



by Dr. André dos Santos

ABOUT THIS COURSE



PART 1

- Bias
- Fairness



PART 2

- RAI Definitions
- XAI



PART 3

- Regulations
- Standards
- Guidelines

SETTING EXPECTATIONS



FIRST STEPS

RAI deployment is a progressive and evolving practice



FREE RESOURCES

Open source and freely available tools



ONGOING JOURNEY

Continuous effort is required to remain ethical and effective



NO ONE-SIZE-FITS-ALL

RAI solutions are highly domain dependent

PART 1

Bias & Fairness

LINKS



<https://www.wired.com>

OpenAI's Long-Term AI Risk Team Has Disbanded



AI could pose 'extinction-level' threat to humans and the US must intervene, State Dept.-commissioned report warns

<https://www.cnn.com>

CNN Business Markets Tech Media Calculators Videos

40,003.59 0.34%
S&P 500 5,303.27 0.12%
Fear & Greed Index → Green is driving the US market
Latest Market News → Disneyland character performers vote to
Minnesota lawmakers strike minimum t

<https://www.nytimes.com>

The New York Times

Artificial Intelligence > When A.I. Takes Your Voice Google's A.I. Evolution A.I.'s 'Her' Era Arrives OpenAI's Old-Fashioned Library Faces Q

How to Tell if Your A.I. Is Conscious

might indicate the presence of some presence in a machine.

Share full article



653

AI IS MORE THAN GENERATIVE AI

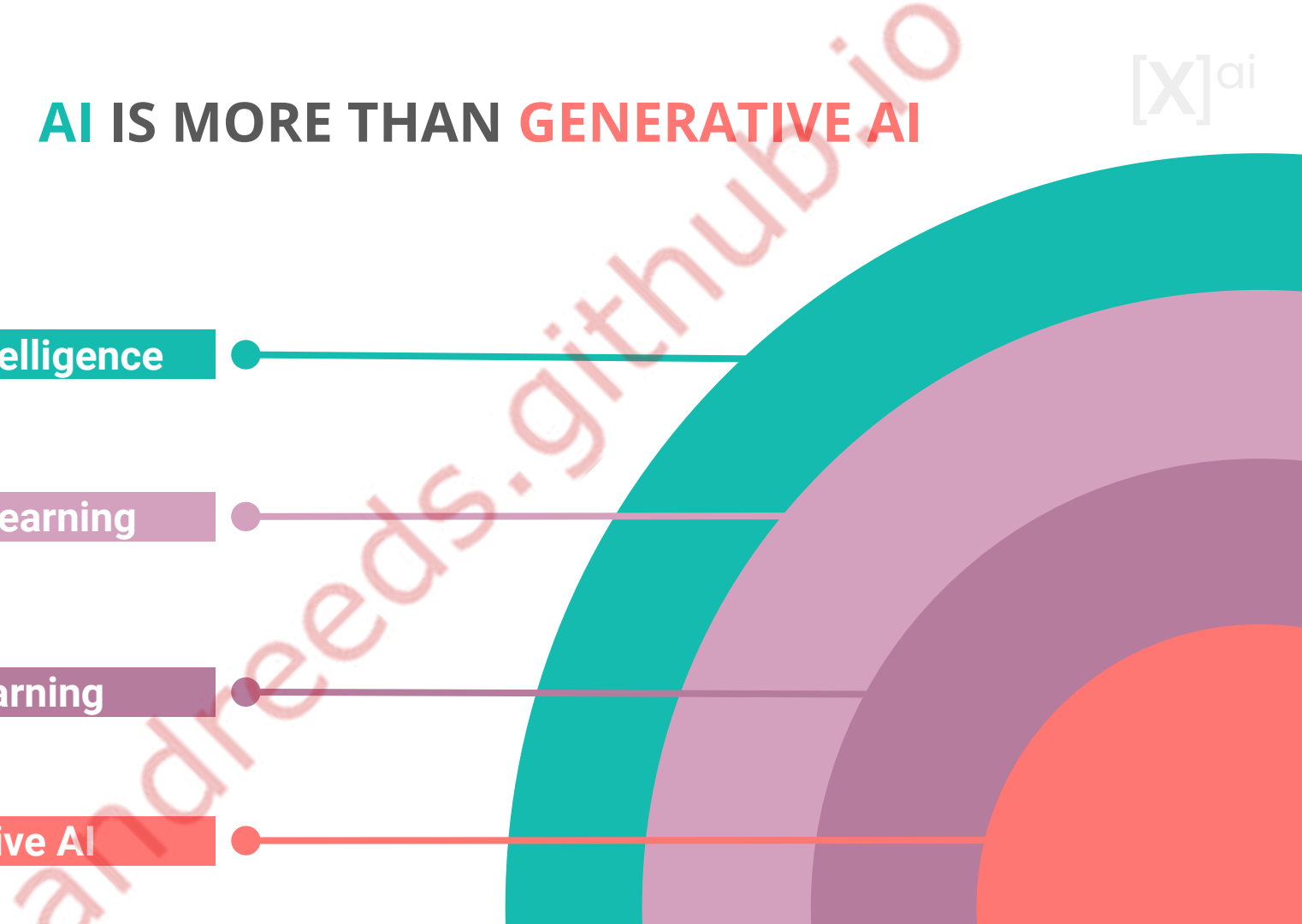
[X]^{ai}

Artificial Intelligence

Machine Learning

Deep Learning

Generative AI



ARTIFICIAL INTELLIGENCE

This new data poisoning tool lets artists fight back against generative AI

The tool, called Nightshade, messes up training data in ways that could cause serious damage to image-generating AI models.

