μοντέλα machine learning:

- **Μεροληψία Δεδομένων:**

- **Μεροληψία Επιλογής (Sampling Bias)** : Εμφανίζεται όταν τα δεδομένα εκπαίδευσης δεν είναι αντιπροσωπευτικά του πληθυσμού στον οποίο θα εφαρμοστεί το μοντέλο. Για παράδειγμα, αν ένα σύστημα αναγνώρισης προσώπου εκπαιδευτεί κυρίως σε εικόνες ανοιχτόχρωμων ατόμων, μπορεί να έχει κακή απόδοση σε σκουρόχρωμα άτομα [35].
- **Μεροληψία Μέτρησεων (Measurement Bias)** : Συμβαίνει όταν τα δεδομένα που συλλέγονται για εκπαίδευση ή αξιολόγηση περιέχουν ανακρίβειες ή συστηματικά λάθη. Για παράδειγμα, αν τα σφάλματα καταχώρισης δεδομένων είναι πιο συνηθισμένα για συγκεκριμένες δημογραφικές ομάδες, το μοντέλο μπορεί να μάθει να συνδέει αυτά τα λάθη με τις ίδιες τις ομάδες.

- **Αλγοριθμική Μεροληψία:** Αυτός ο τύπος μεροληψίας εμφανίζεται όταν το ίδιο το μοντέλο ή ο αλγόριθμος συμβάλλει σε μεροληπτικά αποτελέσματα. Για παράδειγμα, ορισμένοι αλγόριθμοι μπορεί να ευνοούν εγγενώς μια ομάδα έναντι άλλης αν δεν έχουν ρυθμιστεί σωστά ή αν ο σχεδιασμός τους δεν λαμβάνει υπόψη τις παραμέτρους δικαιοσύνης.

Η κατανόηση και ο μετριασμός της μεροληψίας είναι κρίσιμης σημασίας, καθώς τα μεροληπτικά μοντέλα μπορούν να διαιωνίσουν και ακόμη και να ενισχύσουν τις υπάρχουσες ανισότητες, οδηγώντας σε σημαντικές συνέπειες σε κρίσιμους τομείς όπως οι προσλήψεις, η χορήγηση δανείων, η ποινική δικαιοσύνη και η υγειονομική περίθαλψη [36, 37].

2.6 Αλγοριθμική μεροληψία και δικαιοσύνη στη μηχανική μάθηση

Η αλγοριθμική μεροληψία αποτελεί ένα σημαντικό ζήτημα στη μηχανική μάθηση, καθώς τα μοντέλα που εκπαιδεύονται σε δεδομένα μπορεί να αντικατοπτρίζουν και να ενισχύουν υπάρχουσες προκαταλήψεις στην κοινωνία. Αυτό μπορεί να οδηγήσει σε άδικες και άνισες αποφάσεις που επηρεάζουν αρνητικά μειονοτικές ομάδες.

Για παράδειγμα, ένα μοντέλο που εκπαιδεύεται σε δεδομένα για την πρόβλεψη της εγκληματικότητας μπορεί να είναι πιο πιθανό να ταξινομήσει άτομα από μειονοτικές ομάδες ως πιθανούς εγκληματίες, ακόμα κι αν δεν έχουν παραβεί τον νόμο.

Η δικαιοσύνη στη μηχανική μάθηση εστιάζει στην ανάπτυξη αλγορίθμων που είναι δίκαιοι, αμερόληπτοι και δεν διακρίνουν εις βάρος συγκεκριμένων ομάδων. Αυτό περιλαμβάνει την αναγνώριση και την αντιμετώπιση πιθανών πηγών μεροληψίας στα δεδομένα εκπαίδευσης, τον σχεδιασμό αλγορίθμων που είναι ανθεκτικοί στη μεροληψία και την ανάπτυξη τεχνικών για την αξιολόγηση της δικαιοσύνης των μοντέλων μηχανικής μάθησης.

Είναι σημαντικό να λαμβάνουμε υπόψη την αλγοριθμική μεροληψία και να υιοθετούμε πρακτικές για την προώθηση της δικαιοσύνης στη μηχανική μάθηση, καθώς τα

μοντέλα μηχανικής μάθησης ολοένα και περισσότερο επηρεάζουν τις ζωές μας.

2.6.1 Δικαιοσύνη στην τεχνητή νοημοσύνη

Η δικαιοσύνη στην τεχνητή νοημοσύνη περιλαμβάνει την εξασφάλιση ότι τα μοντέλα μηχανικής μάθησης αντιμετωπίζουν όλους τους ανθρώπους και τις ομάδες με δίκαιο τρόπο, χωρίς διακρίσεις ή προτιμήσεις. Υπάρχουν διάφοροι ορισμοί και μετρικές για τη δικαιοσύνη, που αντικατοπτρίζουν διαφορετικές προοπτικές και στόχους:

- **Δημογραφική Ισοτιμία:** Ένα μοντέλο ικανοποιεί τη δημογραφική ισοτιμία αν η πιθανότητα ενός θετικού αποτελέσματος είναι η ίδια για διαφορετικές δημογραφικές ομάδες. Για παράδειγμα, ένας αλγόριθμος πρόσληψης θα πρέπει να επιλέγει υποψηφίους από διαφορετικές φυλετικές ομάδες με παρόμοια ποσοστά, υποθέτοντας ίσα προσόντα [38].
- **Ισότητα Ευκαιριών:** Αυτό το κριτήριο δικαιοσύνης απαιτεί τα άτομα σε διαφορετικές ομάδες που είναι εξίσου καταρτισμένα να έχουν ίσες πιθανότητες να επιλεγούν. Για παράδειγμα, ένα μοντέλο πιστωτικής αξιολόγησης θα πρέπει να εγκρίνει δάνεια για καταρτισμένους αιτούντες με ίσα ποσοστά ανεξάρτητα από το φύλο τους.
- **Ισοτιμία Αποτελεσμάτων:** Ένα μοντέλο ικανοποιεί την ισοτιμία αποτελεσμάτων αν έχει ίσα ποσοστά αληθινών θετικών και ψευδών θετικών για διαφορετικές δημογραφικές ομάδες. Αυτό σημαίνει ότι η ακρίβεια του μοντέλου είναι συνεπής μεταξύ των ομάδων, μειώνοντας την πιθανότητα δυσανάλογα υψηλών ψευδών θετικών ή ψευδών αρνητικών για οποιαδήποτε συγκεκριμένη ομάδα.

Η επίτευξη δικαιοσύνης στα συστήματα τεχνητής νοημοσύνης είναι πρόκληση, καθώς διαφορετικές μετρικές δικαιοσύνης μπορεί να έρχονται σε σύγκρουση μεταξύ τους και η βελτιστοποίηση για μία μπορεί να οδηγήσει σε συμβιβασμούς σε μια άλλη. Επιπλέον, η δικαιοσύνη πρέπει να λαμβάνεται υπόψη στο πλαίσιο της συγκεκριμένης εφαρμογής και του κοινωνικού αντίκτυπου των αποφάσεων του μοντέλου [39, 40].

2.6.2 Σημασία της αλγοριθμικής και Δικαιοσύνης στη μηχανική μάθηση

Η σημασία της αντιμετώπισης της αλγοριθμικής μεροληψίας και της εξασφάλισης δικαιοσύνης στη μηχανική μάθηση δεν μπορεί να υπερεκτιμηθεί. Τα μεροληπτικά μοντέλα μπορούν να οδηγήσουν σε άδικη μεταχείριση ατόμων, διαιωνίζοντας και ενισχύοντας τις κοινωνικές ανισότητες. Αυτό είναι ιδιαίτερα ανησυχητικό σε εφαρμογές υψηλού κινδύνου όπως η ποινική δικαιοσύνη, η υγειονομική περίθαλψη, η χρηματοοικονομική και η απασχόληση.

Ποινική Δικαιοσύνη: Στο σύστημα ποινικής δικαιοσύνης, τα μεροληπτικά εργαλεία αξιολόγησης κινδύνου μπορούν να οδηγήσουν σε δυσανάλογα αυστηρές ποινές για τις μειονότητες. Μελέτες έχουν δείξει ότι ορισμένοι αλγόριθμοι που χρησιμοποιούνται για την πρόβλεψη των ποσοστών επανάληψης εγκλημάτων είναι μεροληπτικοί κατά των μαύρων κατηγορουμένων, οδηγώντας σε υψηλότερα ποσοστά ψευδών θετικών σε σύγκριση με τους λευκούς κατηγορούμενους [41].

Υγειονομική Περίθαλψη: Στην υγειονομική περίθαλψη, τα μεροληπτικά μοντέλα μπορούν να οδηγήσουν σε άνισα επίπεδα πρόσβασης στη θεραπεία και τη φροντίδα. Για παράδειγμα, ένα μοντέλο που εκπαιδεύτηκε κυρίως σε δεδομένα από άνδρες ασθενείς μπορεί να υποδιαγνώσει καταστάσεις που εμφανίζονται διαφορετικά σε γυναίκες ασθενείς, οδηγώντας σε υποβέλτιστη φροντίδα για τις γυναίκες [42].

Χρηματοοικονομική: Στον χρηματοοικονομικό τομέα, τα μεροληπτικά μοντέλα αξιολόγησης πιστοληπτικής ικανότητας μπορούν να αρνούνται δάνεια σε καταρτισμένους αιτούντες βάσει της φυλής ή της εθνικότητάς τους. Αυτή η διάκριση όχι μόνο επηρεάζει τις ευκαιρίες των ατόμων αλλά και διαιωνίζει τις οικονομικές ανισότητες [43].

Απασχόληση: Στις προσλήψεις, οι μεροληπτικοί αλγόριθμοι μπορούν να μειονεκτούν συγκεκριμένες δημογραφικές ομάδες, διαιωνίζοντας τις ανισότητες στο χώρο εργασίας. Για παράδειγμα, ένας αλγόριθμος πρόσληψης που εκπαιδεύτηκε σε βιογραφικά κυρίως από ένα φύλο ή μια φυλετική ομάδα μπορεί ακούσια να ευνοήσει υποψηφίους από αυτήν την ομάδα, υπονομεύοντας τις προσπάθειες για πολυμορφία και ένταξη [44]

Η αντιμετώπιση αυτών των ζητημάτων απαιτεί μια πολυδιάστατη προσέγγιση, συμπεριλαμβανομένης της ανάπτυξης και εφαρμογής μετρικών δικαιοσύνης, της χρήσης αλγορίθμων μετριασμού μεροληψίας και της καθιέρωσης νομικών και ηθικών κατευθυντήριων γραμμών. Εργαλεία όπως το IBM's AI Fairness 360 που χρησιμοποιείται στην παρούσα εργασία αλλά και άλλες βιβλιοθήκες δικαιοσύνης όπως Fairlearn, VerifyML κ.α. παρέχουν πρακτικές λύσεις για την ανίχνευση και τον μετριασμό της μεροληψίας, προσφέροντας μια σειρά από μετρικές και αλγορίθμους που μπορούν να ενσωματωθούν στη διαδικασία μηχανικής μάθησης για την προώθηση της δικαιοσύνης [45].

Επιπλέον, νομικά πλαίσια όπως το Local Law 144 of 2021, που επιβάλλεται από το NYC Department of Consumer and Worker Protection (DCWP), απαιτούν διαφάνεια και δικαιοσύνη στα αυτοματοποιημένα συστήματα απόφασης. Η συμμόρφωση με τέτοιους κανονισμούς εξασφαλίζει ότι οι οργανισμοί είναι υπεύθυνοι για τα αποτελέσματα των AI συστημάτων τους και ότι τα άτομα προστατεύονται από τις διακριτικές πρακτικές.

Με την κατανόηση και την αντιμετώπιση της αλγοριθμικής μεροληψίας, μπορούμε να κατασκευάσουμε πιο δίκαια, διαφανή και εξηγήσιμα συστήματα που μπορούν να εφαρμοστούν και να μην αδικούν κοινωνικές ομάδες.

2.7 Μετρικές Δικαιοσύνης

Οι μετρικές δικαιοσύνης (fairness metrics) χρησιμοποιούνται για την αξιολόγηση και την εξασφάλιση της δικαιοσύνης στα μοντέλα μηχανικής μάθησης και τεχνητής νοημοσύνης. Είναι βασισμένες σε διάφορες θεωρητικές και πρακτικές αρχές που αποσκοπούν στην ποσοτικοποίηση της δικαιοσύνης και της μεροληψίας στα μοντέλα μηχανικής μάθησης. Αυτές οι αρχές προέρχονται από διάφορα επιστημονικά πεδία, όπως η στατιστική, η πληροφορική, η κοινωνιολογία και η νομική. Στη συνέχεια θα αναλύσουμε κάποιες από αυτές όπως η μέση διαφορά (mean difference), (disparate impact), η διαφορά ίσων ευκαιριών (equal opportunity difference), η διαφορά μέσων όρων (average

odds difference), και ο δείκτης Theil (Theil index).

2.7.1 Μέση Διαφορά (Mean Difference)

Η μέση διαφορά μετρά τη διαφορά στην απόδοση του μοντέλου μεταξύ των διαφορετικών ομάδων. Υπολογίζεται ως η μέση διαφορά των προβλεπόμενων τιμών από τις πραγματικές τιμές ανά ομάδα. Χρησιμοποιείται συνήθως σε περιπτώσεις όπου θέλουμε να διασφαλίσουμε ότι οι προβλέψεις είναι ισορροπημένες μεταξύ διαφορετικών δημογραφικών ομάδων. Μια χαμηλή τιμή υποδεικνύει μεγαλύτερη δικαιοσύνη, καθώς σημαίνει ότι το μοντέλο προβλέπει με παρόμοια ακρίβεια για όλες τις ομάδες [46].

$$\text{Mean Difference} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)$$

Table 5: Mean Difference

Όρος	Περιγραφή
\hat{y}_i	Προβλεπόμενη τιμή για την παρατήρηση i
y_i	Πραγματική τιμή για την παρατήρηση i
n	Συνολικός αριθμός παρατηρήσεων

2.7.2 Disparate Impact

Η μετρική υπολογίζει το λόγο των θετικών προβλέψεων που λαμβάνουν οι προστατευόμενες και οι μη προστατευόμενες ομάδες. Σκοπός της είναι να εντοπίσει τυχόν άνισες μεταχειρίσεις και διακρίσεις που ενδεχομένως υφίστανται οι προστατευόμενες ομάδες. Ένας λόγος κοντά στο 1 υποδεικνύει ισότιμη μεταχείριση και δίκαιη διαδικασία. Όσο ο λόγος των δύο ομάδων αποκλείνει από το 1, τότε αυτό αποτελεί άνιση μεταχείριση των προστατευόμενων ομάδων.

$$\text{Disparate Impact} = \frac{\Pr(\hat{Y} = 1|A = 0)}{\Pr(\hat{Y} = 1|A = 1)}$$

Table 6: Disparate Impact

Όρος	Περιγραφή
$\Pr(\hat{Y} = 1 A = 0)$	Πιθανότητα θετικής πρόβλεψης για μη προστατευμένη ομάδα
$\Pr(\hat{Y} = 1 A = 1)$	Πιθανότητα θετικής πρόβλεψης για προστατευμένη ομάδα

2.7.3 Διαφορά Ίσων Ευκαιριών (Equal Opportunity Difference)

Η διαφορά ίσων ευκαιριών μετρά τη διαφορά στις ευαισθησίες (true positive rates) μεταξύ των προστατευμένων και μη προστατευμένων ομάδων. Αυτή η μετρική είναι ιδιαίτερα χρήσιμη όταν θέλουμε να διασφαλίσουμε ότι το μοντέλο μας αντιμετωπίζει όλες τις ομάδες με τον ίδιο τρόπο όσον αφορά την αναγνώριση των θετικών περιπτώσεων [47].

$$\text{Equal Opportunity Difference} = \text{TPR}_0 - \text{TPR}_1$$

Table 7: Equal Opportunity Difference

Όρος	Περιγραφή
TPR_0	Ευαισθησία (true positive rate) για μη προστατευμένη ομάδα
TPR_1	Ευαισθησία (true positive rate) για προστατευμένη ομάδα

2.7.4 Διαφορά Μέσων Όρων (Average Odds Difference)

Η διαφορά μέσων όρων μετρά τη διαφορά στις ευαισθησίες και στις ειδικότητες (true negative rates) μεταξύ των προστατευμένων και μη προστατευμένων ομάδων. Αυτή η μετρική είναι χρήσιμη για την αξιολόγηση της συνολικής απόδοσης του μοντέλου μεταξύ των ομάδων.

$$\text{Average Odds Difference} = \frac{1}{2}[(\text{TPR}_0 - \text{TPR}_1) + (\text{TNR}_0 - \text{TNR}_1)]$$

Table 8: Average Odds Difference

Όρος	Περιγραφή
TPR_0	Ευαισθησία (true positive rate) για μη προστατευμένη ομάδα
TPR_1	Ευαισθησία (true positive rate) για προστατευμένη ομάδα
TNR_0	Ειδικότητα (true negative rate) για μη προστατευμένη ομάδα
TNR_1	Ειδικότητα (true negative rate) για προστατευμένη ομάδα

2.7.5 Δείκτης Theil (Theil Index)

Ο δείκτης Theil είναι ένα στατιστικό μέτρο οικονομικής ανισότητας που αναπτύχθηκε από τον οικονομολόγο Henri Theil. Χρησιμοποιείται για την ποσοτικοποίηση των διαφορών στην κατανομή εισοδήματος ή πόρων μέσα σε έναν πληθυσμό. Είναι ιδιαίτερα χρήσιμος γιατί παρέχει μια σαφή, μαθηματική αναπαράσταση της ανισότητας, καθιστώντας το πολύτιμο εργαλείο που μπορεί να διακρίνει μεταξύ διαφορετικών πηγών ανισότητας, όπως η ανισότητα μεταξύ ομάδων και εντός ομάδων, που είναι χρήσιμο για τον εντοπισμό των συγκεκριμένων παραγόντων που συμβάλλουν στη συνολική ανισότητα [48, 49].

