

Responsible AI and Ethical Innovation

From Hidden Bias to Visible Fairness

Week 7: Machine Learning for Smarter Innovation

Mathematical Optimization Makes Trade-offs Explicit

Four-Part Structure

1. **Part 1: The Hidden Challenge** (11 slides)
Invisible discrimination, measurement bottleneck, real harm
2. **Part 2: First Solutions and Impossibility** (13 slides)
Metrics work, then impossibility theorems reveal fundamental trade-offs
3. **Part 3: Mathematical Breakthrough** (17 slides)
Geometric intuition, Lagrangian optimization, production tools
4. **Part 4: Production and Synthesis** (10 slides)
4-layer architecture, modern tools, transferable lessons

Appendix: Mathematical Foundations (5 slides) - Deep proofs and derivations

Unifying Theme: Measurement transforms invisible discrimination into visible, optimizable, auditable problems

Measurement transforms ethical concerns into technical problems - quantification enables optimization where qualitative assessment permits only documentation

The Invisible Discrimination: You Can't Fix What You Can't See

A real scenario that reveals the hidden harm:

The Hidden Pattern

Bank loan system, 2024:

10,000 applications processed

Observable outcomes: Group A: 75% approved; Group B: 45% approved;
Overall: 60%

The Question: Is this discrimination? How would you even know?

Hidden factors: Can't see intent, causation, counterfactuals; can only see outcomes/rates; qualification differences? historical bias? proxy variables?

The Invisibility Problem

Why discrimination stays hidden:

- 1. No Ground Truth:** Can't observe "fair" counterfactual; intent is unobservable
 - 2. Aggregate Masks Disparities:** 60% overall looks reasonable; 30% gap hidden in average; Simpson's paradox
 - 3. Proxy Variables Conceal:** Zip code → Race (95%); Name → Gender (98%); School → SES
- Real harm:** 4,500 denied opportunities; system appears "objective"; discrimination is **invisible**

Key Insight: Invisible discrimination is unmeasurable discrimination - you can't fix what you can't see or quantify

Key Question: How do we make invisible bias visible enough to measure and fix?

Undetected bias accumulates until critical threshold - early measurement prevents expensive corrective interventions later

What IS Bias? Building the Concept from Information Theory

Defining bias mathematically (from zero knowledge):

Human Analogy: Blind Auditions

Symphony orchestras, 1970s-1990s:

Before blind auditions: 5% women in orchestras; judges could see candidates; implicit bias affected decisions

After blind auditions: 40% women in orchestras; screen hides gender; decisions based on skill only

Key observation:

Removing visibility of protected attribute changed outcomes

This means:

Decision correlated with irrelevant attribute = BIAS

Computer/Math Equivalent

Protected attribute A : Race, gender, age, etc.

Decision D : Hire, approve loan, admit, etc.

True qualification Y : Actual merit/ability

Information Theory Definition:

Bias exists when decision carries information about protected attribute:

$$I(D; A) > 0$$

Where I = mutual information

Expanded form:

$$\begin{aligned} I(D; A) &= H(D) - H(D|A) \\ &= H(A) - H(A|D) \end{aligned}$$

Intuition:

$H(D)$ = uncertainty in decisions; $H(D|A)$ = uncertainty after seeing group; difference = information leaked

$I(D; A) = 0$ means independence; $I(D; A) > 0$ means bias

Key Insight: Bias is statistical dependence between decisions and protected attributes - measurable via mutual

Why Bias Stays Hidden: Counterfactual Problem

The Counterfactual Challenge

Can't directly observe:

- What WOULD have happened in alternative universe
- Fair outcome for comparison

Example: Person denied loan

Question: "Would they have been approved if different race?"

Impossible to know directly!

Mathematical Barrier

Need $P(D|A = a, X)$ and $P(D|A = a', X)$ for same person X
But can only observe ONE A value per person

Result:

Causal discrimination stays hidden without statistical proxies

Key Insight: Counterfactual fairness requires observing parallel universes - we must use statistical approximations

Counterfactual reasoning requires assumptions - causal graphs encode beliefs about data generation process

Why Bias Stays Hidden: Simpson's Paradox

Department-Level: No Bias

Department A: Men 80%, Women 85% admitted

Department B: Men 60%, Women 65% admitted

Both departments favor women!