By **Melissa Heikkilä**

October 23, 2023



NIGHTSHADE

Poisoned Concept

Related Prompts

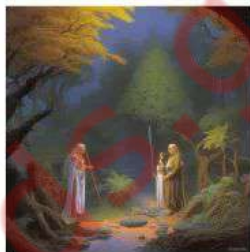
Fantasy art

*A painting by
Michael Whelan*

A dragon

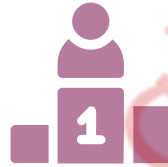
*A castle in the Lord
of the Rings*

**Clean
Model**



**Poisoned
Model**





IDENTIFYING BIAS

“AI bias is the phenomenon that occurs when an AI algorithm produces results that are systemically prejudiced due to erroneous assumptions in the machine learning process.”

—The Internet

AI bias is any systematic error that results in an unfair model

BIAS IN AI

**DATA
AGGREGATION**

**MODEL BUILDING
& IMPLEMENTATION**

BIAS IN AI

FROM DATA AGGREGATION

Historical Bias

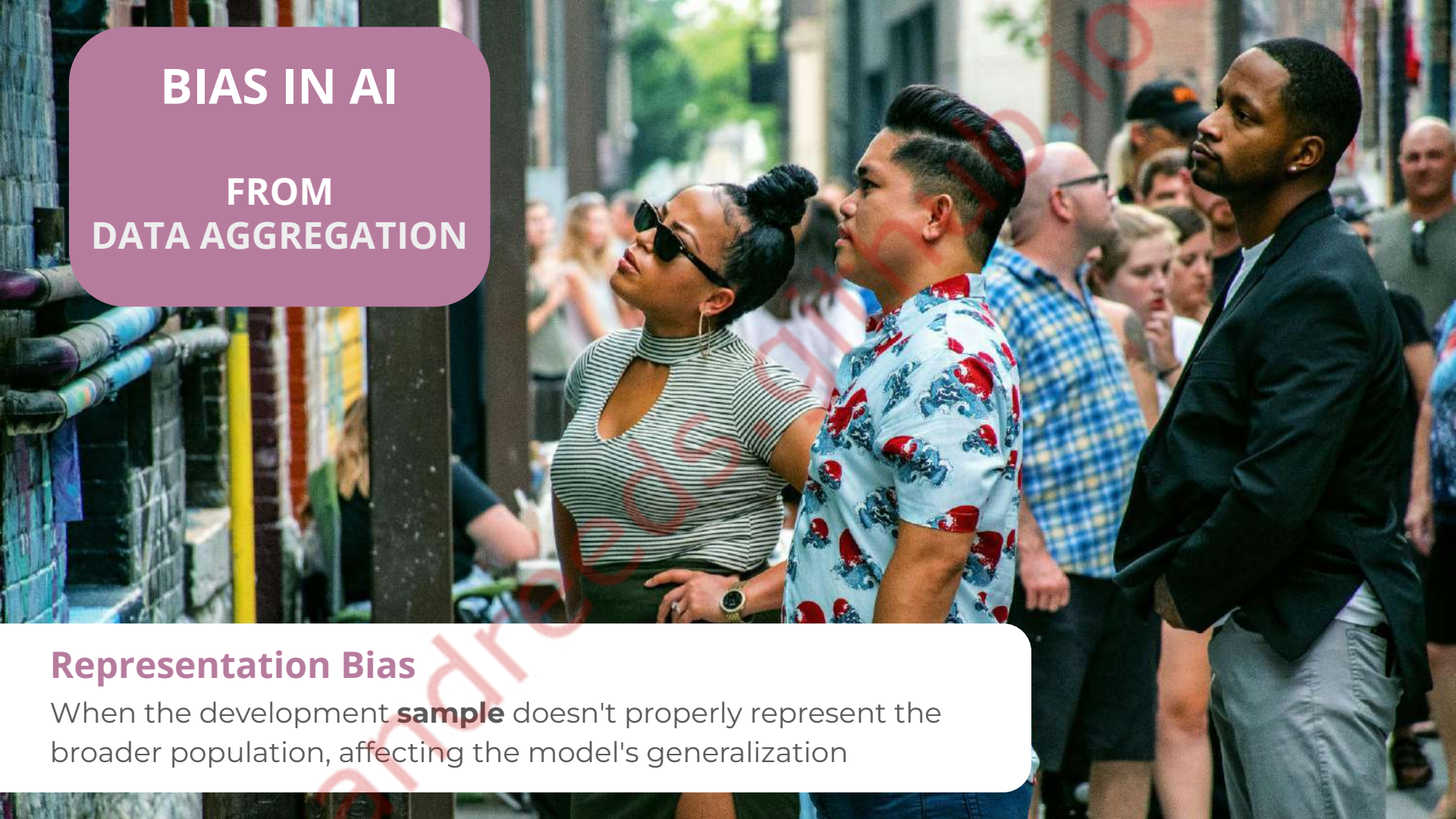
When **data** reflects historical inequalities into model predictions

BIAS IN AI

FROM DATA AGGREGATION

Representation Bias

When the development **sample** doesn't properly represent the broader population, affecting the model's generalization





BIAS IN AI

FROM DATA AGGREGATION

Measurement Bias

When there are flaws in how **data** is **collected** or **measured**, often due to **proxies** that don't accurately capture the desired signals