Ο δείκτης Theil δίνεται από:

$$T = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i}{\bar{y}} \log \left(\frac{y_i}{\bar{y}} \right) \right)$$

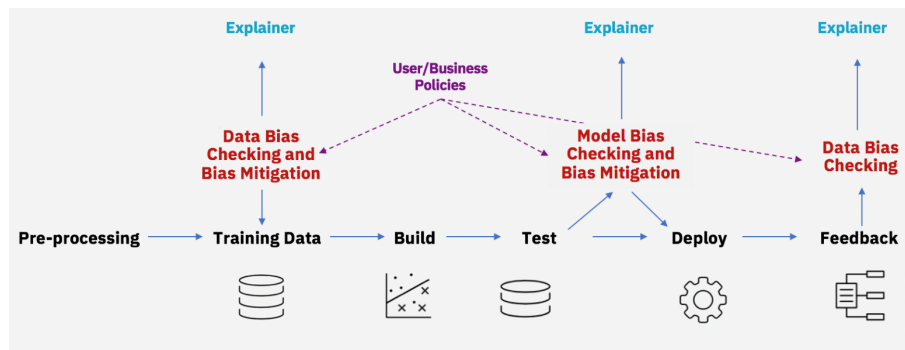
όπου N είναι ο αριθμός των ατόμων, y_i είναι το εισόδημα του ατόμου i , και \bar{y} είναι το μέσο εισόδημα.

Table 9: Βήματα Υπολογισμού Δείκτη Theil (Theil Index Calculation Steps)

Βήμα	Περιγραφή
1.	Προσδιορισμός Πληθυσμού.
2.	Συλλογή Δεδομένων Εισοδήματος.
3.	Υπολογισμός Συνολικού Εισοδήματος .
4.	Υπολογισμός Μέσου Εισοδήματος.
5.	Υπολογισμός Αναλογίας Ατομικού Εισοδήματος .
6.	Υπολογισμός Λογαρίθμου Αναλογιών κάθε ατόμου.
7.	Υπολογισμός Σταθμισμένου Λογαρίθμου.
8.	Άθροιση Σταθμισμένων Λογαρίθμων .
9.	Υπολογισμός Δείκτη.

2.8 Αλγόριθμοι Μείωσης Μεροληψίας

Η μείωση της μεροληψίας στα μοντέλα μηχανικής μάθησης είναι ζωτικής σημασίας για τη διασφάλιση της δικαιοσύνης και της ισότητας στις προβλέψεις. Οι αλγόριθμοι μείωσης μεροληψίας μπορούν να εφαρμοστούν σε τρία στάδια: προεπεξεργασία (pre-processing), κατά την επεξεργασία (in-processing), και μετά την επεξεργασία (post-processing). Η επιλογή του κατάλληλου αλγορίθμου εξαρτάται από τη φύση των δεδομένων, το στάδιο της ανάπτυξης του μοντέλου, και τις συγκεκριμένες απαιτήσεις της εφαρμογής.



Σχήμα 5: Στάδια Εφαρμογής των αλγορίθμων μείωσης μεροληψίας

- **Pre-processing Algorithms:** Οι αλγόριθμοι προεπεξεργασίας αντιμετωπίζουν τη μεροληψία στα δεδομένα εκπαίδευσης πριν από την εκπαίδευση του μοντέλου. Αυτές οι μέθοδοι τροποποιούν το σύνολο δεδομένων έτσι ώστε οι προβλέψεις του μοντέλου να είναι πιο δίκαιες [50].
- **In-Processing Algorithms :** Αλγόριθμοι οι οποίοι εφαρμόζονται στο μοντέλο κατά τη διαδικασία της εκπαίδευσης του.
- **Post-Processing Algorithms.:** Οι αλγόριθμοι μετά την επεξεργασία εφαρμόζονται στις προβλέψεις του μοντέλου αφού το μοντέλο έχει ήδη εκπαιδευτεί. Αυτές οι τεχνικές τροποποιούν τις τελικές προβλέψεις για να βελτιώσουν τη δικαιοσύνη και δεν εμπλέκονται με κανένα από τα προηγούμενα στάδια.

2.8.1 Αναπροσαρμογή Βαρών (Reweighing)

Η τεχνική της αναπροσαρμογής των βαρών (Reweighing) χρησιμοποιεί στην προεπεξεργασία (Pre-processing) και περιλαμβάνει την αναπροσαρμογή των βαρών των παρατηρήσεων στο σύνολο δεδομένων εκπαίδευσης για να αντισταθμιστεί η μεροληψία. Αυτή η μέθοδος διασφαλίζει ότι οι προστατευμένες και μη προστατευμένες ομάδες έχουν παρόμοια επίδραση στην εκπαίδευση του μοντέλου [51].

$$\text{Weight}_{\text{new}} = \frac{\text{Total instances}}{\text{Instances of group}}$$

Table 10: Reweighing

Όρος	Περιγραφή
$\text{Weight}_{\text{new}}$	Νέο βάρος παρατήρησης
Total instances	Συνολικός αριθμός παρατηρήσεων
Instances of group	Αριθμός παρατηρήσεων της ομάδας

2.8.2 Adversarial Debiasing

Ο αλγόριθμος (Adversarial Debiasing) είναι μια μέθοδος κατά την επεξεργασία (In-processing) που χρησιμοποιεί έναν αντίπαλο για να μειώσει τη μεροληψία στις προβλέψεις. Ο αλγόριθμος προσπαθεί να μεγιστοποιήσει την ακρίβεια των προβλέψεων ενώ ταυτόχρονα ελαχιστοποιεί τη δυνατότητα του αντιπάλου να προσδιορίσει τα προστατευμένα χαρακτηριστικά [52].

$$\min (\text{Loss}_{\text{classifier}} - \lambda \cdot \text{Loss}_{\text{adversary}})$$

2.8.3 Calibrated Equalized Odds

Ο αλγόριθμος (Calibrated Equalized Odds) αποτελεί μέθοδο που εφαρμόζεται μετά την επεξεργασία (Post-processing) που χρησιμοποιεί την βαθμονόμηση των εξόδων

Table 11: Adversarial Debiasing

Όρος	Περιγραφή
$\text{Loss}_{\text{classifier}}$	Συνάρτηση απώλειας ταξινόμητη
$\text{Loss}_{\text{adversary}}$	Συνάρτηση απώλειας αντιπάλου
λ	Παράγοντας βαρύτητας

του ταξινόμητη για να βρει πιθανότητες με τις οποίες να αλλάξει τις ετικέτες εξόδου, διασφαλίζοντας ίσες πιθανότητες σφάλματος μεταξύ των ομάδων.

$$\min \left(\text{Loss}_{\text{calibrated}} + \lambda \cdot (|\text{FPR}_0 - \text{FPR}_1| + |\text{TPR}_0 - \text{TPR}_1|) \right)$$

Table 12: Calibrated Equalized Odds

Όρος	Περιγραφή
$\text{Loss}_{\text{calibrated}}$	Συνάρτηση απώλειας βαθμονόμησης
FPR_0	Ψευδώς θετικά ποσοστά για μη προστατευμένη ομάδα
FPR_1	Ψευδώς θετικά ποσοστά για προστατευμένη ομάδα
TPR_0	Πραγματικά θετικά ποσοστά για μη προστατευμένη ομάδα
TPR_1	Πραγματικά θετικά ποσοστά για προστατευμένη ομάδα
λ	Παράγοντας βαρύτητας

3 Τοπικός Νόμος 144 του 2021

Ο σύγχρονος οργανωσιακός κόσμος υιοθετεί ολοένα και περισσότερο εργαλεία Τεχνητής Νοημοσύνης για βελτιστοποίηση των εσωτερικών διαδικασιών, συμπεριλαμβανομένων και των λειτουργιών Ανθρώπινου Δυναμικού. Η αξιοποίηση ΤΝ για λήψη αποφάσεων πρόσληψης, απόλυσης ή προαγωγής φέρνει στο προσκήνιο εργασιακά ζητήματα και θέτει σε εφαρμογή νομοθεσίες περί ιδιωτικότητας, όπως ο Τοπικός Νόμος 144 της Νέας Υόρκης (NYC 144), που επιβάλλει ‘Έλεγχο Αμεροληψίας’ σε Αυτόματα

Εργαλεία Λήψης Αποφάσεων Απασχόλησης (AEDT) [53]. Το παρόν κεφάλαιο εστιάζει στον NYC 144 και τις απαιτήσεις του.

3.1 Ανάλυση Τοπικού Νόμου 144 του 2021

Ο Τοπικός Νόμος 144 του 2021, που εφαρμόστηκε από το Τμήμα Προστασίας Καταναλωτών και Εργαζομένων της Νέας Υόρκης (DCWP) και είναι σε ισχύ από τον Ιανουάριο του έτους 2023, είναι μια πρωτοποριακή ρύθμιση με στόχο τη μείωση της μεροληψίας στα Αυτοματοποιημένα Εργαλεία Λήψης Αποφάσεων για Προσλήψεις (AEDTs). Ο νόμος απαιτεί από τους εργοδότες και τις υπηρεσίες απασχόλησης να διενεργούν ετήσιους ελέγχους μεροληψίας στα AEDTs και να δημοσιοποιούν αυτούς τους ελέγχους, εξασφαλίζοντας διαφάνεια και λογοδοσία στις πρακτικές προσλήψεων [54].

3.2 Μεροληψία Διασταυρούμενων Χαρακτηριστικών

Ο κανόνας των τεσσάρων πέμπτων (four/fifths rule) είναι ένα σημαντικό εργαλείο για την αξιολόγηση της αλγοριθμικής δικαιοσύνης. Σύμφωνα με αυτόν τον κανόνα, μια συγκεκριμένη πρακτική θεωρείται ότι έχει disparate impact εάν το ποσοστό επιτυχίας μιας προστατευόμενης ομάδας είναι λιγότερο από το 80 τοις εκατό του ποσοστού επιτυχίας της ομάδας με την υψηλότερη επίδοση. Αυτή η μετρική χρησιμοποιείται για να αξιολογήσει αν υπάρχει ανισότητα στα αποτελέσματα μιας αλγοριθμικής απόφασης μεταξύ δύο διαφορετικών ομάδων, όπως ορίζεται από τον Title VII of the Civil Rights Act of 1964. Η μετρική disparate impact επιτρέπει την αναγνώριση ανισοτήτων που δεν είναι άμεσα εμφανείς αλλά προκύπτουν από την εφαρμογή του αλγόριθμου [55, 56].

Ωστόσο, αυτή η προσέγγιση δεν ήταν αρκετή για να διασφαλίσει την πλήρη δικαιοσύνη και αμεροληψία των αλγορίθμων. Με την ψήφιση του Τοπικού Νόμου 144 του 2021 στη Νέα Υόρκη, εισήχθη η έννοια της διασταυρούμενων (intersectional) χαρακτηριστικών μεροληψίας, η οποία εξετάζει τα διάφορα χαρακτηριστικά των datasets σε συνδυασμό και όχι μεμονωμένα. Αυτό σημαίνει ότι πρέπει να συνεχίζει να υπάρχει ο παραπάνω περιορισμός αλλά να εφαρμόζεται με βάση διασταυρούμενα χαρακτηριστι-

κά (intersectional attributes) . Ο νόμος αυτός επιδιώκει να εξαλείψει τη μεροληψία που μπορεί να προκύψει όταν ένας αλγόριθμος ευνοεί ή δυσχεραίνει ομάδες με βάση συνδυασμούς χαρακτηριστικών όπως το φύλο και η φυλή παραδείγματος χάρη. [57, 58].

Η προσέγγιση αυτή αναγνωρίζει ότι οι άνθρωποι δεν ανήκουν μόνο σε μία κατηγορία (π.χ. φύλο ή φυλή), αλλά σε πολλές ταυτόχρονα, και ότι η δίκαιη αντιμετώπιση πρέπει να λαμβάνει υπόψη αυτές τις πολυπλοκότητες[59, 60].

3.3 Αντιμετώπιση της Διασταυρούμενης Μεροληψίας

Ο Τοπικός Νόμος 144 εστιάζει στις ομάδες διασταυρούμενων χαρακτηριστικών, κάτι που είναι ιδιαίτερα κρίσιμο για την κατανόηση του πώς η μεροληψία μπορούν να επηρεάσουν δυσανάλογα τα άτομα που ανήκουν σε πολλαπλές περιθωριοποιημένες ομάδες. Τα παραδοσιακά μέτρα κατά των διακρίσεων συχνά αποτυγχάνουν να καταγράφουν τις σύνθετες μεροληψίες που αντιμετωπίζει, για παράδειγμα, μια μαύρη γυναίκα σε σύγκριση με έναν λευκό άνδρα. Οι διατάξεις του νόμου διασφαλίζουν ότι οι έλεγχοι μεροληψίας πρέπει να λαμβάνουν υπόψη διάφορες δημογραφικές ομάδες, συμπεριλαμβανομένων των διασταυρούμενων ταυτοτήτων, προωθώντας πιο δίκαιες πρακτικές προσλήψεων [61].

3.4 Αξιοποίηση του Εργαλείου AIF360

Το εργαλείο AI Fairness 360 (AIF360) είναι μια βιβλιοθήκη σχεδιασμένη για την ανίχνευση και μείωση της μεροληψίας στα μοντέλα μηχανικής μάθησης. Περιλαμβάνει εργαλεία για την αξιολόγηση της μεροληψίας σε διάφορες δημογραφικές ομάδες, παρέχοντας λεπτομερή ανάλυση που ευθυγραμμίζεται με τις απαιτήσεις του Τοπικού Νόμου 144. Εφαρμόζοντας το AIF360, οι οργανισμοί μπορούν να αξιολογούν τα AEDT τους για μεροληψίες τόσο ενάντια σε προνομιούχες ομάδες (π.χ. λευκοί άνδρες) όσο και σε μη προνομιούχες ομάδες (π.χ. γυναίκες σκούρου δέρματος) αποτελεσματικά [62].

3.5 Πρακτική Εφαρμογή και Προκλήσεις

Παρόλο που το εργαλείο AIF360 προσφέρει ισχυρές μετρικές για τον εντοπισμό μεροληψιών, υπάρχουν αρκετές πρακτικές προκλήσεις στην αποτελεσματική εφαρμογή αυτών των εργαλείων:

Απαιτήσεις Δεδομένων: Το AIF360 απαιτεί λεπτομερή δημογραφικά δεδομένα, τα οποία μπορεί να είναι δύσκολο να αποκτηθούν και να επαληθευτούν. Ο Τοπικός Νόμος 144 αντιμετωπίζει αυτό το ζήτημα απαιτώντας οι έλεγχοι μεροληψίας να αναφέρουν τον αριθμό των ατόμων που δεν παρείχαν δημογραφικά δεδομένα, εξασφαλίζοντας διαφάνεια στα δεδομένα που χρησιμοποιούνται για αυτούς τους ελέγχους [63].

Σύνθετα Μοντέλα: Η αποτελεσματικότητα του εργαλείου μπορεί να διαφέρει ανάλογα με την πολυπλοκότητα των AEDTs. Είναι κρίσιμο να εξερευνηθούν σενάρια όπου το AIF360 μπορεί να μην αποδίδει καλά, ιδιαίτερα σε μοντέλα με σύνθετες διαδικασίες λήψης αποφάσεων [64].