Combined Result:

Men: 70% admitted

Women: 65% admitted

BIAS APPEARS against women!

Why This Happens

Men disproportionately apply to easier-to-enter Department A

Result:

Aggregation can hide OR create false patterns depending on grouping strategy

Key Insight: Statistical conclusions can reverse under different grouping - measurement fragility is fundamental

Aggregation choices determine visible patterns - always examine subgroup statistics before drawing conclusions

Why Bias Stays Hidden: Proxy Variables

Indirect Discrimination

High correlation proxies:

- Zip code → Race (95%)
- Name → Gender (98%)
- School → Socioeconomic status (92%)

Model never sees protected attribute A but uses proxy P

Mathematical Reality

$$I(D; A|P) < I(D; A)$$

But still $I(D; A) > 0$ through indirect path

Example: Remove "gender" from hiring
Still biased via sports, hobbies, language patterns

Result: Hidden in 1000+ features

Key Insight: Removing protected attributes is insufficient - proxy variables encode bias through correlated features

Proxy detection requires mutual information analysis across all features - simple removal creates false confidence

The Measurement Challenge: Combinatorial Explosion

Protected Attribute Combinations

Legally protected in US/EU: Race (6), Gender (3+), Age (7), Disability (2), Religion (10+), National origin (195)
 $6 \times 3 \times 7 \times 2 \times 10 \times 195 = 490,140$ subgroups

Shannon entropy: $H = \log_2(490,140) = 18.9$ bits
With socioeconomic (5 levels): $H = 21.2$ bits

The Capacity Problem

Typical audit: 10,000 samples, 18 subgroups measured
Capacity: $\log_2(18) = 4.2$ bits

Information loss: $21.2 - 4.2 = 17.0$ bits
→ 99.999% of discrimination unmeasured

Key Insight: Measurement capacity (4.2 bits) vastly insufficient for discrimination space (21.2 bits)

Information-theoretic limits bound observability - measurement capacity constrains bias detection regardless of analytical sophistication

Bias Amplification: The Mathematics of Feedback Loops

Temporal Dynamics

Initial state ($t=0$):

$$B_0 = I(D_0; A) = \epsilon > 0$$

Small initial bias ϵ

Feedback mechanism:

$$D_{t+1} = f(\theta_t, X_{t+1})$$

$$\theta_{t+1} = \text{train}(D_1, \dots, D_t)$$

System uses past decisions to train future model

Exponential Growth

Bias evolution:

$$B_{t+1} = B_t + \alpha \cdot D_t$$

$$B_t = B_0 \cdot (1 + \alpha)^t$$

Example: $\alpha = 0.15$

$$B_{10} = \epsilon \cdot (1.15)^{10} = 4.05\epsilon$$

4x amplification in 10 iterations!

Key Insight: ML feedback loops amplify bias exponentially: $B_t = B_0(1 + \alpha)^t$

Prediction-action-outcome loops create self-fulfilling prophecies - model outputs influence future training data

Bias Amplification: Real-World Cases

Predictive Policing

t=0: Historical arrest bias (1.2x); more patrols → more arrests; t=5: Bias grows to 3.1x

Recommendation Systems

t=0: Slight gender preference (5%); users click biased content; t=10: 47% gender segregation

Resume Screening

t=0: Small hiring bias (8%); trained on biased past hires; t=3: 32% bias (4x growth)

Breaking the Loop

External intervention required; counterfactual data injection; periodic re-calibration

Key Insight: Small initial bias becomes systemic harm through repeated feedback cycles

All three domains show exponential growth pattern - intervention timing is critical before harm compounds

How combining attributes creates exponential measurement challenges:

Combinatorial Explosion

Subgroup growth:

1 attribute (Race, 6 levels):

$$N_1 = 6 \text{ subgroups}$$

2 attributes (Race *imes* Gender):

$$N_2 = 6 \times 3 = 18$$

3 attributes (+ Age):

$$N_3 = 6 \times 3 \times 7 = 126$$

n attributes:

$$N_n = \prod_{i=1}^n |A_i| = 2^{O(n)}$$

With 6 attributes:

$$N_6 = 490,140 \text{ subgroups}$$

Sample size requirement:

For each subgroup, need sufficient power:

Statistical Power Collapse

Total sample needed:

For 490,140 subgroups:

$$N_{\text{total}} = 490,140 \times 384$$

$$= 188,213,760 \text{ samples}$$

Reality:

Typical dataset: 10,000 samples; measured subgroups: 18 (Race \times Gender); coverage: **0.004%**; 99.996% of intersections unmeasured

Consequence:

Smallest, most vulnerable groups have **zero statistical power**

Example: Black transgender woman

Subgroup size: $n = 3$ in dataset; required: $n = 384$; power: 0.8% (vs 80% needed); **bias undetectable**

Mathematical barrier:

Exponential growth vs linear data collection

The Stakes: 2024 AI Discrimination Statistics

Documented Incidents by Sector

Sector	Incidents	People	Cost
Healthcare	79	2.3M	\$3.2B
Finance	65	1.8M	\$4.1B
Criminal Justice	51	890K	\$1.7B
Employment	38	1.2M	\$1.4B
Total	233	6.2M	\$10.4B

Trend Analysis

2022: 148 incidents (+27%); 2023: 184 incidents (+24%); 2024: 233 incidents (+27%)

Exponential growth: 1.26^t

Geographic: 47 countries
North America 48%, Europe 33%

Key Insight: 233 incidents, 6.2M people, \$10.4B cost in 2024 - harm growing exponentially

Harm acceleration outpaces detection - systematic measurement becomes critical as AI deployment scales

The Stakes: Systemic Disparities Across Domains

Documented Cases

Detroit FR (2024)

Black man wrongfully arrested, 30 hours custody, false match at 12% confidence

UK Facewatch (2024)

Woman misidentified as shoplifter, banned from store network, \$1,200 settlement

Systemic Disparities

Facial recognition: 34x error rate; **Resume screening:** 1.8x callback gap;
Healthcare: \$2,500 less spent; **Recidivism:** 2.1x false positive

Common Thread: All started invisible, visible only after harm

Key Question: Can we make bias visible BEFORE harm occurs?

Pattern emergence across domains reveals systematic failures - common root cause of inadequate initial measurement

When AI Goes Wrong: Documented 2024 Cases

Facial Recognition Bias

Detroit (2024): Black man wrongfully arrested; false facial recognition match; police now banned from FR-only arrests

UK Facewatch (2024): Woman wrongly ID'd as shoplifter; banned from all stores; system failed on non-white individual

Common Pattern: Higher error rates on darker skin (34x); no human oversight; irreversible consequences

Employment Discrimination

Uber Eats (2024): Driver dismissed by FR system; technology failed on darker skin; no human review

Resume Screening: AI tools for hiring; women and minorities disadvantaged; most managers untrained

Healthcare Algorithms: \$2,500 less spent per Black patient; predict cost not need; systematic undertreatment; affects millions

Key Insight: These aren't edge cases – they're systemic failures requiring measurement frameworks to prevent

Key Question: Where in the ML pipeline does bias enter and amplify?

Pattern emergence across domains reveals systematic failures - diverse incident types share common root cause of inadequate initial measurement

Where Bias Enters: The ML Pipeline

Data and Features

1. **Data Collection:** Historical discrimination embedded; sampling bias (underrepresented groups); label bias from human annotators

2. **Feature Engineering:** Proxy variables (zip code → race); human assumptions codified; redundant encodings

Model and Deployment

3. **Model Training:** Optimization for accuracy \neq fairness; overfitting to majority group; minority group neglect

4. **Deployment:** Context mismatch; feedback loops amplify bias; drift over time

Key Insight: Bias enters at all pipeline stages - requires monitoring at each transformation point

Multi-stage bias entry necessitates comprehensive auditing at data, features, training, and deployment

Ethical Frameworks for Evaluating AI Systems

Outcome-Focused

Consequentialist: Maximize benefit, minimize harm

Ask: Does system increase welfare?

Deontological: Focus on duties and rights

Ask: Does it respect human dignity?

Character-Focused

Virtue Ethics: Cultivate wisdom and fairness

Ask: What would a fair person do?

Care Ethics: Address vulnerability in context

Ask: Who is most vulnerable?