BIAS IN AI

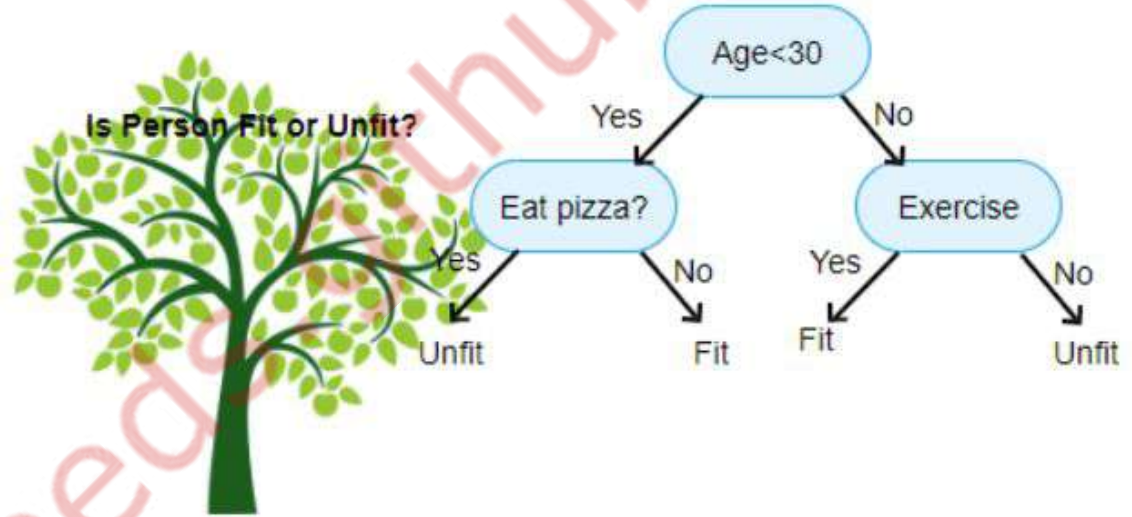
FROM MODEL BUILDING & IMPLEMENTATION

Aggregation Bias

When diverse **groups** are inappropriately **combined** into a single analysis, ignoring meaningful distinctions among them

BIAS IN AI

FROM MODEL BUILDING & IMPLEMENTATION



Learning Bias

When the choice of ML algorithms and their settings may not treat all groups in the data equally

BIAS IN AI

FROM MODEL BUILDING & IMPLEMENTATION



Evaluation Bias

When the data used to evaluate and **benchmark** a model does not represent the actual diversity of the user population

BIAS IN AI

FROM MODEL BUILDING & IMPLEMENTATION

Deployment Bias

When a model is used in real-world applications for which it was not specifically intended





FAIRNESS METRICS

“Everyone agrees that fairness involves treating equal persons equally, and unequal persons unequally, but they do not agree on the standard by which to judge individuals as being equally (or unequally) worthy or deserving.”

—Aristotle

INTRO TO ALGORITHMIC FAIRNESS



Type of Fairness

Within the context of an AI/ML model used for classification



Classification

One type of AI/ML task



Fairness Metrics

Quantitative measures used to assess the fairness of AI models



Pure Math

Understanding the nuances of fairness metrics is essential

BASIC FAIRNESS METRIC TERMINOLOGY



Sensitive attribute

Attribute needing special ethical, legal, or social consideration



Proxy Attribute

Attribute correlated with a sensitive attribute



Parity

Observational measure ensuring metrics are independent of defined groups



Confusion Matrix

Tool measuring accuracy in predictions

CONFUSION METRIC

		True condition			
		Condition positive	Condition negative	Prevalence $= \frac{\Sigma \text{Condition positive}}{\Sigma \text{Total population}}$	Accuracy (ACC) = $\frac{\Sigma \text{True positive} + \Sigma \text{True negative}}{\Sigma \text{Total population}}$
Predicted condition	Total population				
	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\Sigma \text{True positive}}{\Sigma \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{False positive}}{\Sigma \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\Sigma \text{False negative}}{\Sigma \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\Sigma \text{True negative}}{\Sigma \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\Sigma \text{True positive}}{\Sigma \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\Sigma \text{False positive}}{\Sigma \text{Condition negative}}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) $= \frac{\text{LR+}}{\text{LR-}}$
		False negative rate (FNR), Miss rate $= \frac{\Sigma \text{False negative}}{\Sigma \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\Sigma \text{True negative}}{\Sigma \text{Condition negative}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$	
					F ₁ score = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

LINKS



Datapoint editor

Performance & Fairness

Features

500 datapoints loaded



Configure

Ground Truth Feature
over_50k

WHAT IS GROUND TRUTH?

The feature that your model is trying to predict. [More](#)

Cost Ratio (FP/FN)

1

WHAT IS COST RATIO?

The cost of false positives relative to false negatives. Required for optimization. [More](#)

Slice by

<none>

WHAT DOES SLICING DO?

Shows the model's performance on datapoints grouped by each value of the selected feature.

Fairness

Apply an optimization strategy

Select a strategy to automatically set classification thresholds, based on the set cost ratio and data slices. Manually altering thresholds or changing cost ratio will revert the strategy to 'custom thresholds'.

- ☒ Custom thresholds ⓘ
- ☐ Single threshold ⓘ
- ☐ Demographic parity ⓘ
- ☐ Equal opportunity ⓘ
- ☐ Equal accuracy ⓘ
- ☐ Group thresholds ⓘ

Explore overall performance ⓘ

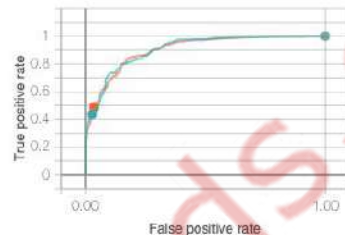
Sort by

Count

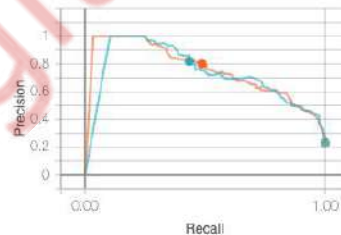


Feature Value	Count	Model	Threshold ⓘ	False Positives (%)	False Negatives (%)	Accuracy (%)	F1
▼ All datapoints	500	1		2.2	13.0	84.8	0.57
	2	2		2.8	11.8	85.4	0.61