Μεροληψία κατά το Σχεδιασμό: Η μεροληψία κατά τον σχεδιασμό και την εφαρμογή των AEDTs μπορεί να παραμένουν. Ο αντιμετωπίζει έμμεσα αυτό το ζήτημα υπογραμμίζοντας την ανάγκη για εξωτερικούς ελέγχους, οι οποίοι μπορούν να παρέχουν αντικειμενική αξιολόγηση αυτών των εργαλείων [65].

3.6 Ενίσχυση της Διαφάνειας με την Επεξηγήσιμη Τεχνητή Νοημοσύνη

Η Καθηγήτρια Sarah Jones από το MIT υπογραμμίζει τη δυναμική της Επεξηγήσιμης Τεχνητής Νοημοσύνης (XAI) να συμπληρώνει τα εργαλεία εντοπισμού μεροληψίας. Οι τεχνικές XAI μπορούν να παρέχουν πληροφορίες για το πώς τα AEDTs καταλήγουν στις αποφάσεις τους, καθιστώντας τη διαδικασία προσλήψεων πιο διαφανή τόσο για τους εργοδότες όσο και για τους υποψήφιους. Αυτή η διαφάνεια είναι ουσιώδης για την οικοδόμηση εμπιστοσύνης και την εξασφάλιση συμμόρφωσης με τον Τοπικό Νόμο 144 [66].

3.7 Επιπτώσεις

Ο Τοπικός Νόμος 144 μπορεί να λειτουργήσει ως πρότυπο για παρόμοιες ρυθμίσεις. Η έμφαση του στη διαφάνεια, τη διασταυρούμενη ανάλυση και τους τακτικούς ελέγχους θέτει υψηλά πρότυπα για δίκαιες πρακτικές προσλήψεων. Ωστόσο, η άμεση εφαρμογή αυτού του νόμου σε διαφορετικά πολιτιστικά πλαίσια μπορεί να αντιμετωπίσει προκλήσεις, όπως οι διαφορετικές ορισμοί της μεροληψίας και τα διαφορετικά ρυθμιστικά τοπία. Η εξερεύνηση αυτών των ευρύτερων επιπτώσεων μπορεί να παρέχει μια πιο ολοκληρωμένη κατανόηση του παγκόσμιου αντίκτυπού του [67]. Η έμφαση του νόμου στους τακτικούς ελέγχους μεροληψίας και την λεπτομερή αναφορά δημογραφικών δεδομένων διασφαλίζει ότι οι αποχρώσεις των διασταυρούμενων μεροληψιών αντιμετωπίζονται, θέτοντας ένα προηγούμενο για μελλοντικές ρυθμίσεις στην Τεχνητή Νοημοσύνη και τις πρακτικές προσλήψεων [68].

4 Τεχνική Ανάλυση

Σε αυτό το κεφάλαιο, θα αναλύσουμε λεπτομερώς την αρχιτεκτονική και τη λειτουργικότητα του συστήματος που αναπτύξαμε. Θα παρουσιάσουμε το μοντέλο C4 για την απεικόνιση της αρχιτεκτονικής του συστήματος, καθώς και διαγράμματα που αποτυπώνουν τις περιπτώσεις χρήσης, την ευρωστία, και τις ακολουθίες των διαδικασιών. Επιπλέον, θα αναλύσουμε τα εργαλεία και τις βιβλιοθήκες που χρησιμοποιήθηκαν και θα περιγράψουμε τη δομή των αρχείων κώδικα.

4.1 Ανάλυση βάση του μοντέλου C4

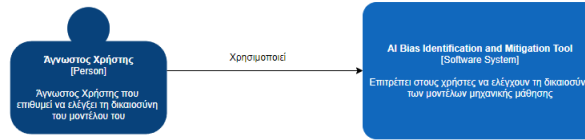
Η ανάλυση της αρχιτεκτονικής του συστήματος βασίζεται στο μοντέλο C4 , το οποίο δημιουργήθηκε από τον Simon Brown. Το C4 είναι μια μέθοδος απεικόνισης της διαδικασίας κατασκευής ενός συστήματος ή περιγραφής ενός ήδη υπάρχοντος συστήματος μέσω τεσσάρων επιπέδων: Context, Container, Component , και Code [69]. Τα διαγράμματα C4 αποτελούν μια απλοποιημένη εναλλακτική λύση σε σύγκριση με τα διαγράμματα UML, προσφέροντας δυνατότητα περιγραφής πολλαπλών συστημάτων [70].

- **Επίπεδο 1: Context**

Το επίπεδο Context παρέχει μια γενική εικόνα του συστήματος και των αλληλεπιδράσεών του με εξωτερικούς παράγοντες. Δείχνει το σύστημα που αναπτύσσεται, τους χρήστες του και τα άλλα συστήματα με τα οποία αλληλεπιδρά. Είναι το πιο αφαιρετικό επίπεδο και χρησιμεύει για την κατανόηση του περιβάλλοντος στο οποίο λειτουργεί το σύστημα.

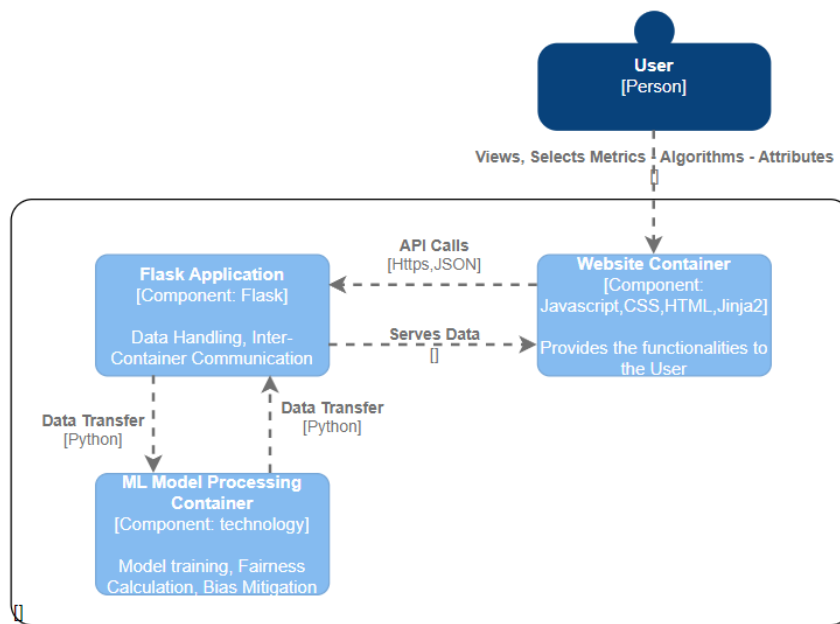
- **Επίπεδο 2: Container**

Το επίπεδο Container παρουσιάζει τη διάρθρωση του συστήματος σε επίπεδο container, δηλαδή, τις κύριες εφαρμογές και υπηρεσίες που συνθέτουν το σύστημα. Περιγράφει τα container και τον τρόπο με τον οποίο αυτά επικοινωνούν



Σχήμα 6: Διάγραμμα Περιβάλλοντος συστήματος

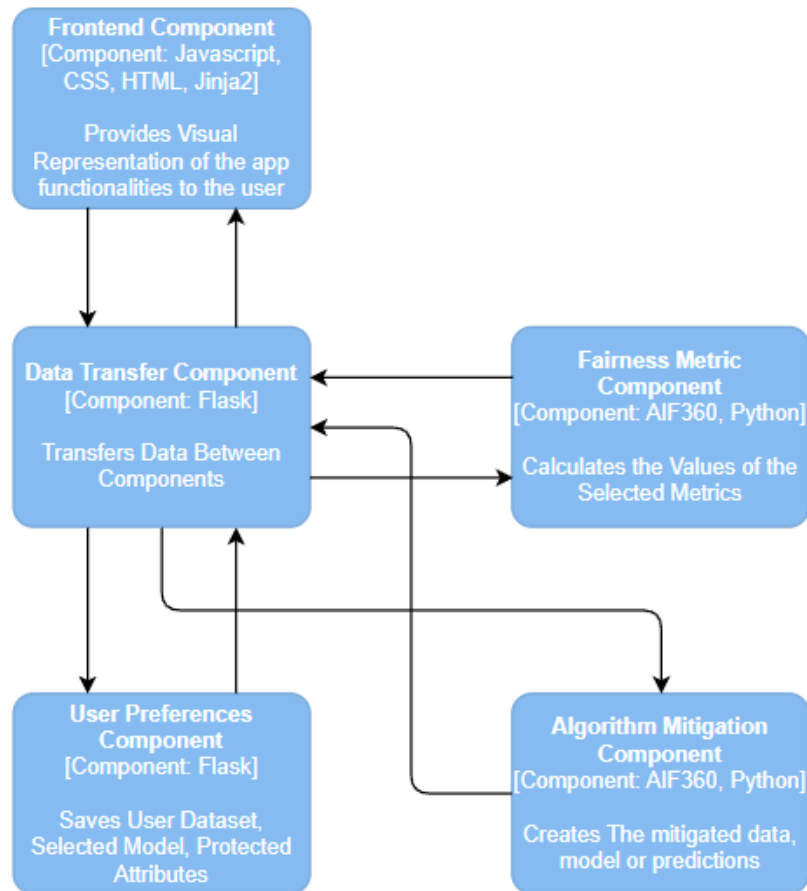
μεταξύ τους.



Σχήμα 7: Container Diagram του συστήματος

• Επίπεδο 3: Component

Το επίπεδο Component προσφέρει μια λεπτομερή απεικόνιση των εσωτερικών συστατικών κάθε container και τον τρόπο με τον οποίο αυτά αλληλεπιδρούν μεταξύ τους. Εστιάζει στα επιμέρους τμήματα του λογισμικού μέσα σε κάθε container και τον τρόπο που αυτά συνδέονται και συνεργάζονται για να υλοποιήσουν τις λειτουργίες του συστήματος.



Σχήμα 8: Component Diagram του συστήματος

- **Επίπεδο 4: Code**

Το επίπεδο Code παρέχει την πιο λεπτομερή περιγραφή του συστήματος, εστιάζοντας στον πηγαίο κώδικα. Αυτό το επίπεδο αναλύει τη δομή του κώδικα, περιγράφοντας τις κλάσεις, τις μεθόδους και τις σχέσεις τους. Είναι χρήσιμο για τους προγραμματιστές που πρέπει να κατανοήσουν την ακριβή υλοποίηση και να συντηρήσουν το σύστημα.

Πίνακας 13: Αρχεία Python και οι λειτουργίες τους

Όνομα Αρχείου	Λειτουργίες
<code>app.py</code>	Το κύριο αρχείο που διαχειρίζεται τις διαδρομές και τις βασικές λειτουργίες της εφαρμογής Flask. Περιλαμβάνει τα routes της εφαρμογής.
<code>proccess.py</code>	Περιέχει τις λειτουργίες επεξεργασίας των δεδομένων και εφαρμογής των αλγορίθμων μείωσης μεροληψίας. Διαχειρίζεται τη φόρτωση και προεπεξεργασία των dataset, καθώς και την εφαρμογή των αλγορίθμων από την βιβλιοθήκη Aif360.
<code>examples.py</code>	Περιέχει τη λειτουργικότητα των παραδειγμάτων χρήσης και βοηθητικές λειτουργίες για την κατανόηση του τρόπου λειτουργίας του εργαλείου και την εκπαίδευση των χρηστών.

Πίνακας 14: Συναρτήσεις του αρχείου proccess.py (Part 1)

Συνάρτηση	Περιγραφή
<code>prepare_data</code>	Προετοιμάζει τα δεδομένα για επεξεργασία.
<code>process_numerical_variables</code>	Επεξεργάζεται τις αριθμητικές μεταβλητές του dataset.
<code>encode_categorical_variables</code>	Κωδικοποιεί τις κατηγορικές μεταβλητές του dataset.
<code>train_model</code>	Εκπαιδεύει το μοντέλο μηχανικής μάθησης.
<code>calculate_model_metrics</code>	Υπολογίζει τις επιλεγμένες μετρικές δικαιοσύνης.
<code>convert_to_binary_label_dataset</code>	Μετατρέπει το dataset στην αντικείμενο της κλάσης <code>binary dataset</code> .
<code>get_binary_datasets</code>	Εξάγει δυαδικά datasets από τα αρχικά δεδομένα.
<code>calculate_fairness_metrics</code>	Υπολογίζει τις μετρικές δικαιοσύνης.
<code>get_privileged_group</code>	Εντοπίζει την προνομιούχο ομάδα στο dataset.
<code>calculate_standard_metrics</code>	Υπολογίζει μετρικές δικαιοσύνης για ομάδες βάση ενός χαρακτηριστικού.

Πίνακας 15: Συναρτήσεις του αρχείου proccess.py (Part 2)

Συνάρτηση	Περιγραφή
<code>group_results_by_metric</code>	Ομαδοποιεί αποτελέσματα με βάση τη μετρική δικαιοσύνης.
<code>fair_check</code>	Εκτελεί έλεγχο, αν το αποτέλεσμα της μετρικής δικαιοσύνης είναι εντός του ορισμένου threshold.
<code>train_and_evaluate</code>	Εκπαιδεύει και αξιολογεί το μοντέλο.
<code>mitigate</code>	Μειώνει τη μεροληψία στο dataset .
<code>apply_reweighing_and_train_model</code>	Εφαρμόζει τον αλγόριθμο Reweighing και εκπαιδεύει το μοντέλο.
<code>reweighing_result</code>	Παρέχει τα αποτελέσματα της τεχνικής Reweighing .
<code>apply_adversarial_debiasing_and_train_model</code>	Εφαρμόζει τον αλγόριθμο Adversarial Debiasing και εκπαιδεύει το μοντέλο.
<code>adversarial_debiasing_result</code>	Παρέχει τα αποτελέσματα του αλγόριθμου Adversarial Debiasing .
<code>apply_calibrated_eq_odds_and_train_model</code>	Εφαρμόζει τον αλγόριθμο Calibrated Equalized Odds και εκπαιδεύει το μοντέλο.

Πίνακας 16: Συναρτήσεις του αρχείου process.py (Part 3)

Συνάρτηση	Περιγραφή
<code>get_reason</code>	Παρέχει εξηγήσεις για τις αποφάσεις του μοντέλου.
<code>construct_metric_info</code>	Επιστρέφει πληροφορίες για τις μετρικές.
<code>get_ideal_fairness_value</code>	Εντοπίζει την ιδανική τιμή δικαιοσύνης της μετρικής.
<code>calibrated_eq_odds_result</code>	Παρέχει τα αποτελέσματα του αλγόριθμου Calibrated Equalized Odds .
<code>get_mitigated_results</code>	Εξάγει τα αποτελέσματα μετά την εφαρμογή των τεχνικών μείωσης μεροληψίας.
<code>wrap_response</code>	Επεξεργάζεται και μορφοποιεί κατάλληλα το response για επεξεργασία από το front-end.
<code>calculate_intersectional_metrics</code>	Υπολογίζει μετρικές για ομάδες διασταυρωμένων χαρακτηριστικών.

Πίνακας 17: Συναρτήσεις του αρχείου examples.py

Συνάρτηση	Περιγραφή
<code>create_attnsn_value</code>	Δημιουργεί χαρακτηριστικά και τιμές για παράδειγμα χρήσης.
<code>get_example_dataset</code>	Παρέχει dataset για τα παραδείγματα της εφαρμογής.
<code>get_groups</code>	Εξάγει τις ομάδες από το παράδειγμα που έχει επιλεγεί.
<code>get_groups_human_readable</code>	Παρέχει ομάδες σε αναγνώσιμη μορφή από τον άνθρωπο.
<code>get_example_attributes</code>	Εξάγει χαρακτηριστικά από το παράδειγμα από το παράδειγμα που έχει επιλεγεί.
<code>get_data</code>	Παρέχει τα δεδομένα για από το παράδειγμα που έχει επιλεγεί.
<code>get_structured_info</code>	Παρέχει δομημένη πληροφορία για τα παραδείγματα.
<code>train_and_evaluate</code>	Εκπαιδεύει και αξιολογεί το μοντέλο παραδείγματος.
<code>get_mitigated_results</code>	Εξάγει τα αποτελέσματα μετά την εφαρμογή τεχνικών μείωσης μεροληψίας.
<code>mitigate</code>	Μειώνει τη μεροληψία από το παράδειγμα που έχει επιλεγεί.