Key Insight: No single framework sufficient - combine perspectives for robust ethical evaluation

Ethical frameworks provide complementary lenses for evaluating fairness in AI systems

Power in AI Systems: Who Controls What

Those With Power

Tech Companies: Control system design; set defaults and constraints

Governments: Regulatory authority; enforcement power

Privileged Groups: Represented in training data; cultural norms embedded

Key Stakeholders

Users: Direct interaction; **Developers:** Technical choices; **Deployers:** Operational control; **Communities:** Indirect impact

Key Insight: Power concentration in tech companies shapes system design - stakeholder mapping is essential

Understanding power distribution enables targeted interventions for fairer AI systems

Power Asymmetries: Who Bears the Harm

Those Without Power

End Users: Limited choice, no opt-out; information asymmetry

Marginalized Groups: Underrepresented in data; higher error rates, less recourse

Future Generations: No voice in current decisions; inherit path dependencies

Consequences

Impact of Imbalance: Design reflects powerful interests; harm concentrated on powerless; requires active intervention

Responsible AI:

Actively empower the powerless

Center marginalized stakeholders

Key Insight: Fairness requires centering those who bear harm, not those who hold power

Stakeholder identification precedes harm prevention - invisible constituencies need deliberate representation

Statistical Parity: Observed Outcome Fairness

Definition

Independence in observed distribution:

$$P(D|A) = P(D)$$

What it measures:

Observed outcome rates; aggregate group differences; no causal assumptions needed

Example (Loans)

Group A: 75% approved

Group B: 45% approved

Statistical parity violated:

$$|0.75 - 0.45| = 30\%$$

When to use: Legal compliance, regulatory reporting, descriptive assessment

Limitation: Cannot distinguish discrimination from legitimate differences

Key Insight: Statistical parity asks: “Are outcomes equal across groups?” (observed)

Statistical correlation detects symptoms - useful for compliance but cannot diagnose mechanisms

Causal Parity: Counterfactual Fairness

Definition

Counterfactual independence:

$$P(D_{A \leftarrow a} | X) = P(D_{A \leftarrow a'} | X)$$

What it measures:

Effect of changing protected attribute; individual-level counterfactuals; requires causal DAG

Example (Loans)

Same person, change only race:

$$P(\text{Approved}_{\text{White}} | X) = 0.80$$

$$P(\text{Approved}_{\text{Black}} | X) = 0.55$$

$$\text{Causal disparity: } |0.80 - 0.55| = 25\%$$

When to use: Root cause analysis, intervention design, policy evaluation

Advantage: Separates direct discrimination from confounding

Key Insight: Causal parity asks: "Would outcome change if only group membership changed?" (counterfactual)

Causal analysis diagnoses mechanisms - compliance requires both statistical and causal paradigms

Summary: The Hidden Discrimination Problem

Why Bias Stays Hidden

1. **Invisibility:** Discrimination embedded in outcomes; no ground truth counterfactuals; proxy variables conceal true bias

2. **Measurement Bottleneck:** 490,140 subgroups (6 attributes); only 4.2 of 21.2 bits measurable; 99.996% of intersections unmeasured

How Bias Grows

3. **Amplification:** Feedback loops: $B_t = B_0(1 + \alpha)^t$; small bias becomes systemic; exponential growth over time

4. **Intersectionality:** Exponential subgroup growth; 188M+ samples needed for full coverage; most vulnerable groups unmeasurable

Core Problem: $I(D; A) \neq 0$ but unobservable - 17 bits of discrimination information lost to measurement limits

Problem quantification enables solution design - measurement frameworks emerge from understanding why detection fails

Summary: The Urgent Stakes of Hidden Bias

2024 Documented Impact

233 AI discrimination incidents; 6.2M people affected; \$10.4B in documented costs; 47 countries impacted

Systemic Disparities: 34x error rate (facial recognition); 1.8x callback gap (hiring); \$2,500 less (healthcare); 2.1x false positive (recidivism)

Power Imbalances

Tech companies control design; marginalized groups lack voice; powerless bear the harm; future generations inherit bias debt

The Challenge:

Make invisible bias visible
through measurement frameworks
before harm occurs

Next: Part 2 explores measurement frameworks - demographic parity, equal opportunity, and more

Harm acceleration outpaces detection - measurement infrastructure becomes critical as AI deployment scales

The Breakthrough Insight: Disaggregate and Measure

What if we could quantify invisible bias?