ROC curve (AUCs: 0.90, 0.90) ⓘ



PR curve (AUCs: 0.71, 0.74) ⓘ



Confusion Matrix ⓘ

	Predicted Yes	Predicted No	Total
1			
Actual Yes	10.0% (50)	13.0% (65)	23.0% (115)
Actual No	2.2% (11)	74.8% (374)	77.0% (385)
Total	12.2% (61)	87.8% (439)	
2			
Actual Yes	11.2% (56)	11.8% (59)	23.0% (115)
Actual No	2.8% (14)	74.2% (371)	77.0% (385)
Total	14.0% (70)	86.0% (430)	

LINKS



FAIRNESS METRICS EXAMPLES



UNAWARENESS

Ignoring sensitive attributes to achieve fairness in decisions



DEMOGRAPHIC PARITY

Equal decision rates across groups regardless of outcome



EQUAL OPPORTUNITY & EQUALIZED ODDS

Fairness through equal true positive rates and error rates across groups



PREDICTIVE VALUE PARITY

Equal predictive accuracy across different groups



INDIVIDUAL FAIRNESS

Similar treatment for individuals with similar attributes.

Simulating loan decisions for different groups

Drag the black threshold bars left or right to change the cut-offs for loans.

Click on different preset loan strategies.

Loan Strategy

Maximize profit with:

MAX PROFIT

No constraints.

GROUP UNAWARE

Blue and orange thresholds are the same.

DEMOGRAPHIC PARITY

Same fractions blue / orange loans

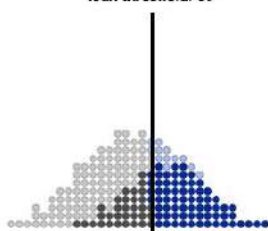
EQUAL OPPORTUNITY

Same fractions blue / orange loans to people who can pay them off

Blue Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 59

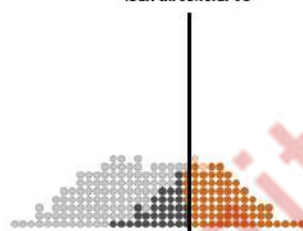


denied loan / would default
denied loan / would pay back

Orange Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 53



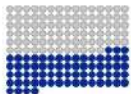
denied loan / would default
denied loan / would pay back

Total profit = 30400

Equal Opportunity

Among people who would pay back a loan, blue and orange groups do equally well. This choice is almost as profitable as demographic parity, and about as many people get loans overall.

Correct 78%
loans granted to paying applicants and denied to defaulters

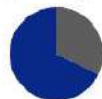


Incorrect 22%
loans denied to paying applicants and granted to defaulters



True Positive Rate 68%

percentage of paying applications getting loans



Profit: 11700

Positive Rate 40%

percentage of all applications getting loans



Correct 83%
loans granted to paying applicants and denied to defaulters



Incorrect 17%
loans denied to paying applicants and granted to defaulters



True Positive Rate 68%

percentage of paying applications getting loans



Profit: 18700

Positive Rate 35%

percentage of all applications getting loans



LINKS



OPEN SOURCE LIBRARIES

Here are a few common open-source libraries and tools on AI Fairness:

AIF360 [Python/R]

- ❑ IBM Research
- ❑ Last update: 1 month ago
- ❑ Bias mitigation algorithms
- ❑ Fairness metrics
- ❑ Video tutorial

Fairlearn [Python]

- ❑ Microsoft, now community driven
- ❑ Last update: 1 week ago
- ❑ Bias mitigation algorithms
- ❑ Fairness metrics
- ❑ Complete Documentation - Notebook Examples

Aequitas [Python]

- ❑ Carnegie Mellon University
- ❑ Last update: 2 weeks ago
- ❑ Complete Toolkit
- ❑ Complete Documentation - Notebook Examples
- ❑ Aequitas's license does not allow commercial use.

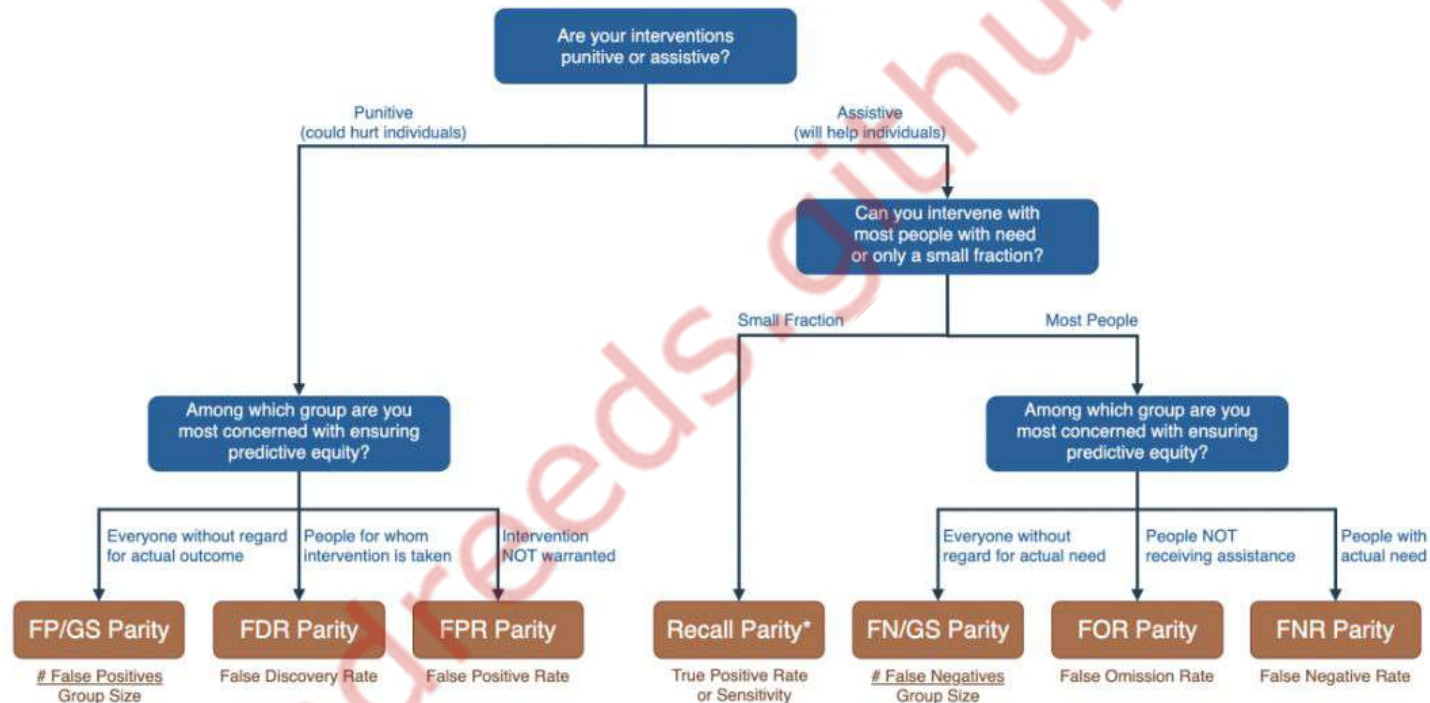
VerifyML [Python]