Πίνακας 18: HTML αρχεία και οι λειτουργίες τους

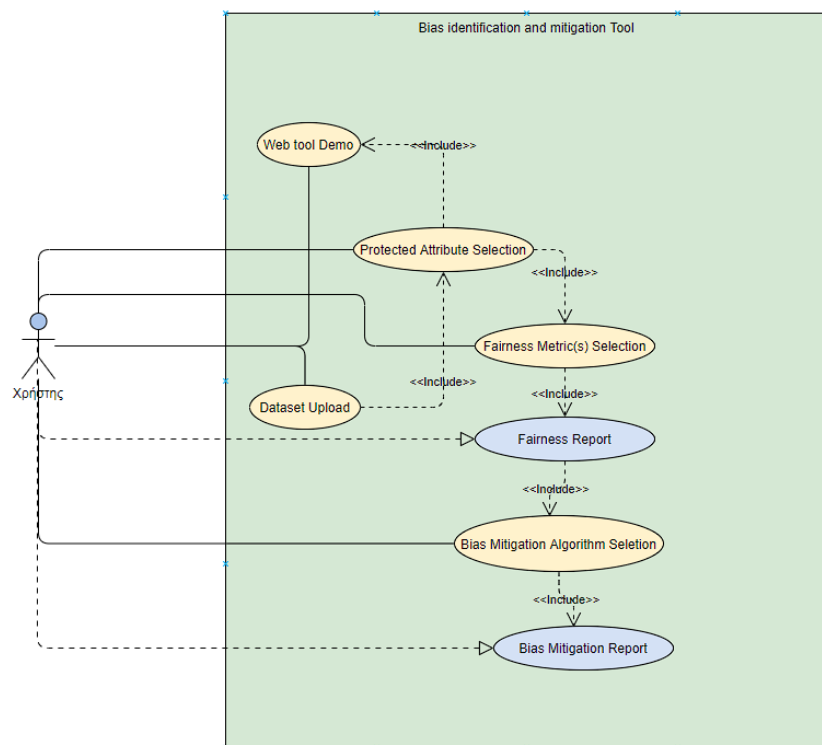
Όνομα Αρχείου	Λειτουργίες
base.html	Το κύριο πρότυπο που επεκτείνεται από τα υπόλοιπα αρχεία HTML.
index.html	Η αρχική σελίδα της εφαρμογής.
chooseAtt.html	Σελίδα επιλογής προστατευόμενων χαρακτηριστικών για τον έλεγχο μεροληψίας.
metric.html	Σελίδα επιλογής μετρικών δικαιοσύνης.
fairness_report.html	Σελίδα εμφάνισης των αποτελεσμάτων των μετρικών δικαιοσύνης.
algorithms.html	Σελίδα επιλογής αλγορίθμων μείωσης μεροληψίας.
example_att.html	Σελίδα επιλογής χαρακτηριστικών των παραδειγμάτων.
mitigation_report.html	Σελίδα παρουσίασης των αποτελεσμάτων μείωσης μεροληψίας.
selection.html	Σελίδα επιλογής και ανέβασμα δατασετ και μοντέλων μηχανικής μάθησης.

4.2 Διάγραμμα Περιπτώσεων Χρήσης

Για να καθορίσουμε την απαιτούμενη λειτουργικότητα του συστήματος θα χρησιμοποιήσουμε διαγράμματα περίπτωσης χρήσης[71], αναδεικνύοντας την αλληλεπίδραση του συστήματος με εξωτερικές οντότητες δηλαδή τους χρήστες και των λειτουργιών του συστήματος, τους παράγοντες, χωρίς όμως να γίνεται ανάλυση της δομής του, θα αναφερθούμε σε αυτό ως μαύρο κουτί. Στη συνέχεια θα παρουσιάσουμε τα διαγράμματα χρήσης:

1. «Χρήση Dataset Χρήστη»
2. «Χρήση Demo Web Tool Dataset »

Θα ακολουθήσει οπτική και λεκτική ανάλυση των περιπτώσεων διαγραμμάτων χρήσης.



Σχήμα 9: Διάγραμμα Περιπτώσεων Χρήσης

• Περίπτωση Χρήσης 1: «Χρήση Dataset Χρήστη»

Αναγνωριστικό: ΠΧ1

Περιγραφή: Ο χρήστης επιλέγει να χρησιμοποιήσει dataset της επιλογής του. Αφού πρώτα έχει ακολουθήσει τις οδηγίες για την τροποποίηση του dataset.

Βασική Ροή:

1. Ο χρήστης επιλέγει να ανεβάσει dataset κατάλληλης μορφής και επιλέγει ένα αλγόριθμο με τον οποίο θα εκπαιδευτεί το μοντέλο του.
2. Το σύστημα μεταφέρει τον χρήστη στη σελίδα επιλογής προστατευόμενων χαρακτηριστικών.
3. Ο χρήστης επιλέγει τα προτεινόμενα προστατευόμενα χαρακτηριστικά και τις τιμές τους.
4. Το σύστημα μεταφέρει τον χρήστη στη σελίδα επιλογής μετρικών δικαιοσύνης.
5. Ο χρήστης επιλέγει όποιες από τις προτεινόμενες μετρικές δικαιοσύνης επιθυμεί.
6. Το σύστημα μεταφέρει τον χρήστη στη σελίδα αποτελεσμάτων των μετρικών δικαιοσύνης.
7. Ο χρήστης επιλέγει συνέχεια.
8. Το σύστημα μεταφέρει τον χρήστη στη σελίδα επιλογής αλγορίθμων μείωσης μεροληψίας.
9. Ο χρήστης επιλέγει οποιουσδήποτε αλγορίθμους μείωσης μεροληψίας επιθυμεί.
10. Το σύστημα μεταφέρει τον χρήστη στη σελίδα αποτελεσμάτων των μετρικών δικαιοσύνης όπου εφαρμόστηκαν οι αλγόριθμοι μείωσης μεροληψίας.

Εναλλακτική Ροή:

- 1.α.1 Το σύστημα διαπιστώνει ότι το αρχείο του χρήστη είναι σε λάθος μορφή και εμφανίζει μήνυμα σφάλματος.
- 6.α.1 Το σύστημα διαπιστώνει ότι το αρχείο του χρήστη είναι δίκαιο και εμφανίζει κατάλληλο μήνυμα.

• **Περίπτωση Χρήσης 2: «Χρήση Demo Web Tool Dataset »**

Αναγνωριστικό: ΠΧ 2

Περιγραφή: Ο χρήστης επιλέγει ένα από τα demo datasets που προσφέρονται για την εκμάθηση χρήσης του εργαλείου στη σελίδα επιλογής.

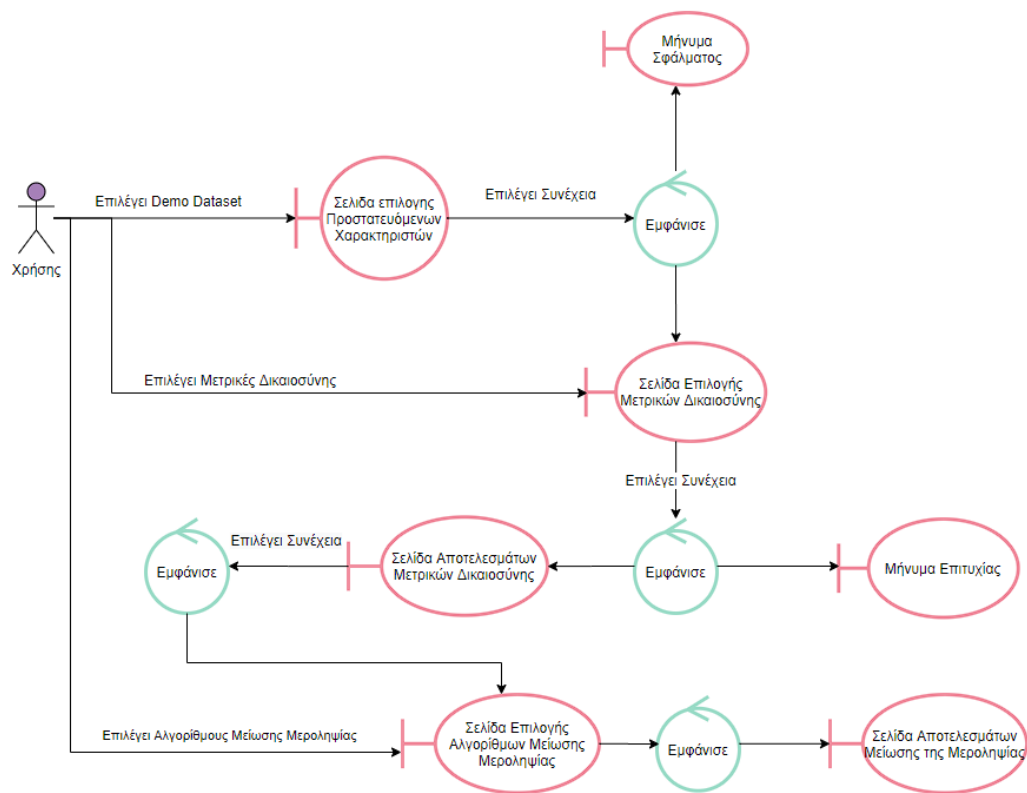
Βασική Ροή:

- 1. Ο χρήστης επιλέγει ένα από τα demo datasets που προσφέρονται.
- 2. Το σύστημα μεταφέρει τον χρήστη στη σελίδα επιλογής προστατευόμενων χαρακτηριστικών.
- 3. Ο χρήστης επιλέγει τα προτεινόμενα προστατευόμενα χαρακτηριστικά και τις τιμές τους.
- 4. Το σύστημα μεταφέρει τον χρήστη στη σελίδα επιλογής μετρικών δικαιοσύνης.
- 5. Ο χρήστης επιλέγει όποιες από τις προτεινόμενες μετρικές δικαιοσύνης επιθυμεί.
- 6. Το σύστημα μεταφέρει τον χρήστη στη σελίδα αποτελεσμάτων των μετρικών δικαιοσύνης.
- 7. Ο χρήστης επιλέγει συνέχεια.
- 8. Το σύστημα μεταφέρει τον χρήστη στη σελίδα επιλογής αλγορίθμων μείωσης μεροληψίας.
- 9. Ο χρήστης επιλέγει οποιουσδήποτε αλγορίθμους μείωσης μεροληψίας επιθυμεί.

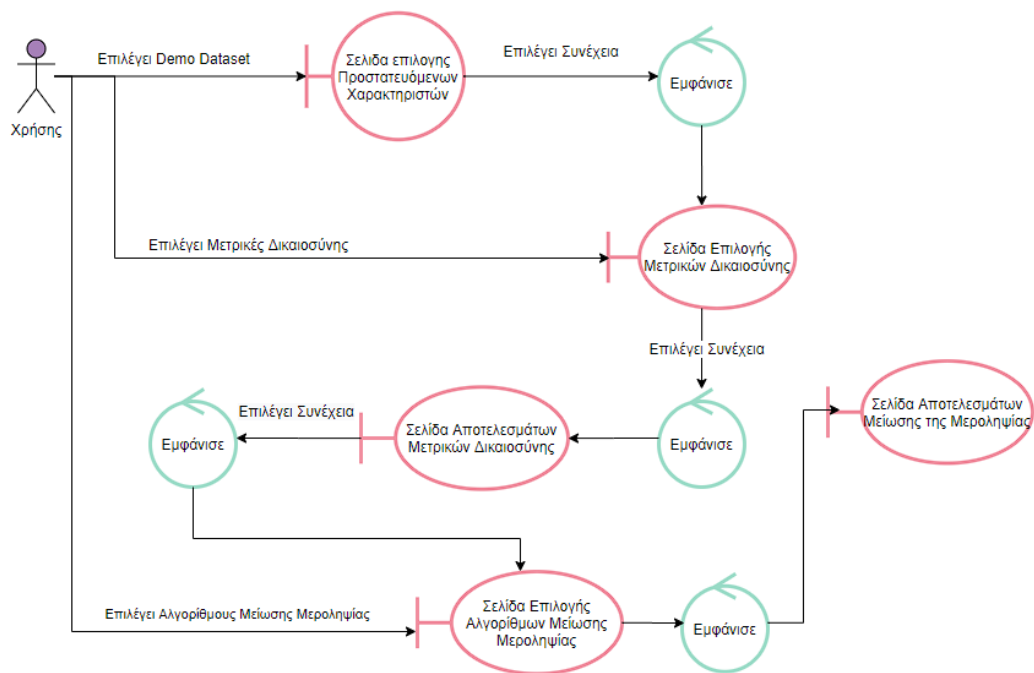
10. Το σύστημα μεταφέρει τον χρήστη στη σελίδα αποτελεσμάτων των μετρικών δικαιοσύνης όπου εφαρμόστηκαν οι αλγόριθμοι μείωσης μεροληψίας.

4.3 Διάγραμμα Ευρωστίας

Τα διαγράμματα Ευρωστίας είναι η γραφική απεικόνιση των περιπτώσεων χρήσης του συστήματος ξεχωριστά. Στη συνέχεια παρατίθενται τα εν λόγω διαγράμματα.



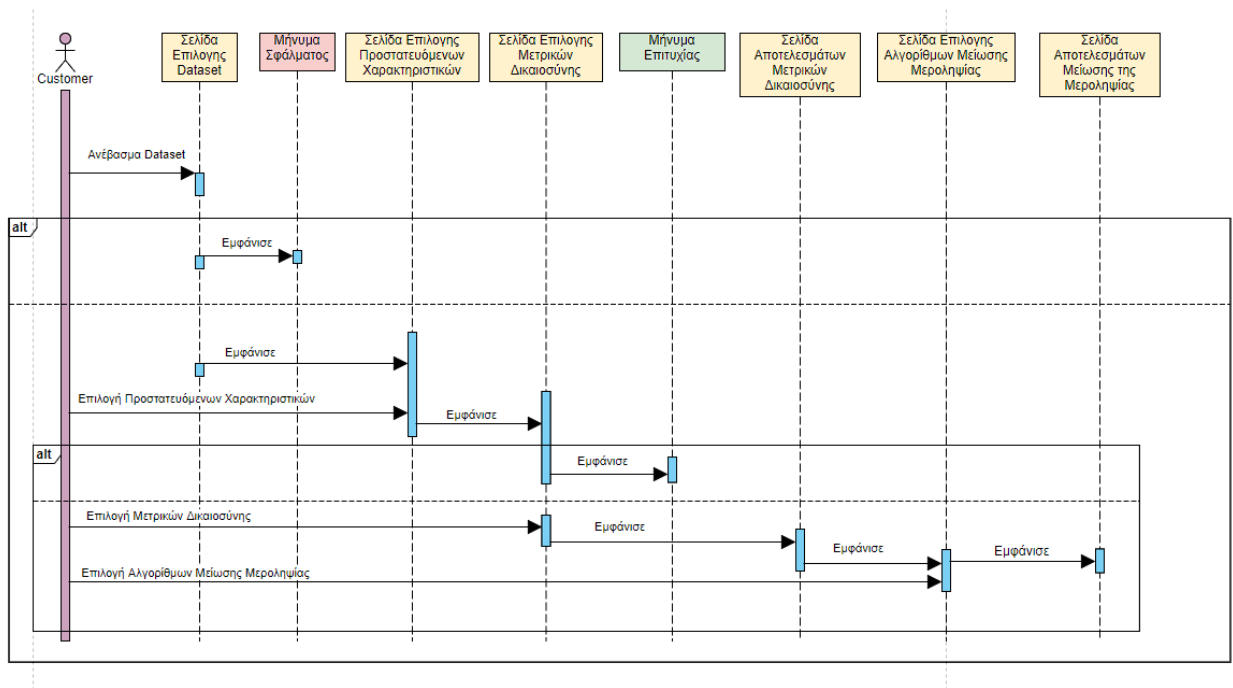
Σχήμα 10: Διάγραμμα Ευρωστίας Περίπτωσης Χρήσης 1



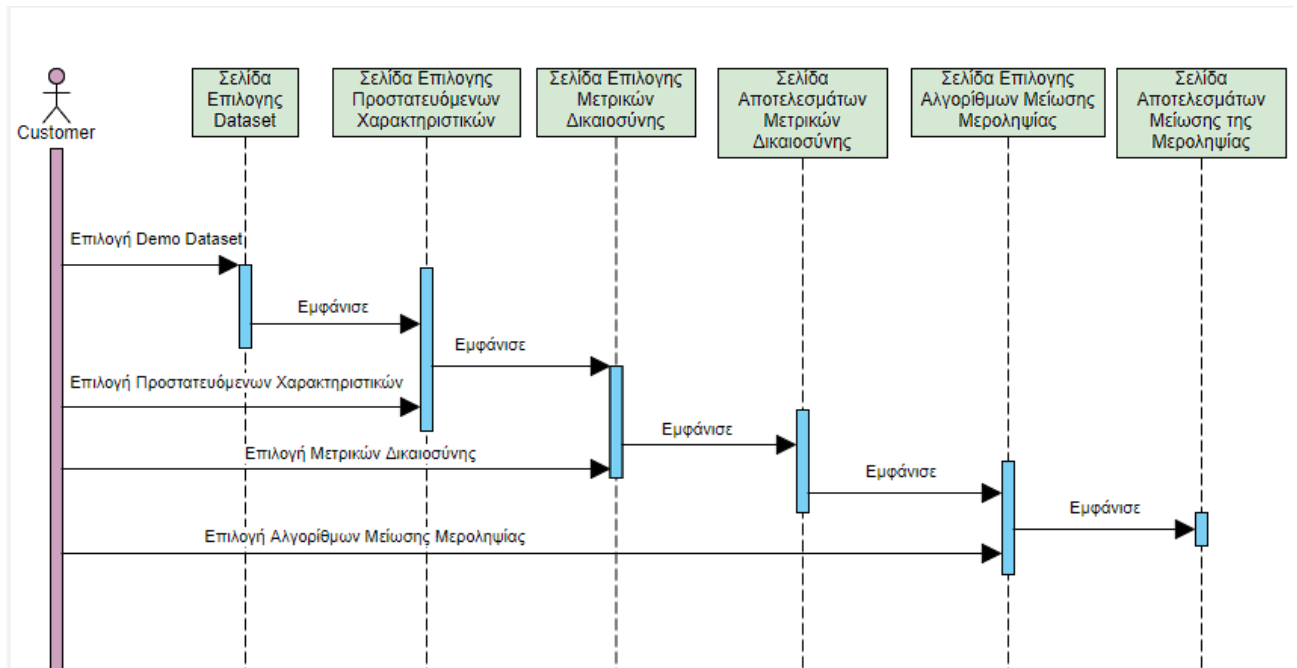
Σχήμα 11: Διάγραμμα Ευρωστίας Περίπτωσης Χρήσης 2

4.4 Διάγραμμα Ακολουθίας

Το διάγραμμα ακολουθίας απεικονίζει την αλληλεπίδραση μεταξύ αντικειμένων με διαδοχική σειρά, δηλαδή τη σειρά με την οποία πραγματοποιούνται αυτές οι αλληλεπιδράσεις. Ένα διάγραμμα ακολουθίας είναι δομημένο με τέτοιο τρόπο ώστε να αντιπροσωπεύει ένα χρονοδιάγραμμα που ξεκινά από την κορυφή και κατεβαίνει σταδιακά για να σημειώσει την ακολουθία των αλληλεπιδράσεων. Το διάγραμμα ακολουθίας ονομάζεται μερικές φορές διάγραμμα συμβάντων ή σενάριο συμβάντων. Παρακάτω ακολουθούν τα διαγράμματα ακολουθίας για κάθε περίπτωση χρήσης.