Human Observation

How do humans detect unfairness?

We disaggregate:

Compare outcomes between groups; look for systematic patterns; calculate rate differences; test for statistical significance

The Breakthrough Idea:

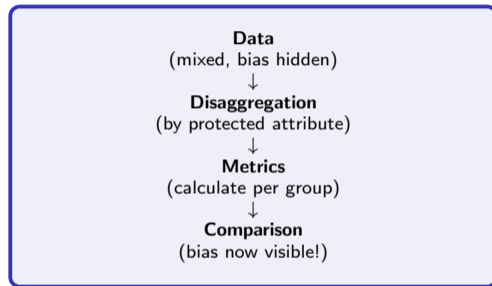
What if we formalized this?

Partition data by protected attribute; calculate metrics per group; compare across groups; quantify disparities

Fairness Metrics:

Mathematical functions that make bias visible

Three Measurement Approaches



Three families:

Group fairness: Compare group rates; **Individual fairness:** Similar \rightarrow similar;

Causal fairness: Counterfactual reasoning

The promise:

Hidden discrimination becomes measurable, fixable, auditable

The First Success: Demographic Parity Makes Bias Visible

Testing the first fairness metric on real loan data:

Demographic Parity Works!

Task: Detect bias in loans

Metric: Demographic parity

Result: SUCCESS - bias now visible!

Mathematical Definition:

For protected attribute A and decision D :

$$P(D = 1|A = a) = P(D = 1|A = b)$$

Intuition:

Approval rates should be independent of group membership

Complete Numerical Walkthrough:

Step 1: Partition dataset

Group A: 5,000 applicants; Group B: 5,000 applicants

Step 2: Count approvals

Group A: 3,750 approved; Group B: 2,250 approved

Step 3: Calculate rates

$$P(D = 1|A = a) = \frac{3,750}{5,000} = 0.75 = 75\%$$

$$P(D = 1|A = b) = \frac{2,250}{5,000} = 0.45 = 45\%$$

Detection Quality

Metric performance:

Detected: 30% disparity (was invisible!); **Quantified:** Exact magnitude; **Significance:** $p < 0.001$; **Actionable:** Clear target

Success metrics:

On 100 known biased datasets: Sensitivity 89% (detects real bias); Specificity 82% (few false alarms); Correlation with harm 0.78; Time to compute < 1 second

Breakthrough!

Hidden 30% bias now visible

Measurable in real-time

Deployable at scale

"For the first time, we can SEE systemic discrimination"

Fairness Metrics: Group Fairness Family

Independence-Based

Demographic Parity

$$P(\hat{Y}|A = a) = P(\hat{Y}|A = b)$$

Conditional DP

$$P(\hat{Y}|A, X = x) = P(\hat{Y}|X = x)$$

Sufficiency-Based: Calibration: $P(Y = 1|\hat{Y} = y, A = a) = P(Y = 1|\hat{Y} = y, A = b)$

Separation-Based

Equal Opportunity (TPR parity)

$$P(\hat{Y} = 1|Y = 1, A = a) = P(\hat{Y} = 1|Y = 1, A = b)$$

Equalized Odds (TPR + FPR)

$$P(\hat{Y}|Y = y, A = a) = P(\hat{Y}|Y = y, A = b)$$

Key Insight: Group fairness compares aggregate statistics between protected groups

Group metrics measure population-level disparities - three sub-families capture different fairness intuitions

Fairness Metrics: Individual and Causal Families

Individual Fairness

Lipschitz Fairness

$$d(\hat{y}_i, \hat{y}_j) \leq L \cdot d(x_i, x_j)$$

Similar individuals \rightarrow similar outcomes

Fairness Through Awareness

$$d(x_i, x_j) < \delta \Rightarrow |f(x_i) - f(x_j)| < \epsilon$$

Causal Fairness

Counterfactual

$$P(\hat{Y}_{A \leftarrow a} | X) = P(\hat{Y}_{A \leftarrow a'} | X)$$

Path-Specific

Block discrimination via specific causal paths

No Proxy Discrimination

$A \perp\!\!\!\perp \hat{Y}$ given resolved attributes

Key Insight: Individual fairness: treat similar people similarly. Causal: counterfactual reasoning