- ❑ Cylynx
- ❑ Winner Global Veritas Challenge
- ❑ Complete Toolkit
- ❑ Last update: 2 years ago
- ❑ Code Demo

LINKS



FAIRNESS TREE (Zoomed in)



LINKS



Source: Aequitas

Bias and Fairness Audit Toolkit

The Bias Report is powered by [Aequitas](#), an open-source bias audit toolkit for machine learning developers, analysts, and policymakers to audit machine learning models for discrimination and bias, and make informed and equitable decisions around developing and deploying predictive risk-assessment tools.



See an [example report](#) on COMPAS risk assessment scores.

Or try out the audit tool using your own data or one of our sample data sets.

[Get Started!](#)

LINKS



PART 2

Responsible AI

LINKS



RESPONSIBLE AI DEFINITIONS

Whats is the difference between
Understandability, Interpretability,
Comprehensibility, Explainability,
and **Transparency?**



RESPONSIBLE AI KEY CONCEPTS



UNDERSTANDABILITY

What?



COMPREHENSIBILITY

How?



INTERPRETABILITY

Why? When?
Who? Which?



EXPLAINABILITY

Why? What?
Who? When?
How? Which?



TRANSPARENCY

No questions
needed.



UNDERSTANDABILITY

You know to get a *latte* from the machine



COMPREHENSIBILITY

You know that to get a *latte* from the machine, it'll have to:

Step 1: Brew and Pour the Espresso

Step 2: Steam the Milk

Step 3: Pour the Milk

Step 4: Top Off with Foam



INTERPRETABILITY

You know that if your *latte* is not great, it's because the machine got the milk too hot - 40 degrees for small drinks up to 155 for big ones. Likely they do not account for the lag with the thermometer.



EXPLAINABILITY (XAI)

You know that to get a *latte* from the machine, you see it

Step 1: Brew and pour 2 shots of espresso.

Step 2: Steam 1/2 cup milk to 150 F

Step 3: Pour the steamed milk over the espresso, using a spoon to hold back the foam.

Step 4: Top off the drink with the reserved foam



TRANSPARENCY

You ask the barista for a *latte*.

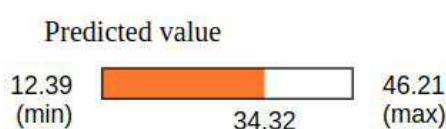




XAI METHODS

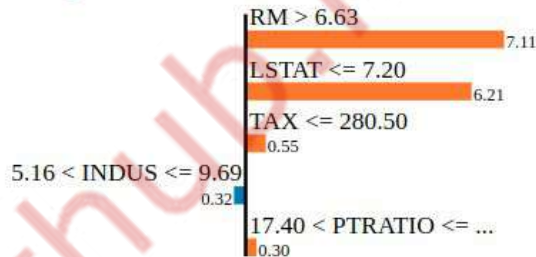
LIME

Tabular



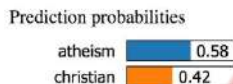
negative

positive



Feature	Value
RM	7.27
LSTAT	6.05
TAX	254.00
INDUS	6.41
PTRATIO	17.60

Text



atheism

christian

Posting 0.15

Host 0.14

NNTP 0.11

edu 0.04

have 0.01

There 0.01

Text with highlighted words

From: johnchad@triton.unm.edu (jchadwic)
 Subject: Another request for Darwin Fish
 Organization: University of New Mexico, Albuquerque
 Lines: 11
 NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the DARWIN fish.
 This is the same question I have and I have not seen an answer on the net. If anyone has a contact please post on the net or email me.