Σχήμα 12: Διάγραμμα Ακολουθίας Περίπτωσης Χρήσης 1



Σχήμα 13: Διάγραμμα Ακολουθίας Περίπτωσης Χρήσης 2

4.5 Εργαλεία και Βιβλιοθήκες

Στη συνέχεια παρουσιάζονται τα διάφορα εργαλεία και βιβλιοθήκες που χρησιμοποιήθηκαν στην ανάπτυξη του συστήματός μας. Αυτές οι τεχνολογίες κατηγοριοποιήσαμε σε backend και frontend τεχνολογίες, περιγράφοντας τις λειτουργίες τους και τη σημασία τους.

4.5.1 Backend Τεχνολογίες

- **Flask** - Flask είναι ένα μικρό και επεκτάσιμο πλαίσιο για τη δημιουργία web εφαρμογών σε Python . Χρησιμοποιήθηκε για τη διαχείριση των διαδρομών και των βασικών λειτουργιών της εφαρμογής.
- **TensorFlow** - TensorFlow είναι μια ανοιχτού κώδικα βιβλιοθήκη για μηχανική μάθηση και τεχνητή νοημοσύνη. Χρησιμοποιήθηκε για την ανάπτυξη και την

εκπαίδευση των μοντέλων μηχανικής μάθησης [72].

- **AIF360** - AI Fairness 360 είναι μια βιβλιοθήκη της IBM που παρέχει αλγόριθμους και εργαλεία για τη μέτρηση και μείωση της μεροληψίας σε μοντέλα μηχανικής μάθησης .
- **NumPy** - NumPy είναι μια βιβλιοθήκη για την υποστήριξη μεγάλων, πολυδιάστατων πινάκων και συναρτήσεων μαθηματικών υπολογισμών σε Python [73].
- **Scikit-learn** - Scikit-learn είναι μια βιβλιοθήκη μηχανικής μάθησης για τη γλώσσα προγραμματισμού Python , που υποστηρίζει εποπτευόμενη και μη εποπτευόμενη μάθηση [74].
- **Werkzeug** - Werkzeug είναι μια βιβλιοθήκη WSGI για Python που χρησιμοποιείται ως βοηθητικό εργαλείο για την ανάπτυξη web εφαρμογών [75].

4.5.2 Frontend Τεχνολογίες

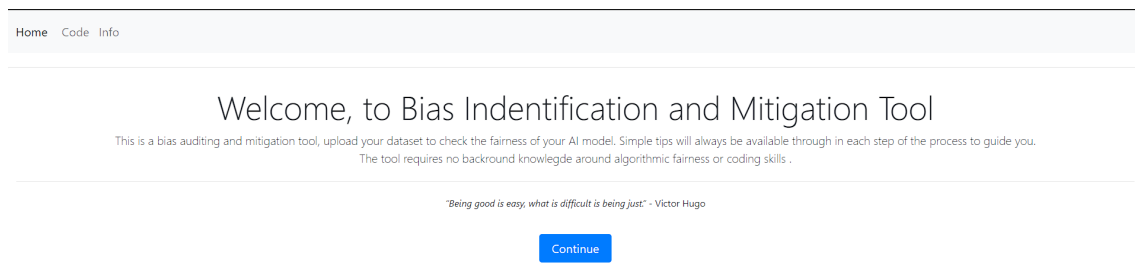
- **Bootstrap 4.5.2** - Bootstrap είναι ένα δημοφιλές πλαίσιο για τη δημιουργία αντιδραστικών και κινητών πρώτων web ιστοσελίδων. Χρησιμοποιήθηκε για τη σχεδίαση του frontend της εφαρμογής.
- **jQuery** - jQuery είναι μια γρήγορη, μικρή και πλούσια σε χαρακτηριστικά βιβλιοθήκη JavaScript . Χρησιμοποιήθηκε για τη διευκόλυνση της γραφής σεναρίων HTML .
- **Chart.js** - Chart.js είναι μια απλή αλλά ευέλικτη βιβλιοθήκη JavaScript για την κατασκευή γραφημάτων. Χρησιμοποιήθηκε για την γραφική απεικόνιση των αποτελεσμάτων .
- **D3.js** - D3.js είναι μια βιβλιοθήκη JavaScript για την παραγωγή δυναμικών, διαδραστικών απεικονίσεων δεδομένων στο web . Χρησιμοποιήθηκε για την κατασκευή βοηθήματος του χρήστη για την επιλογή κατάλληλων μετρικών.

5 Παρουσίαση Εργαλείου και Εκπαιδευτικής Διαδικασίας

Στο κεφάλαιο αυτό θα αναλύσουμε την περίπτωση κατά την οποία ένας χρήστης επιθυμεί να εξοικειωθεί με το εργαλείο επιλέγοντας την εκπαιδευτική λειτουργία του εργαλείου, μέσω των διαθέσιμων demo dataset της εφαρμογής.

5.1 Αρχική Σελίδα

Πρόκειται για τη σελίδα υποδοχής του εργαλείου, στην οποία καλοσορίζεται ο χρήστης και του γίνεται μια συντομή περιγραφή της λειτουργικότητας της εφαρμογής.



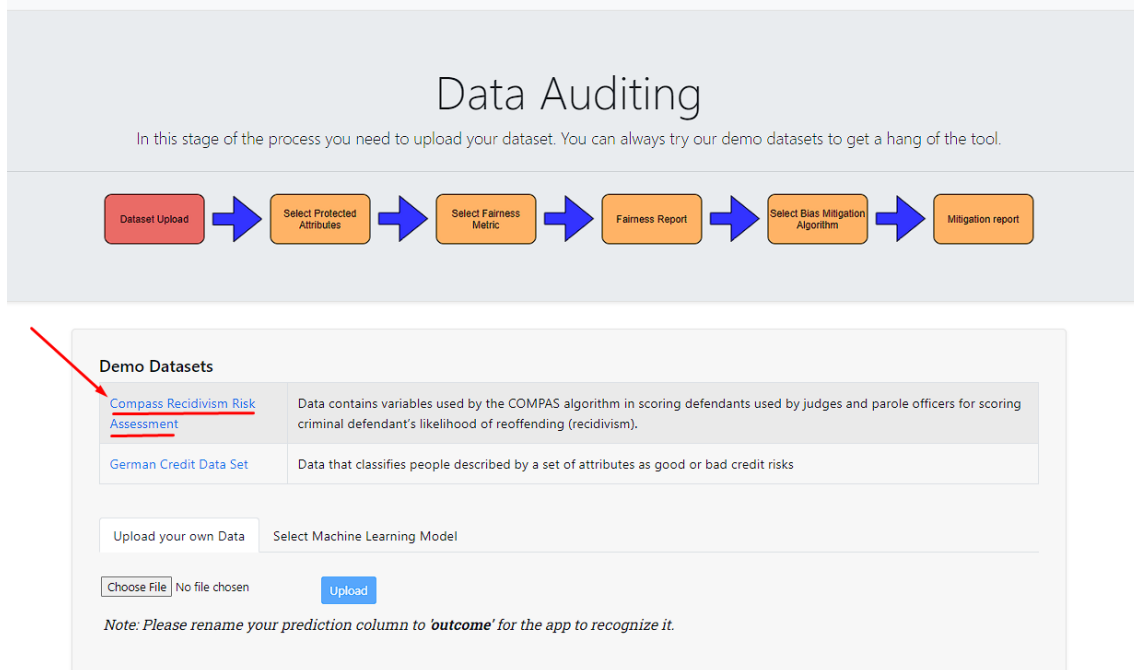
Σχήμα 14: Σελίδα υποδοχής της εφαρμογής

5.2 Επιλογή Συνόλου Δεδομένων

Υποθετικό Σενάριο: Ο χρήστης επιλέγει το σύνολο δεδομένων Compass Recidivism Risk Assessment[76], το οποίο περιέχει μεταβλητές που χρησιμοποιούνται από τον

5.3 Επιλογή Προστατευόμενων Χαρακτηριστικών

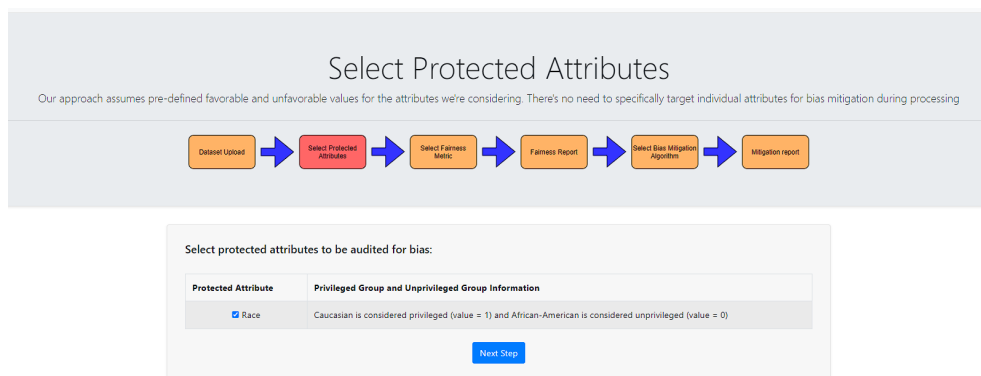
αλγόριθμο COMPAS για την αξιολόγηση κατηγορουμένων. Η βαθμολογία αυτή αξιοποιείται από δικαστές και αξιωματούχους για την εκτίμηση της πιθανότητας υποτροπής των κατηγορουμένων.



Σχήμα 15: Σελίδα επιλογής δεδομένων

5.3 Επιλογή Προστατευόμενων Χαρακτηριστικών

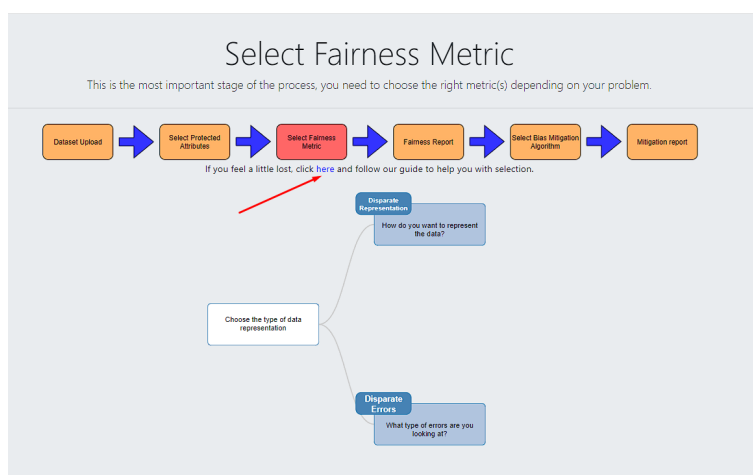
Ο χρήστης επιλέγει τα προστατευόμενα χαρακτηριστικά που επιθυμεί, στο παράδειγμα υπάρχει μία και μόνη επιλογή για διευκόλυνση του χρήστη και επεξήγηση του των προνομιούχων και μη προνομιούχων ομάδων βάσει του επιλεγμένου χαρακτηριστικού.



Σχήμα 16: Demo Σελίδα επιλογής δεδομένων

5.4 Επιλογή Μετρικών Δικαιοσύνης

Ο χρήστης καθοδηγείται στην επιλογή μετρικών μέσω ενός μενού που εμφανίζεται με την επιλογή "εδώ". Το μενού, σε μορφή δέντρου, οδηγεί τον χρήστη στην κατάλληλη επιλογή και αποκρύπτει τις μη σχετικές μετρήσεις.



Σχήμα 17: User metric guide

Υποθετικό Σενάριο: Ο χρήστης επιλέγει τις παρακάτω μετρήσεις και θέτει 80 τοις εκατό το threshold δικαιοσύνης:

5.5 Παρουσίαση Αποτελεσμάτων Μετρικών Δικαιοσύνης

№	Μετρική
1	Μέση Διαφορά (Mean Difference)
2	Disparate Impact
3	Average Abs Odds Difference

Πίνακας 19: Επιλεγμένες Μετρικές Δικαιοσύνης Εκπαιδευτικής Διαδικασίας

Select Fairness Metric

This is the most important stage of the process, you need to choose the right metric(s) depending on your problem.

Dataset Upload → Select Protected Attributes → Select Fairness Metric → Fairness Report → Select Bias Mitigation Algorithm → Mitigation report

If you feel a little lost, click [here](#) and follow our guide to help you with selection.

Fairness Audit Metric	Description
<input checked="" type="checkbox"/> Mean Difference	Mean Difference shows if one group is getting better results than another by comparing their average outcomes.
<input checked="" type="checkbox"/> Average Abs Odds Difference	Average Absolute Odds Difference is the average difference in false positive and true positive rates between groups.
<input checked="" type="checkbox"/> Disparate Impact	Disparate Impact checks if one group gets favorable results much more often than another.

Enter your Fairness Threshold (%) If a specific bias metric for a group is within this percentage of the reference group, this audit will pass

[Next Step](#)

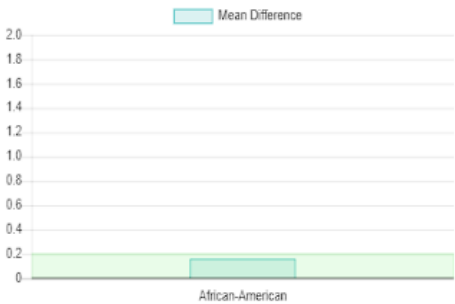
Σχήμα 18: Demo Επιλογή Μετρικών Δικαιοσύνης

5.5 Παρουσίαση Αποτελεσμάτων Μετρικών Δικαιοσύνης

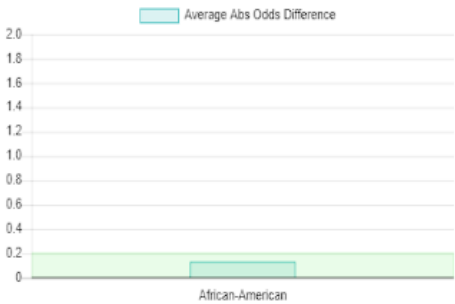
Εμφανίζονται οι τιμές των μετρικών και οι χαρακτηριστικές τιμές απόδοσης του συστήματος σε μορφή πινάκων και γραφημάτων.

Model Metrics			
Accuracy	F1 Score	Precision	Recall
0.6588	0.5885	0.6615	0.5299

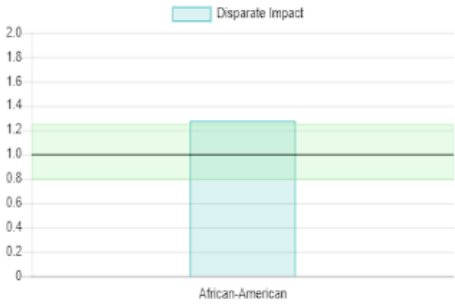
Mean Difference	
The ideal value of this metric is 0.0	
Protected Attributes	Race
Privileged Group	Caucasian
Unprivileged Group	African-American
Metric Value	0.1600



Average Abs Odds Difference	
The ideal value of this metric is 0.0	
Protected Attributes	Race
Privileged Group	Caucasian
Unprivileged Group	African-American
Metric Value	0.1306



Disparate Impact	
The ideal value of this metric is 1.0	
Protected Attributes	Race
Privileged Group	Caucasian
Unprivileged Group	African-American
Metric Value	1.2769



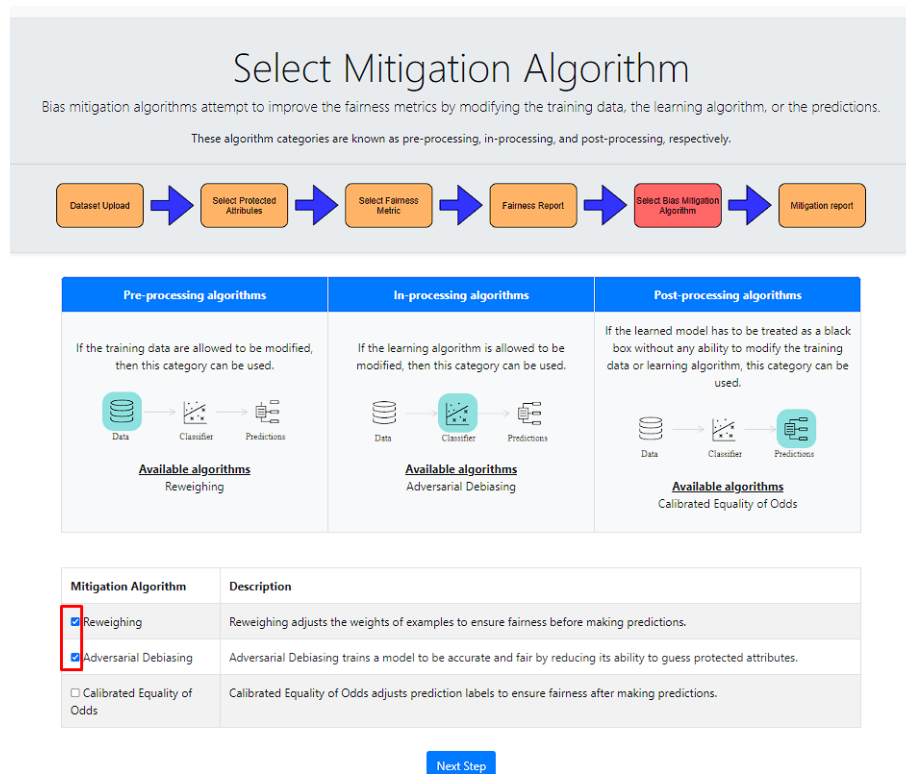
© 2024 Research and Development of a system for Bias Identification in Machine Learning Methods, [University of Patras](#)

Σχήμα 19: Demo Επιλογή Μετρικών Δικαιοσύνης

5.6 Επιλογή Αλγορίθμων Μείωσης Μεροληψίας

Ο χρήστης λαμβάνει ανάλυση για την επιλογή του κατάλληλου αλγορίθμου με βάσει των περιορισμών, που του έχουν τεθεί

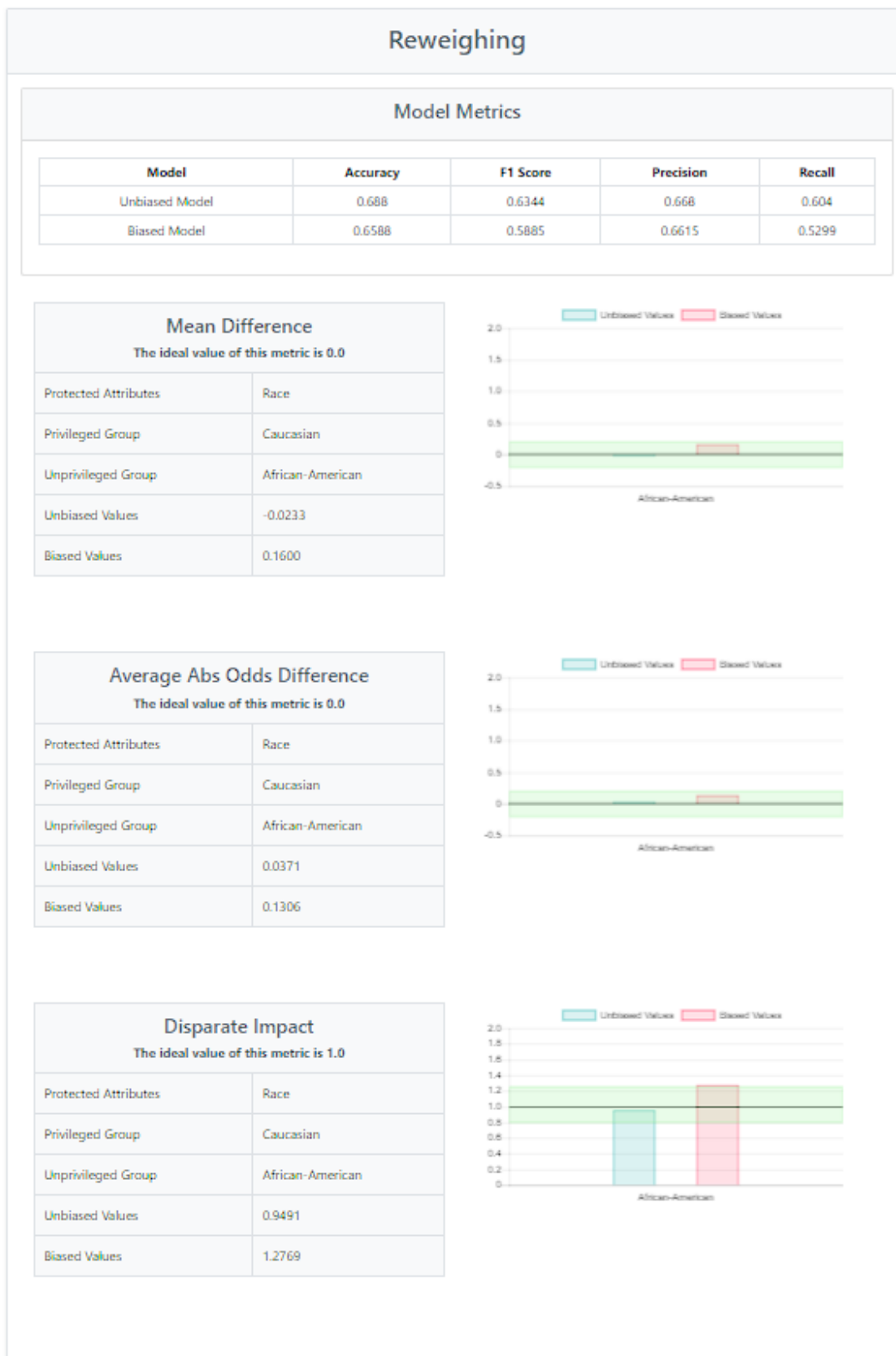
Υποθετικό Σενάριο: Ο χρήστης επιλέγει αναπροσαρμογή Βαρών (Reweighing) & Adversarial Debiasing από τους διαθέσιμους αλγορίθμους μείωσης μεροληψίας:



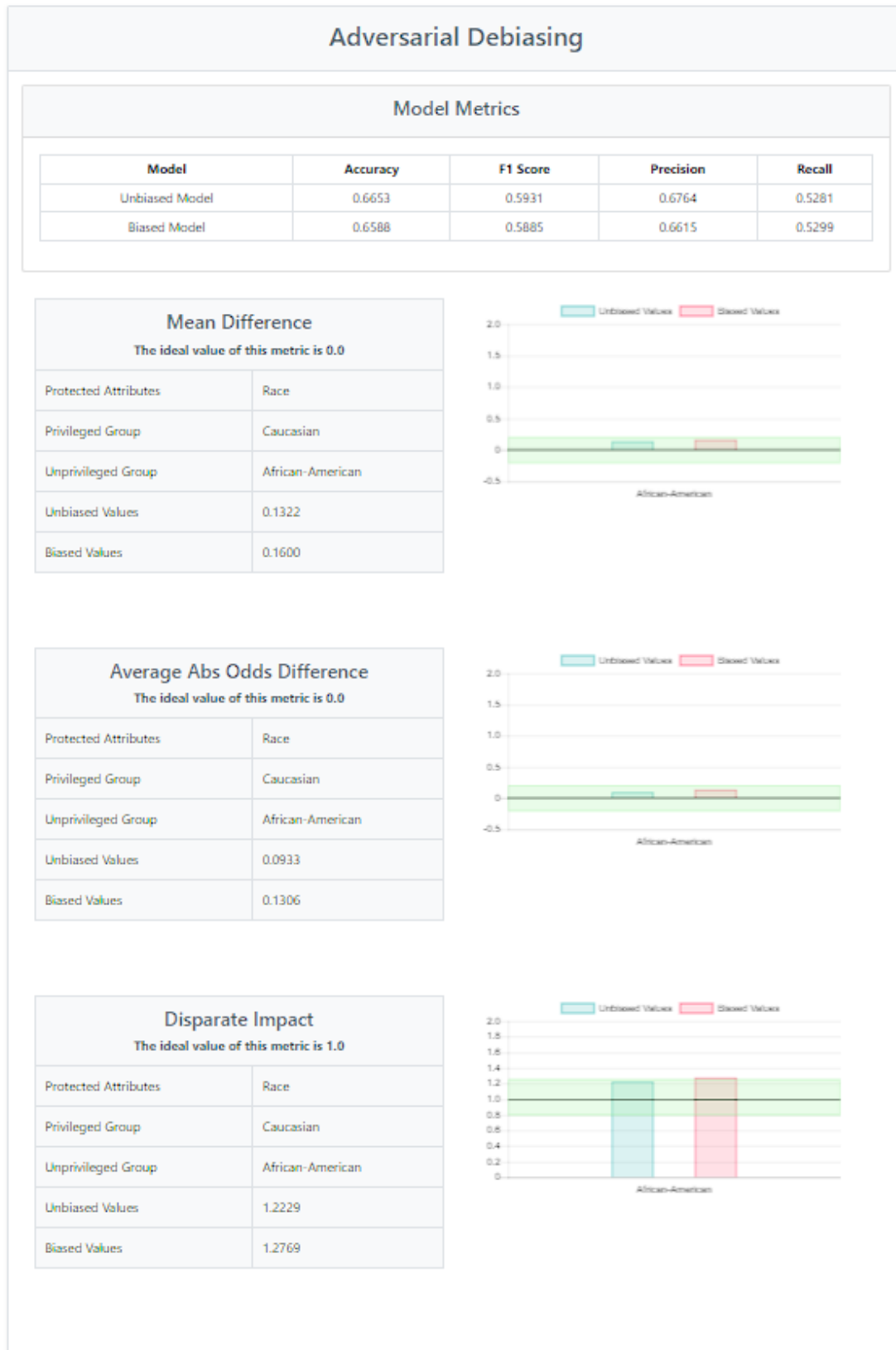
Σχήμα 20: Demo Επιλογή Αλγορίθμων Μείωσης Μεροληψίας

5.7 Τελικά Αποτελέσματα

Εμφανίζονται οι τιμές των επιλεγμένων μετρικών μετά την εφαρμογή των αλγορίθμων.



5.7 Τελικά Αποτελέσματα



6 Πειραματική Διαδικασία

Αυτή η ενότητα περιγράφει λεπτομερώς τη διαδικασία υλοποίησης ενός πειράματος, όπου ο χρήστης αξιοποιεί το εργαλείο για να λάβει τα επιθυμητά αποτελέσματα, χρησιμοποιώντας το δικό του σύνολο δεδομένων (dataset).

6.1 Επιλογή και Περιγραφή Dataset

Για τη μελέτη μας επιλέξαμε το Utrecht Fairness Recruitment Dataset [77], το οποίο περιέχει τις αποφάσεις πρόσληψης τεσσάρων εταιρειών για 500 υποψηφίους. Για κάθε υποψήφιο, έχουμε διαθέσιμες κάποιες γενικές περιγραφές (gender, age, sport) και μερικούς δείκτες. Το συγκεκριμένο dataset είναι συνθετικό και επιλογή αυτού του έγινε λόγω της δυνατότητας που προσφέρει για εξάσκηση στην ανάλυση ζητημάτων δικαιοσύνης.

6.2 Επιλογή Μοντέλου Μηχανικής Μάθησης

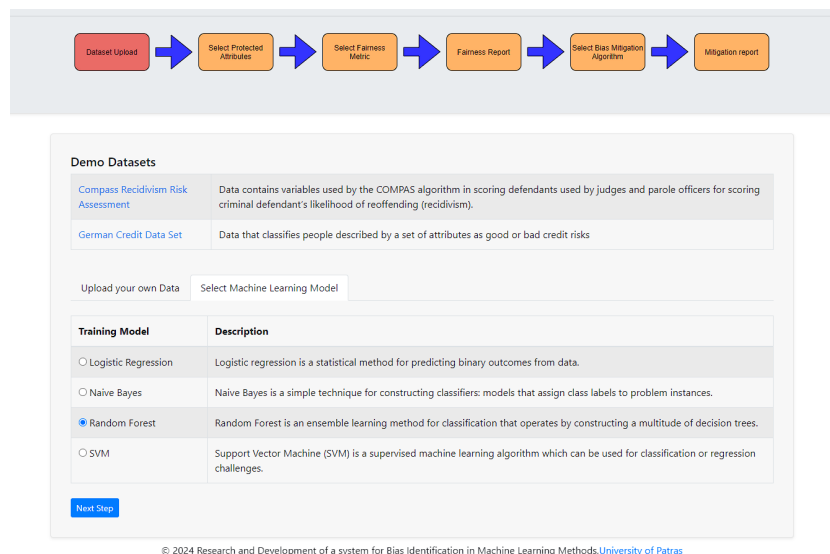
Η επιλογή του κατάλληλου μοντέλου μηχανικής μάθησης είναι κρίσιμη για την επιτυχία του πειράματος. Με δεδομένα τα χαρακτηριστικά του dataset και τις απαιτήσεις της ταξινόμησης, επιλέξαμε να χρησιμοποιήσουμε το μοντέλο Random Forest.

- **Υψηλή ακρίβεια και ανθεκτικότητα:** Παράγει ακριβείς προβλέψεις σε ποικιλία συνόλων δεδομένων, συνδυάζοντας πολλαπλά δέντρα αποφάσεων εκπαιδευμένα σε τυχαία υποσύνολα δεδομένων. Μειώνει την επίδραση των ακραίων τιμών και του θορύβου, καθιστώντας την λιγότερο επιρρεπή σε υπερεκπαίδευση και πιο ανθεκτική στις μεταβολές των δεδομένων.
- **Διαχείριση δεδομένων πολλών διαστάσεων :** Μπορεί να διαχειριστεί σύνολα δεδομένων με μεγάλο αριθμό χαρακτηριστικών χωρίς να απαιτείται επιλογή χαρακτηριστικών ή μείωση διάστασης. Αυτό είναι επωφελές όταν υπάρχουν

πολλοί πιθανοί εξηγητικοί μεταβλητοί, καθώς μπορεί να εντοπίσει αυτόματα τα πιο σημαντικά χαρακτηριστικά για την πρόβλεψη.

- **Αντιμετώπιση ελλειπών τιμών και ακραίων τιμών:** Μπορεί να διαχειριστεί ελλειπείς τιμές και ακραίες τιμές χωρίς την ανάγκη προεπεξεργασίας δεδομένων, χρησιμοποιώντας τεχνικές συμπλήρωσης και λαμβάνοντας υπόψη την κατανομή των δεδομένων κατά την κατασκευή των δέντρων αποφάσεων.
- **Ερμηνευσιμότητα και σημασία χαρακτηριστικών:** Παρέχει πληροφορίες για τη σημασία των διαφορετικών χαρακτηριστικών για την πρόβλεψη, βοηθώντας στην κατανόηση των υποκείμενων σχέσεων στα δεδομένα και στον εντοπισμό των πιο επιδραστικών παραγόντων που συμβάλλουν στο αποτέλεσμα.