Image



(a) Original Image



(b) Explaining Electric guitar



(c) Explaining Acoustic guitar



(d) Explaining Labrador

LINKS



SHAP



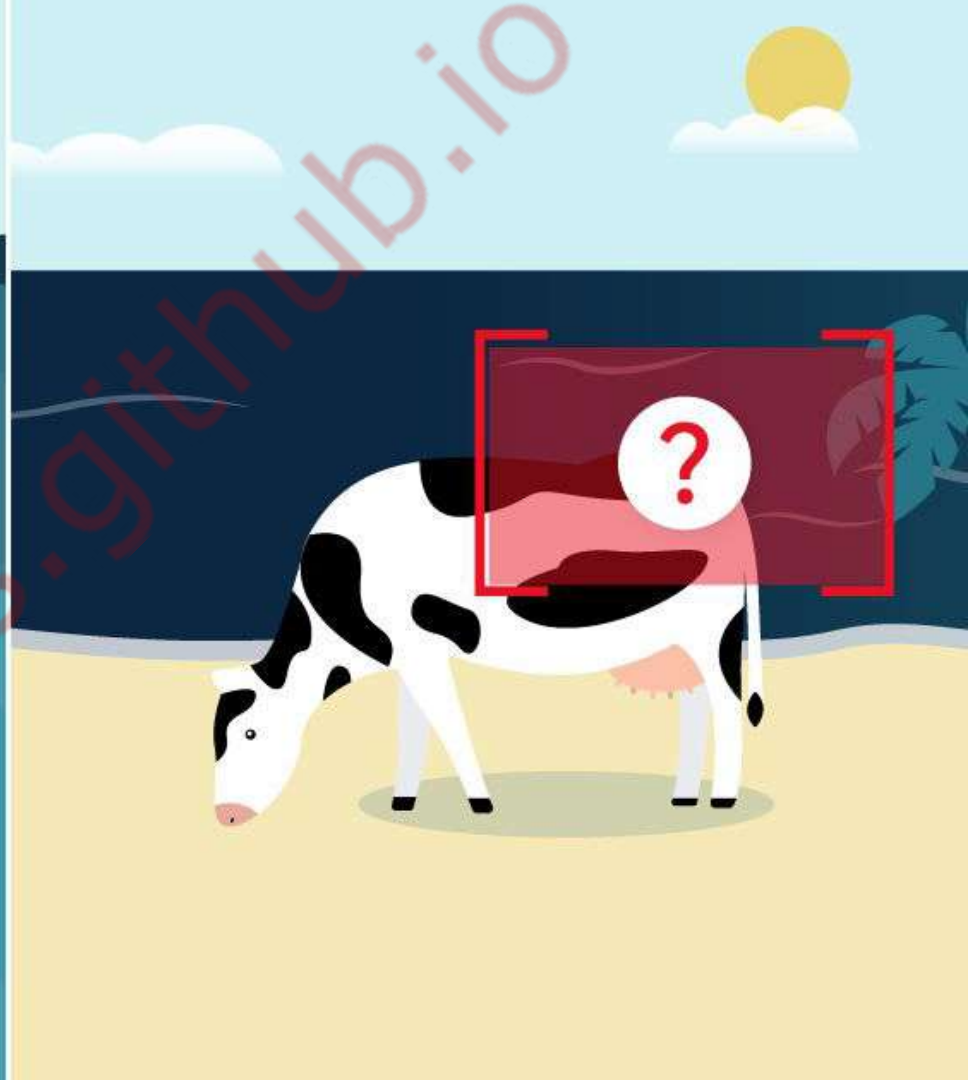
LINKS



**TRANSPARENT
AI**

VERSUS

DL + XAI



COFFEE BREAK

PART 3

AI Regulations, Standards & Guidelines

LINKS





REGULATIONS



United Kingdom

Legislation & regulation: Online Safety Bill (2022, draft); Data Protection and Digital Information Bill (2023, draft)

Standards: Algorithmic Transparency Standard (Central Digital Data Office, 2021)

Principles: A pro-innovation approach to AI regulation (2023)



Canada

Legislation & regulation: Directive on Automated Decision-Making (2019) Bill C-27; Digital Charter

Implementation Act, including AI and Data Act (AIDA) (2022, draft)

Standards: CAN-ASC-6.2: Accessible and Equitable AI Systems (2023, draft)

Principles: Canada's Digital Charter (2019)

Oversight: Minister; Proposed AI and Data Commissioner



European Union

Legislation & regulation: Proposed EU AI Act (2021, draft); Updates to the EU Product Liability Directive (2022, draft); AI Liability Directive (2022, draft); EU's Digital Services Act

Standards: CEN/CENELEC standards for AI and related data (forthcoming)

Principles: Ethics guidelines on AI (2018)

Oversight: Proposed European Artificial Intelligence Board



United States

Legislation & regulation: Federal Trade Commission Act, for deceptive practices from deepfakes or chatbots (1914); Algorithmic Accountability Act (US AAA) (2022, draft)

Standards: NIST AI Risk Management Framework (2023)

Principles: Blueprint for an AI Bill of Rights (2023)



China

Legislation & regulation: Chinese Internet Information Service Algorithmic Recommendation Provisions (2021); Opinion on Strengthening the Ethics and Governance of Science and Technology (2022)

Standards: National Standards for Autonomous Vehicle Testing (2018)

Principles: New Generation AI Ethics Specifications (2019); New Generation AI Code of Ethics (2021); Internet Information Service Algorithmic Recommendation Management Provisions (2021)



Brazil

Legislation & regulation: Report and proposed substitute text for draft bills 5051/2019, 21/2020 and 872/2021 (2022, draft); Bill 705 on the compatibility of AI use in the public sector with ESG practices (2022, draft)

Standards: Incorporation of international standards National standards by the Brazilian Association of Technical Norms (ABNT)

Principles: Art. 3 of the proposed substitute text for draft bills 5051/2019, 21/2020 and 872/2021 (2022, draft)



Intergovernmental Organisations

Legislation & regulation: Council of Europe Convention on AI, Human Rights, Democracy and the Rule of Law (2023, draft)

Standards: ISO 31000 Risk management (2009, 2018); ISO/IEC 23053:2022 Framework for AI Systems Using Machine Learning (ML) (2022)

Principles: OECD Recommendation of the Council on AI (2019); UNESCO Recommendation on the Ethics of AI (2021)

LINKS



EU

AI ACT

AI systems are broadly defined, with a focus on autonomy.