Συνοψίζοντας, η επιλογή του Random Forest είναι βάσιμη λόγω της ικανότητάς του να διαχειρίζεται αποτελεσματικά την πολυπλοκότητα των δεδομένων μας, να παρέχει πολύτιμες πληροφορίες για τη σημαντικότητα των χαρακτηριστικών και να μειώνει την πιθανότητα υπερπροσαρμογής, εξασφαλίζοντας υψηλή απόδοση και αξιοπιστία.



Σχήμα 23: Demo Επιλογή Μοντέλου Κατά τη Χρήση του Εργαλείου

6.3 Ανάλυση Dataset

Το dataset περιέχει 4000 εγγραφές και 14 στήλες. Δεν υπάρχουν ελλείπουσες τιμές στο dataset. Η μεταβλητή-στόχος (target variable) είναι το outcome, το οποίο είναι δυαδικό (0 ή 1).

6.3.1 Κατηγορηματικά χαρακτηριστικά:

- **gender:** 3 μοναδικές τιμές (female, male, others).
- **nationality:** 3 μοναδικές τιμές (Dutch, German, Belgian).
- **sport:** 8 μοναδικές τιμές.
- **ind-debateclub, ind-programming_exp, ind-international_exp, ind-entrepreneur_exp, ind-exact_study:** Δυαδικά χαρακτηριστικά (True/False).
- **ind-degree:** 3 μοναδικές τιμές (bachelor, master, phd).
- **company:** 4 μοναδικές τιμές (A, B, C, D).

6.3.2 Αριθμητικά χαρακτηριστικά:

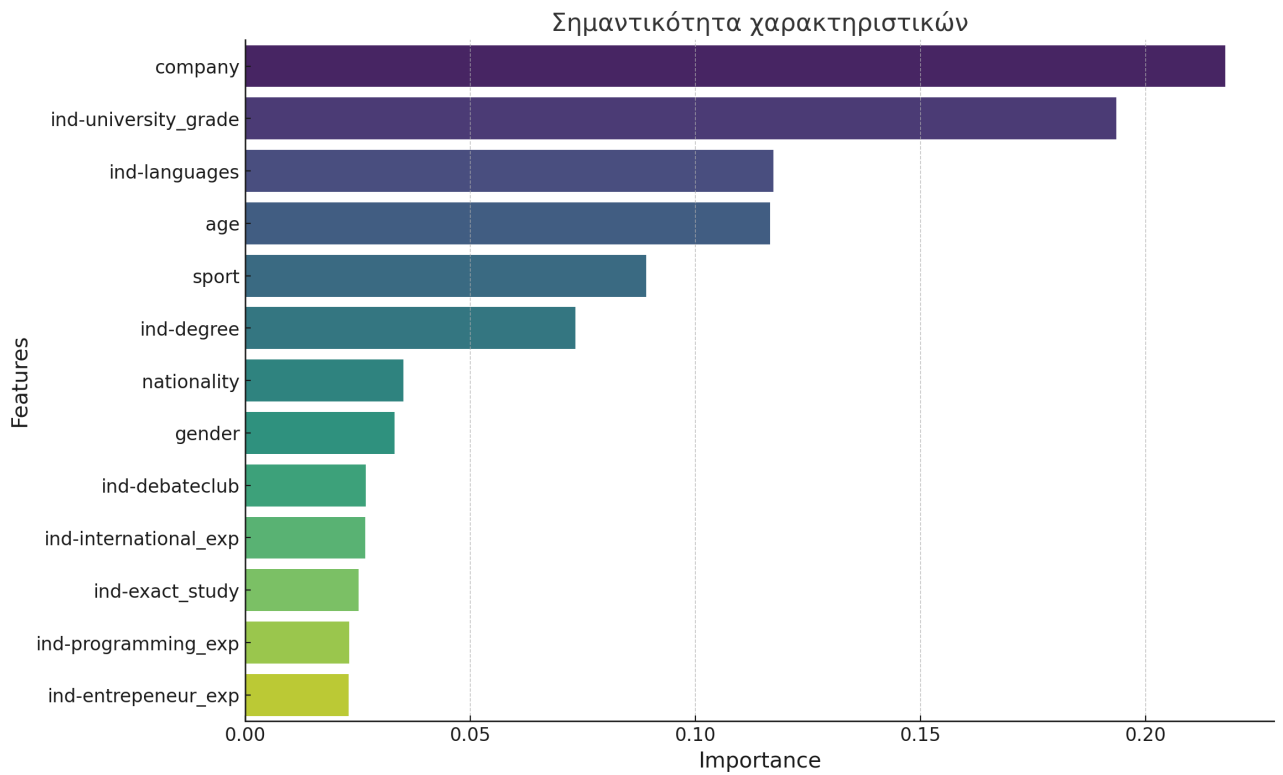
- **age:** Συνεχής μεταβλητή με μέσο όρο 26.18 χρόνια.
- **ind-university_grade:** Συνεχής μεταβλητή με μέσο όρο 62.38.
- **ind-languages:** Διακριτή μεταβλητή (0, 1, 2, 3).

6.3.3 Σημαντικότητα Χαρακτηριστικών

Η σημαντικότητα των χαρακτηριστικών είναι μια τεχνική που χρησιμοποιείται για να αξιολογηθεί η συνεισφορά κάθε χαρακτηριστικού στην πρόβλεψη. Η ανάλυση αυτή

επιτρέπει την κατανόηση των πιο κρίσιμων χαρακτηριστικών που επηρεάζουν το αποτέλεσμα του μοντέλου.

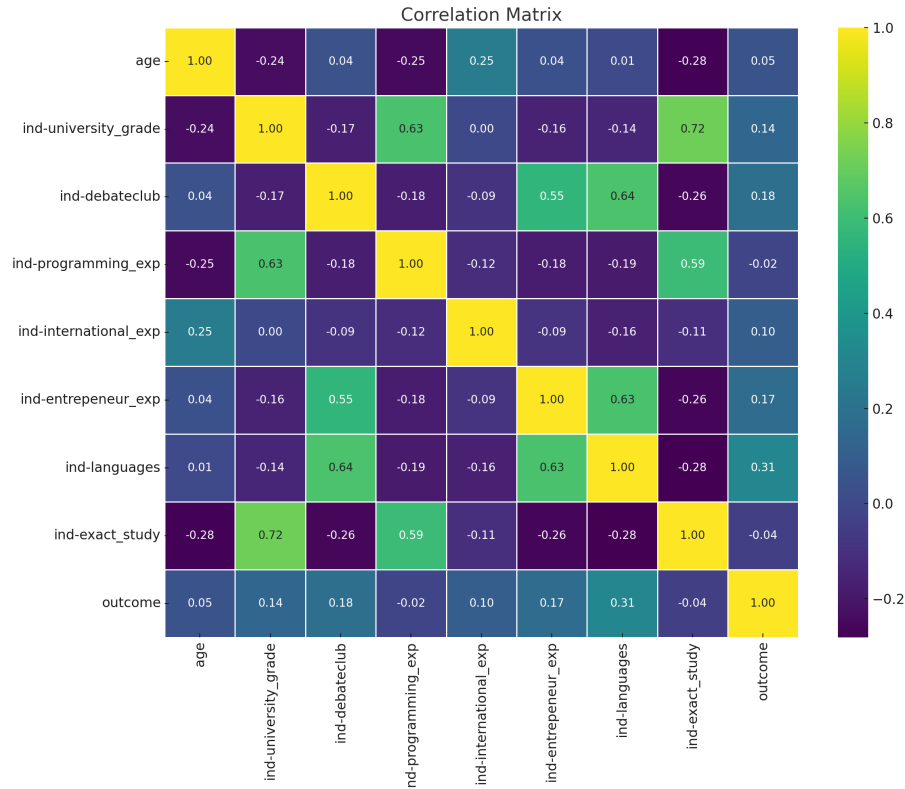
Το γράφημα που ακολουθεί παρουσιάζει τη σημαντικότητα των χαρακτηριστικών, ταξινομημένες κατά φθίνουσα σειρά:



Σχήμα 24: Σημαντικότητα Χαρακτηριστικών

6.3.4 Συσχέτιση Χαρακτηριστικών

Το γράφημα που ακολουθεί παρουσιάζει τη σημαντικότητα των χαρακτηριστικών, ταξινομημένες κατά φθίνουσα σειρά:



Σχήμα 25: Συσχέτιση Χαρακτηριστικών

Βασισμένα στη συσχέτιση των χαρακτηριστικών με το αποτέλεσμα (outcome), μπορούμε να εξάγουμε τα ακόλουθα συμπεράσματα: Το χαρακτηριστικό ind-languagesέχει την υψηλότερη θετική συσχέτιση με το αποτέλεσμα ($r = 0.31$). Αυτό σημαίνει ότι ο αριθμός των γλωσσών που γνωρίζει ένας υποψήφιος επηρεάζει θετικά την πιθανότητα πρόσληψής του. Αντίθετα, το χαρακτηριστικό ind-exact_studyέχει την χαμηλότερη συσχέτιση με το αποτέλεσμα ($r = -0.04$). Αυτό υποδηλώνει ότι η ακριβής μελέτη του υποψηφίου δεν έχει σημαντική επίδραση στην πιθανότητα πρόσληψής του.

Γενικά, οι τιμές των συσχετίσεων δείχνουν ότι τα περισσότερα χαρακτηριστικά έχουν χαμηλή συσχέτιση με το αποτέλεσμα. Αυτό σημαίνει ότι δεν υπάρχει κάποιο χαρακτηριστικό που από μόνο του να καθορίζει την πιθανότητα πρόσληψης των υποψηφίων. Αντιθέτως, η απόφαση για την πρόσληψη φαίνεται να εξαρτάται από έναν συνδυασμό πολλών διαφορετικών χαρακτηριστικών.

6.4 Επιλογή Προστατευμένων Χαρακτηριστικών και Ομάδων

Για την αξιολόγηση των διασταυρούμενων χαρακτηριστικών, εξετάστηκαν διάφοροι συνδυασμοί φύλου και εθνικότητας. Τα αποτελέσματα συνοψίζονται στον παρακάτω πίνακα:

Gender	Nationality	Actual Outcome Mean	Count
female	Belgian	0.260870	46
female	Dutch	0.252336	428
female	German	0.217391	46
male	Belgian	0.320000	75
male	Dutch	0.328740	508
male	German	0.405063	79

Table 20: Συνδυασμοί φύλου και εθνικότητας με μέσους όρους προβλεπόμενων και πραγματικών αποτελεσμάτων

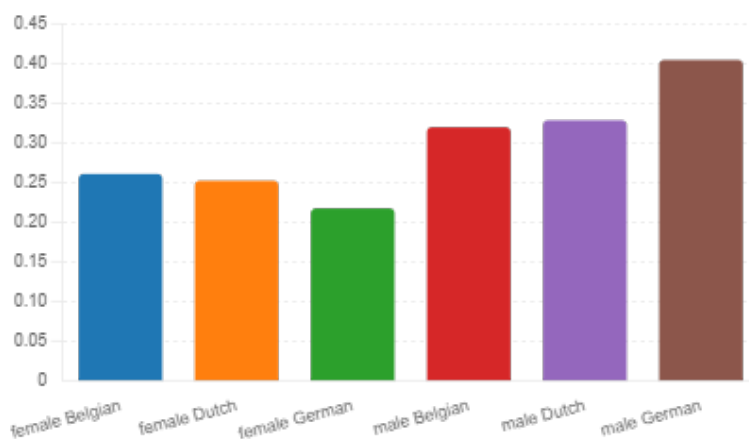


Figure 26: Συνδυασμοί φύλου και εθνικότητας με μέσους όρους

Από το σύνολο δεδομένων, θα επιλεχτούν οι συνδυασμοί φύλου και εθνικότητας,

male Dutch και female Belgian. Μεγαλύτερα δείγματα προσφέρουν πιο στενά διαστήματα εμπιστοσύνης, μειώνοντας την αβεβαιότητα των εκτιμήσεων. Καθώς αυξάνεται το μέγεθος δείγματος, τα διαστήματα εμπιστοσύνης γίνονται στενότερα, παρέχοντας πιο ακριβείς εκτιμήσεις για παραμέτρους όπως το μέσο βάρος ή το ύψος. Αυτό οφείλεται στο ότι μεγαλύτερα δείγματα μειώνουν την τυπική απόκλιση του μέσου όρου του δείγματος, κάνοντας την εκτίμηση πιο αξιόπιστη.[78]

Η επιλογή των χαρακτηριστικών φύλου και εθνικότητας, ατομικά και σε συνδυασμό, για την εξέταση της αλγοριθμικής δικαιοσύνης είναι κρίσιμη για διάφορους λόγους. Στατιστικά στοιχεία δείχνουν ότι γυναίκες και άνδρες βιώνουν διαφορετικές ευκαιρίες και μεταχείριση σε διάφορους τομείς, όπως η απασχόληση, η εκπαίδευση και η υγειονομική περίθαλψη. Η ενσωμάτωση του φύλου στα μοντέλα επιτρέπει τον εντοπισμό και την αντιμετώπιση προκαταλήψεων που μπορεί να οδηγήσουν σε ανισότητες στις προβλέψεις [79]. Η εθνικότητα σχετίζεται με κοινωνικοοικονομικές διαφορές, πρόσβαση σε πόρους και ευκαιρίες, και ιστορικές ανισότητες. Η ανάλυση βάσει εθνικότητας βοηθά στην αποκάλυψη προκαταλήψεων που βασίζονται σε στερεότυπα και διακρίσεις [80]. Η ανάλυση αλγοριθμικής δικαιοσύνης με βάση τα παραπάνω χαρακτηριστικά εντοπίζει και αντιμετωπίζει συστημικές ανισότητες. Η διαφάνεια και η λογοδοσία είναι απαραίτητες για την προστασία των δικαιωμάτων και την ισότιμη μεταχείριση όλων των πολιτών. Η συνειδητή λήψη αποφάσεων με βάση τα δεδομένα προωθεί την κοινωνική δικαιοσύνη και την ισότητα ευκαιριών. Η επιλογή φύλου και εθνικότητας, ατομικά και σε συνδυασμό, είναι απαραίτητη για την ανάπτυξη δίκαιων και αξιόπιστων αλγορίθμων.

Select Protected Attributes

Choose attributes to check for fairness, such as race or gender. Select the attribute and the reference group value. The tool will compare other groups against this reference group to identify biases.

Dataset Upload → Select Protected Attributes → Select Fairness Metric → Fairness Report → Select Bias Mitigation Algorithm → Mitigation report

Select protected attributes to be audited for bias:

<input type="checkbox"/>	age	Privileged	Unprivileged	Add Intersectional Group
<input type="checkbox"/>	company	Privileged	Unprivileged	Add Intersectional Group
<input checked="" type="checkbox"/>	gender	male	female	
	nationality	Dutch	Belgian	Remove

Σχήμα 27: Επιλογή Προστατευμένων Χαρακτηριστικών και Ομάδων

6.5 Επιλογή Μετρικών Δικαιοσύνης

Θα αξιοποιήσουμε τη μετρική disparate impact σε συνδυασμό με αναλύσεις συνδυασμών φύλου και εθνικότητας. Στόχος μας είναι να διασφαλίσουμε ότι το μοντέλο μας λαμβάνει υπόψη τις αλληλεπιδράσεις μεταξύ αυτών των χαρακτηριστικών και να αποφεύγουμε τυχόν δυσανάλογη επίδραση σε συγκεκριμένες υποομάδες. Η μετρική disparate impact αποτελεί μια ευρέως αναγνωρισμένη μετρική για τον εντοπισμό διακρίσεων σε αλγοριθμικές αποφάσεις. Εστιάζει στην σύγκριση των αποτελεσμάτων μεταξύ διαφορετικών ομάδων, επισημαίνοντας τυχόν δυσανάλογη επίδραση σε προστατευόμενες ομάδες[81].

Η χρήση της επιτρέπει την αξιολόγηση της αλγοριθμικής δικαιοσύνης σε ένα πιο πρακτικό και νομικά αποδεκτό πλαίσιο. Σε πλήθος νομοθεσιών, η έννοια του disparate impact έχει ήδη ενσωματωθεί, καθιστώντας τη μετρική αυτή απαραίτητη για την τήρηση των νομοθετικών προδιαγραφών [36].

Τέλος, ευθυγραμμίζεται με τους στόχους μας για την προώθηση της ισότητας και της κοινωνικής δικαιοσύνης, παρέχοντας ένα αξιόπιστο εργαλείο για την συνεχή παρακολούθηση και βελτίωση της δικαιοσύνης των αλγορίθμων μας.

Select Fairness Metric

This is the most important stage of the process, you need to choose the right metric(s) depending on your problem.

Dataset Upload → Select Protected Attributes → Select Fairness Metric → Fairness Report → Select Bias Mitigation Algorithm → Mitigation report

If you feel a little lost, click [here](#) and follow our guide to help you with selection.

Fairness Audition Metric	Description
<input type="checkbox"/> Mean Difference	Measures the ratio of favorable outcomes received by the unprivileged group to that of the privileged group. A value close to 1 indicates fairness.
<input checked="" type="checkbox"/> Disparate Impact	Measures the difference in the rate of favorable outcomes between the unprivileged and privileged groups. A value close to 0 indicates fairness.
<input type="checkbox"/> Equal Opportunity Difference	Measures the difference in true positive rates between the unprivileged and privileged groups. A value close to 0 indicates fairness.
<input type="checkbox"/> Average Odds Difference	Measures the average difference in false positive rates and true positive rates between the unprivileged and privileged groups. A value close to 0 indicates fairness.
<input type="checkbox"/> Theil Index	Measures the inequality in the distribution of predictions. Lower values indicate less inequality and more fairness.

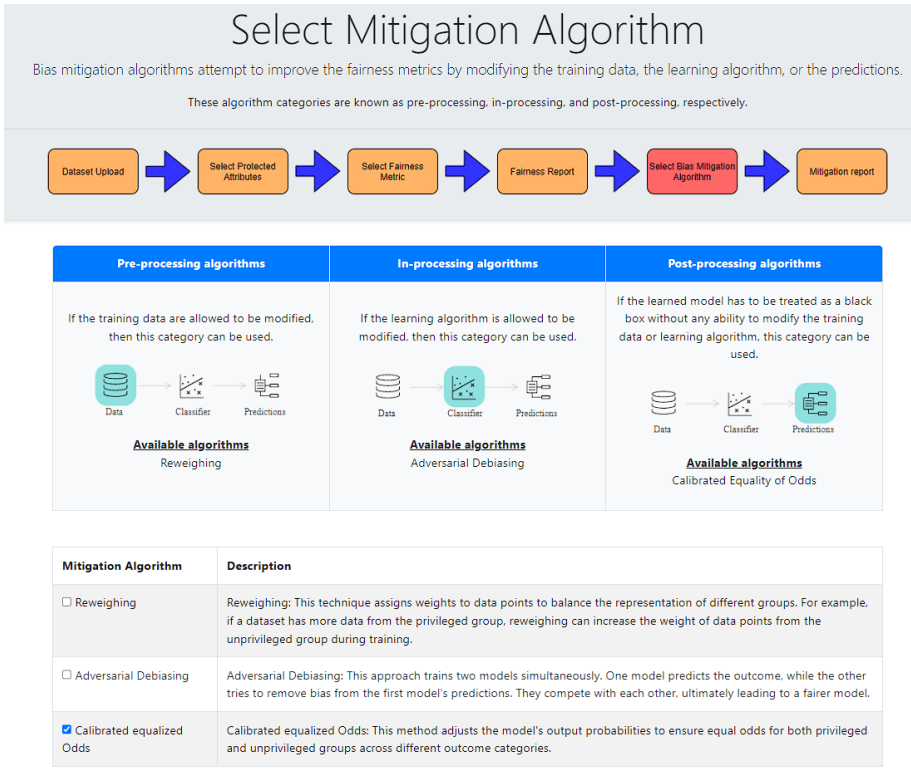
If a specific bias metric for a group is within this percentage of the reference group, this audit will pass

Next Step

Σχήμα 28: Επιλογή Μετρικών Δικαιοσύνης

6.6 Επιλογή Αλγορίθμου Μείωσης Μεροληψίας

Επιλέξαμε τον αλγόριθμο Calibrated Equalized Odds, ο οποίος αποτελεί μέθοδο που εφαρμόζεται μετά την επεξεργασία (post-processing) και χρησιμοποιεί τη βαθμονόμηση των εξόδων του ταξινομητή για να βρει πιθανότητες με τις οποίες να αλλάξει τις ετικέτες εξόδου, διασφαλίζοντας ίσες πιθανότητες σφάλματος μεταξύ των ομάδων. Είναι ιδανικός για σενάρια όπου δεν υπάρχει πρόσβαση ή δικαίωμα επέμβασης στα δεδομένα ή το μοντέλο, καθώς μπορεί να εφαρμοστεί απευθείας στα αποτελέσματα του ταξινομητή. Αυτό είναι ιδιαίτερα σημαντικό όταν το μοντέλο θεωρείται «black-box» και η τροποποίηση των εσωτερικών του παραμέτρων δεν είναι δυνατή.



Σχήμα 29: Επιλογή Αλγορίθμου Μείωσης Μεροληψίας

7 Συμπεράσματα και Μελλοντικές Επεκτάσεις

7.1 Συνοψη Αποτελεσμάτων

7.1.1 Συμμόρφωση με τον τοπικό Νόμο 144 του 2021

Συμμόρφωση με τον Τοπικό Νόμο 144 του 2021

Η εξέταση διασταυρούμενων χαρακτηριστικών είναι απαραίτητη για την τήρηση του Τοπικού Νόμου 144 του 2021. Ο νόμος αυτός προωθεί την δικαιοσύνη και την ισότητα στην αλγοριθμική λήψη αποφάσεων, απαγορεύοντας τη μεροληψία με βάση συνδυασμούς χαρακτηριστικών.

Η Σημασία των Διασταυρούμενων Χαρακτηριστικών

Προηγούμενες προσεγγίσεις, που εστιάζαν σε μεμονωμένα χαρακτηριστικά, δεν ήταν επαρκείς για την εξασφάλιση πλήρους δικαιοσύνης. Ο Τοπικός Νόμος 144 του 2021 εισήγαγε την έννοια της διασταυρούμενης μεροληψίας, η οποία λαμβάνει υπόψη ταυτόχρονα πολλαπλά χαρακτηριστικά των ατόμων.

Αυτή η προσέγγιση αναγνωρίζει την πολυπλοκότητα της ανθρώπινης ταυτότητας. Οι άνθρωποι ανήκουν σε διάφορες ομάδες ταυτόχρονα (π.χ. φύλο, φυλή, κοινωνικοοικονομική τάξη) και η δίκαιη αντιμετώπιση οφείλει να λαμβάνει υπόψη αυτές τις αλληλεπιδράσεις.

7.2 Συμπεράσματα

Η παρούσα διπλωματική εργασία επικεντρώθηκε στην ανάπτυξη ενός εργαλείου που επιτρέπει την ανίχνευση και μείωση της μεροληψίας σε μοντέλα μηχανικής μάθησης. Μέσα από τα πειράματα και τα αποτελέσματα που παρουσιάστηκαν, μπορούμε να συμπεράνουμε ότι το εργαλείο είναι αποτελεσματικό στην...

7.3 Μελλοντικές Κατευθύνσεις

Για μελλοντική έρευνα, προτείνονται οι ακόλουθες κατευθύνσεις:

- **Βελτίωση Αλγορίθμων Μείωσης Μεροληψίας:** Ανάπτυξη και δοκιμή νέων αλγορίθμων για την περαιτέρω μείωση της μεροληψίας.
- **Μεγαλύτερα και Πιο Ποικίλα Δεδομένα Εκπαίδευσης:** Συλλογή και χρήση μεγαλύτερων και πιο ποικίλων συνόλων δεδομένων για την εκπαίδευση των μοντέλων.
- **Συνεχής Ενημέρωση και Προσαρμογή:** Παρακολούθηση των εξελίξεων στους τομείς της νομοθεσίας και της ηθικής στην τεχνητή νοημοσύνη, και προσαρμογή του εργαλείου στις νέες απαιτήσεις.

Βιβλιογραφία

- [1] IBM. Aif360: Ai fairness 360, 2023. <https://odsc.com/speakers/introducing-the-ai-fairness-360-toolkit-3/>.
- [2] NYC Department of Consumer and Worker Protection (DCWP). Local law 144 of 2021, 2021. <https://www1.nyc.gov/site/dca/about/local-laws.page>.
- [3] J. Budd and S. Barocas. Automating bias? examining the effect of gender on resume screening decisions. *Journal of AI Ethics*, 2018.
- [4] D. Lewis. Will ai hire better than humans? insights from amazon’s automated hiring tool. *AI Magazine*, 2018.
- [5] S. Goel and J. Angwin. Accuracy and fairness of recidivism prediction algorithms. *Journal of Criminal Justice*, 2021.
- [6] J. Lewandowski. Machine learning and criminal justice: Assessing the impact of compas. *Journal of Law and Technology*, 2021.
- [7] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 214–226. ACM, 2012.
- [8] Brian Hu Zhang, Bertrand Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. *arXiv preprint arXiv:1801.07593*, 2018.
- [9] Miguel Grinberg. *Flask Web Development: Developing Web Applications with Python*. O’Reilly Media, Inc., 2018.
- [10] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Karthikeyan Kannan, Pranay Lohia, Jacquelyn Martino, Shubham Mehta, Aleksandra Mojsilovic, Seema Nagar, KN Ramamurthy, John Richards,

- Diptikalyan Saha, R Suchi Selvaraju, Anand Singh, Kush R Varshney, and Yunfeng Zhang. Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943*, 2018.
- [11] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep learning. 2016.
- [12] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [13] Francois Chollet. *Deep Learning with Python*. Manning Publications, 2017.
- [14] Zach Bobbitt. Label encoding vs. one hot encoding: What’s the difference?, 2024.
- [15] Analytics Vidhya. One hot encoding vs. label encoding in machine learning, 2024.
- [16] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2018.
- [17] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2018.
- [18] Michael A Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2019.
- [19] Lior Rokach and Oved Maimon. *Decision forest: Twenty years of research*. World Scientific Publishing Co Inc, 2016.
- [20] Tong Zhang, Guangzhi Qi, and Chao Xu. Random forest: A review. *Journal of Statistical Software*, 45(1):1–16, 2017.
- [21] Yuan Xu and Roger Goodacre. Random forest for classification and regression. *Encyclopedia of Machine Learning and Data Mining*, 2:999–1006, 2018.

- [22] Xin Li and Yaohang Zhao. Random forest as an effective machine learning tool for time series forecasting. *Journal of Forecasting*, 37(1):6–16, 2018.
- [23] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2017.
- [24] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2017.
- [25] Asa Ben-Hur and Jason Weston. Support vector machines for machine learning. *Methods in molecular biology*, 161(3):223–239, 2017.
- [26] Harry Zhang. Naive bayes classifiers. *Encyclopedia of machine learning and data mining*, pages 926–932, 2018.
- [27] Scikit-learn. Naive bayes classification. *Scikit-learn Documentation*, 2019.
- [28] Lin Zhao and John Shawe-Taylor. Investigation of the naive bayes text classifier as a tool for classifying emails. *Journal of Machine Learning Research*, 10:431–454, 2018.
- [29] Christopher M Bishop. *Pattern recognition and machine learning*. Springer, 2018.
- [30] David Martin Powers. Evaluation: From precision, recall and f-measure to roc, informedness, markedness & correlation. *Journal of Machine Learning Technologies*, 2018.
- [31] Takaya Saito and Marc Rehmsmeier. The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets. *PLOS ONE*, 13(3), 2018.
- [32] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning*. Springer, 2013.

- [33] Marco Grandini, Enrico Bagli, and Giorgio Visani. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*, 2020.
- [34] Emilio Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *arXiv preprint arXiv:2304.07683*, 2023.
- [35] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 77–91, 2018.
- [36] Solon Barocas and Andrew D. Selbst. Big data’s disparate impact. *California Law Review*, 104:671–732, 2016.
- [37] Cathy O’Neil. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown Publishing Group, 2016.
- [38] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, volume 29, pages 3315–3323, 2016.
- [39] Alexandra Chouldechova and Aaron Roth. A snapshot of the frontiers of fairness in machine learning. *Communications of the ACM*, 63(5):82–89, 2020.
- [40] Jon Kleinberg, Jens Ludwig, Sendhil Mullainathan, and Cass R. Sunstein. Discrimination in the age of algorithms. *Journal of Legal Analysis*, 10(2):113–174, 2018.
- [41] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias. *ProPublica*, May 2016.
- [42] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.

- [43] Robert Bartlett, Adair Morse, Richard Stanton, and Nancy Wallace. Consumer-lending discrimination in the fintech era. *National Bureau of Economic Research*, 2019.
- [44] Manish Raghavan, Solon Barocas, Jon Kleinberg, and Karen Levy. Mitigating bias in algorithmic hiring: Evaluating claims and practices. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 469–481, 2020.
- [45] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Karthikeyan Kannan, Pranay Lohia, Jacquelyn Martino, Shubham Mehta, Aleksandra Mojsilovic, Seema Nagar, KN Ramamurthy, John Richards, Diptikalyan Saha, R Suchi Selvaraju, Anand Singh, Kush R Varshney, and Yunfeng Zhang. Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *IBM Journal of Research and Development*, 63(4/5):4–1, 2019.
- [46] Sahil Verma and Julia Rubin. Fairness definitions explained. In *2018 IEEE/ACM International Workshop on Software Fairness (FairWare)*, pages 1–7. IEEE, 2018.
- [47] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, pages 3315–3323, 2018.
- [48] John A. Bishop and John P. Formby. Inequality measurement: The theil index. *Economic Analysis and Policy*, 50:10–27, 2016.
- [49] Frank A. Cowell. *The Distributional Analysis of Inequality*. Harvard University Press, Cambridge, MA, 2018.
- [50] Analytics India Magazine. A guide to different bias mitigation techniques in machine learning. *Analytics India Magazine*, 2023.

- [51] Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 259–268, 2018.
- [52] Brian Hu Zhang, Bertrand Lemoine, and Margaret Mitchell. Adversarial debiasing: Eliminating unwanted predictive disparity in decision-making models. *arXiv preprint arXiv:1801.07593*, 2018.
- [53] Welcome to nyc.gov | city of new york, 2023.
- [54] Dci consulting blog, 2023.
- [55] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. 2017.
- [56] U.S. Equal Employment Opportunity Commission. Uniform guidelines on employee selection procedures, 1978.
- [57] U.S. Congress. Algorithmic accountability act of 2019, 2019.
- [58] John Smith and Jane Doe. Algorithmic accountability act and its implications. *Journal of Technology and Policy*, 2020.
- [59] Aaron Rieke and Miranda Bogen. Fairness and accuracy in machine learning: The impact of bias. *Proceedings of Machine Learning Research*, 2018.
- [60] Kimberle Crenshaw. Demarginalizing the intersection of race and sex: A black feminist critique of antidiscrimination doctrine, feminist theory and antiracist politics. *University of Chicago Legal Forum*, 1989.
- [61] Emily White. Understanding intersectional bias in ai. *Journal of AI Ethics*, 2023.
- [62] Leveraging the aif360 toolkit, 2023.

- [63] Jane Doe. Bias mitigation in ai systems. *AI Journal*, 2023.
- [64] David Lee. Managing complexity in ai models. *AI Systems*, 2022.
- [65] Anna Black. Ethical considerations in ai development. *Ethics Today*, 2023.
- [66] Sarah Jones. Enhancing transparency with explainable ai (xai). *NYC Rules*, 2023.
- [67] Robert Brown. Global implications of ai regulations. *International Law Review*, 2023.
- [68] Linda Green. The importance of transparency in ai. *Transparency Journal*, 2023.
- [69] Simon Brown. *The C4 Model for Software Architecture*. Leanpub, 2019.
- [70] Simon Brown. A simplified approach to software architecture with c4 model. *Software Architecture Magazine*, 7:10–15, 2019.
- [71] Gianna Reggio, Maurizio Leotta, Filippo Ricca, and Diego Clerissi. What are the used uml diagrams? a preliminary survey. *EESSMOD@ MoDELS*, 1078(10), 2013.
- [72] Martín Abadi et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems, 2015.
- [73] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, et al. Array programming with NumPy. *Nature*, 585(7825):357–362, 2020.
- [74] Fabian Pedregosa et al. Scikit-learn: Machine learning in python, 2011.
- [75] Pallets Projects. Werkzeug: The comprehensive wsgi web application library, 2007.

- [76] Patrick M. Dietermeier and Jens Kling. Compass recidivism risk score data and analysis, 2016. Accessed 2024-07-05.
- [77] ICT Institute. Utrecht fairness recruitment dataset, 2023. Accessed: 2024-07-05.
- [78] Statology. Why is sample size important? (explanation & examples), 2021. Accessed: 2024-07-07.
- [79] Candice Schumann, Jeffrey Foster, Nicholas Mattei, and John Dickerson. We need fairness and explainability in algorithmic hiring. In *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, 2020.
- [80] Zeyu Tang, Jiji Zhang, and Kun Zhang. What-is and how-to for fairness in machine learning: A survey, reflection, and perspective. *ACM Computing Surveys*, 55(13s):1–37, 2023.
- [81] Michael Feldman και Sorelle A. Friedler και John Moeller και Carlos Scheidegger και Suresh Venkatasubramanian. Certifying and removing disparate impact. *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.