Key takeaways

- Risk-Based Approach
- Banned Practices
- Transparency Obligations
- Market Surveillance

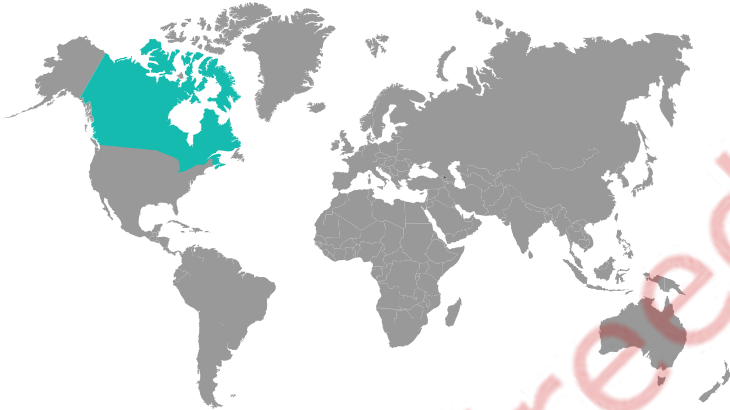


CANADA

BILL C-27

Bill comprised of three acts, one is the AI & Data Act (AIDA)

- Committee Stage
- Government institutions excluded
- Two purposes: regulate trade and prohibit harm
- Prevention of biased **outputs**
- High impact system not clear
- Minister of Innovation Powers





GUIDELINES

EU AI ACT COMPLIANCE CHECKER

LINKS



Compliance Checker

Discover how the AI Act will affect you in 10 minutes by answering a series of straightforward questions.

How will the EU AI Act affect my AI system?

Please complete this form for each individual AI system used in your organisation.

Entity type

Which kind of entity is your organisation?

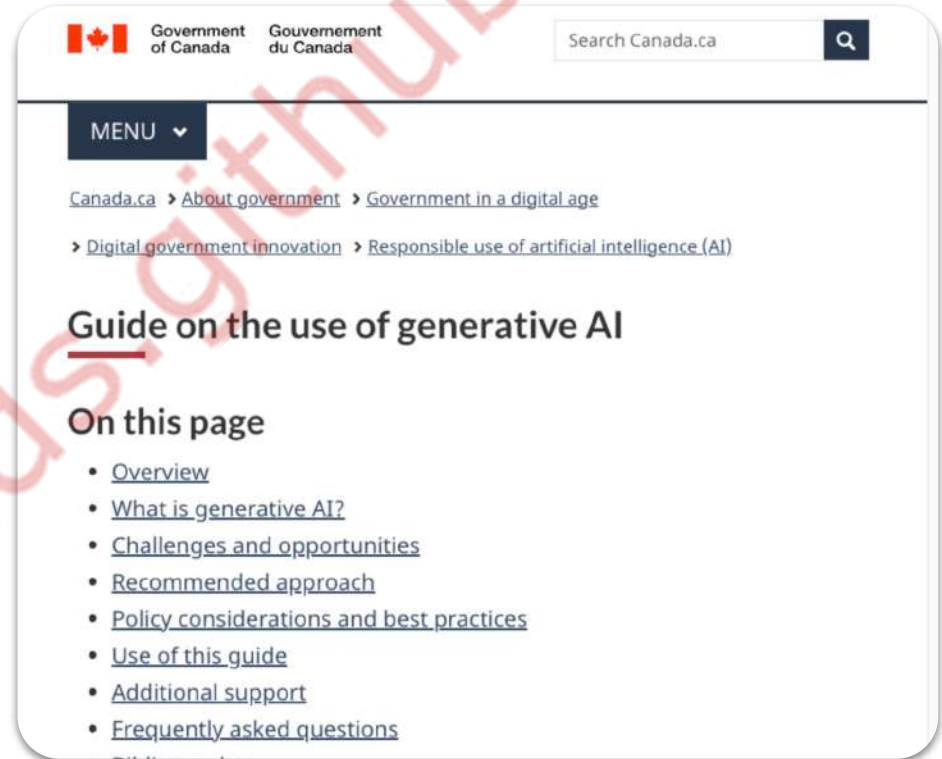
- ☐ Provider
- ☐ Deployer
- ☐ Distributor
- ☐ Importer
- ☐ Product manufacturer
- ☐ Authorised representative

Your results

No results yet - Please complete the web form above to see your results.

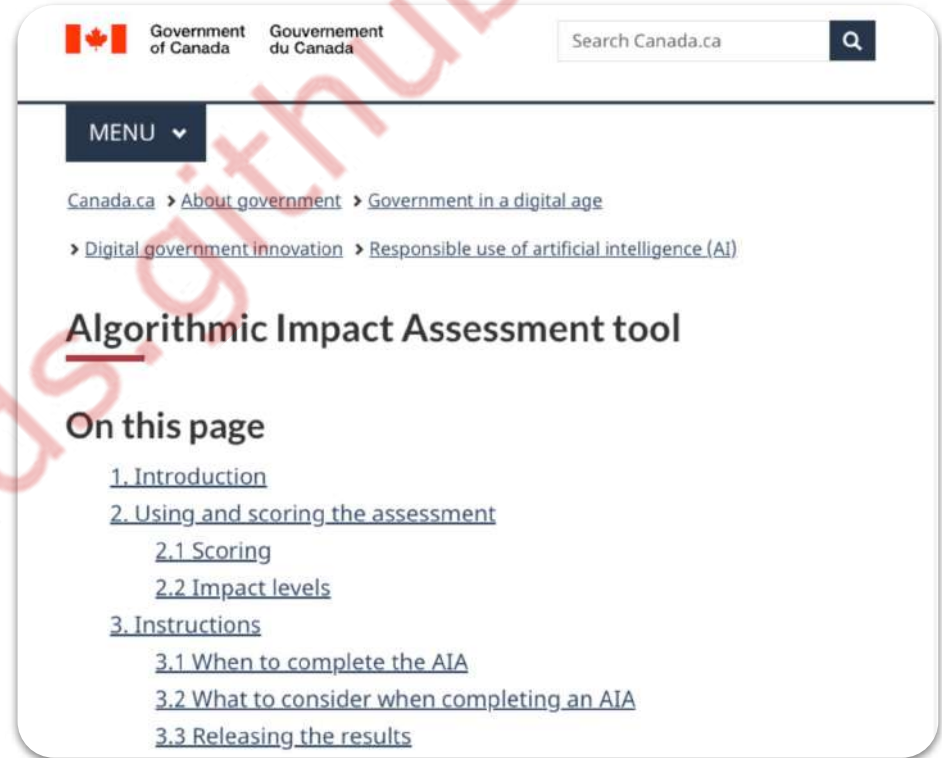
GOVERNMENT OF CANADA'S GUIDE ON THE USE OF GENERATIVE AI

LINKS



GOVERNMENT OF CANADA'S ALGORITHMIC IMPACT ASSESSMENT TOOL

LINKS



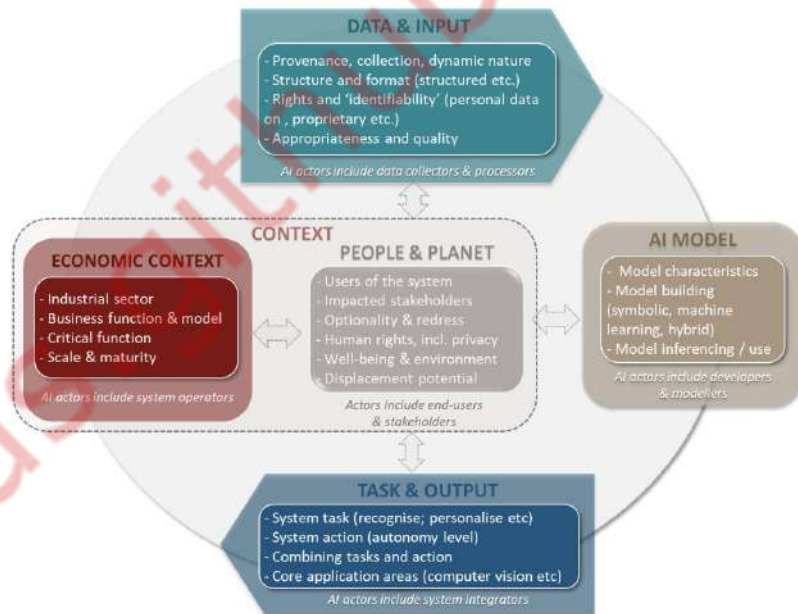
OECD FRAMEWORK FOR THE CLASSIFICATION OF AI SYSTEMS

OECD DIGITAL ECONOMY PAPERS

February 2022 No. 323



Figure 2. Characteristics per classification dimension and key actor(s) involved



Note: Actors are illustrative, non-exhaustive and notably relevant to accountability.

Source: Based on the work of ONE AI and the AI system lifecycle work of AIGO (OECD, 2019_[2]).

LINKS





GPAI / THE GLOBAL PARTNERSHIP ON ARTIFICIAL INTELLIGENCE

Follow us



[HOME](#)

[ABOUT](#)

[COMMUNITY ▾](#)

[OUR WORK ▾](#)

[EVENTS](#)





STANDARDS

IEEE

ETHICALLY ALIGNED DESIGN

LINKS



Digital Governance
Standards Institute

ETHICAL DESIGN AND USE OF AUTOMATED DECISION SYSTEMS



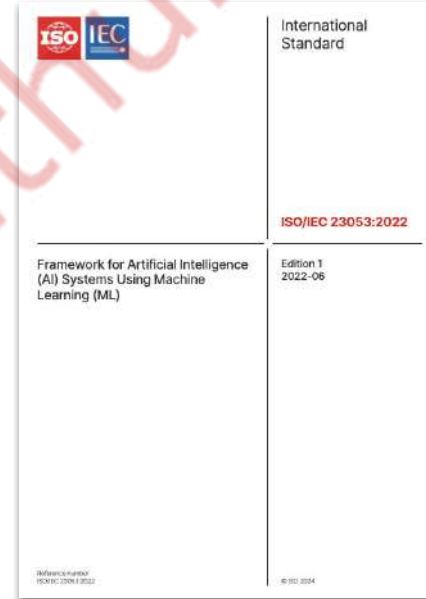
Ethical Design and Use of Automated Decisions Systems

CAN/CIOSC 101:2019 (Reaffirmed 2021-10)

LINKS



ISO/IEC 23053:2022 FRAMEWORK FOR AI SYSTEMS USING ML



LINKS



CONCLUSION



PART 1

- Bias
- Fairness



PART 2

- RAI Definitions
- XAI



PART 3

- Regulations
- Standards
- Guidelines

LINKS

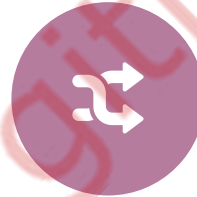


MORE ESSENTIALS



PRIVACY

Adversarial attacks



CAUSALITY

Correlation is not necessarily causation



AI & PEOPLE

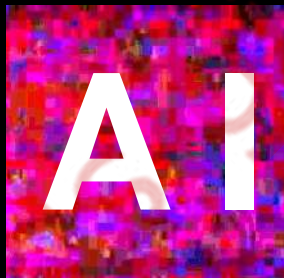
How AI affects and is perceived by people



FRAMEWORKS

Ethics charter

THE
TRUTH
ABOUT



Separating **Fact**
from **Fiction**

by Dr. André dos **Santos**

•

4

7

Re

•

•

[Export BibTeX Citation](#)

Bookmark



BE HUMAN.



[X]ai plain

PROFESSIONAL
DEVELOPMENT

AI EXPLAINED

AI SAFETY

AI LEADERSHIP

EXPERTS
NETWORKING

AI CONSULTING