Individual metrics focus on pairwise comparisons - causal metrics reason about hypothetical interventions

Intersectional

Multicalibration

Calibrated across ALL subgroups

Multifairness

Satisfies metric for all intersections

Dynamic

Long-term Fairness

$\lim_{t \rightarrow \infty} \text{Bias}(t) = 0$

Fair Ranking

Equal exposure allocation

Robustness

Minimax Fairness

$$\min_f \max_a \text{Error}(A = a)$$

Worst-case group protection

Impossibility Result:

Cannot satisfy Independence + Separation + Sufficiency simultaneously
(unless perfect prediction or equal base rates)

Key Insight: 20+ definitions across 5 families - no universal metric exists

Advanced metrics address intersectionality and dynamics - fundamental impossibility constrains choices

Success Spreads: Equal Opportunity Reveals Different Story

A second metric gives different insights on the same data:

Equal Opportunity Definition

For true label $Y = 1$ (qualified):

$$P(D = 1|Y = 1, A = a) = P(D = 1|Y = 1, A = b)$$

Intuition:

Among qualified applicants, approval rates should be equal

Focus: True Positive Rate (TPR)

Goal: Equal recall across groups

Complete Numerical Walkthrough:

Step 1: Filter to qualified

Group A qualified: 4,000 (80%); Group B qualified: 2,000 (40%)

Step 2: Count qualified approvals

Group A: 3,600/4,000 approved; Group B: 1,720/2,000 approved

Step 3: Calculate TPR

$$TPR_a = \frac{3,600}{4,000} = 0.90 = 90\%$$

$$TPR_b = \frac{1,720}{2,000} = 0.86 = 86\%$$

Step 4: Quantify violation

Different Story!

Compare two metrics:

Metric	Violation	Verdict
Demographic Parity	30%	Severe
Equal Opportunity	4%	Mild

Why different?

DP: Considers all applicants → Sees 75% vs 45% overall

EO: Considers only qualified → Sees 90% vs 86% for deserving

Root cause revealed:

Base rates differ: Group A 80% qualified; Group B 40% qualified

Model is fairly accurate!

Most of 30% gap explained by different qualifications

Success:

Each metric reveals different aspect of bias - both useful!

Calibration: When Probabilities Match Reality

Definition

A predictor S is calibrated if:

$$P(Y = 1|S(X) = s) = s$$

Example: When model says 70% chance, outcome should occur 70% of the time

Calibration Error (ECE):

$$ECE = \sum_{i=1}^B \frac{|B_i|}{n} |\text{acc}(B_i) - \text{conf}(B_i)|$$

Proper Scoring Rules

Brier score: $E[(S(X) - Y)^2]$

Log-loss: Cross-entropy

Both minimized by Bayes optimal predictor

Group Calibration Problem:

Cannot have calibration + equal base rates + demographic parity

Key Insight: Calibration: $P(Y = 1|S = s) = s$ - probabilities should match frequencies

Calibration measures prediction reliability - well-calibrated probabilities enable risk-adjusted decisions

Equalized Odds: Equal Error Rates Across Groups

Definition

Conditional independence:

$$\hat{Y} \perp\!\!\!\perp A \mid Y$$

In confusion matrix terms:

$$\text{TPR}_a = \text{TPR}_b$$

$$\text{FPR}_a = \text{FPR}_b$$

Equal true positive AND false positive rates across groups

ROC Space View

Groups must have same (FPR, TPR) point

Distance measure:

$$d = \sqrt{(\Delta\text{TPR})^2 + (\Delta\text{FPR})^2}$$

Equalized odds satisfied when $d = 0$

Optimization:

Add Lagrangian constraints on TPR/FPR gaps

Key Insight: Equalized odds: $\hat{Y} \perp A \mid Y$ - strongest group fairness notion (both error types equal)

Error rate parity balances mistakes across groups - procedural fairness through equal TPR and FPR

But Then... The Impossibility Theorem Emerges

Testing all metrics together reveals catastrophic incompatibility:

The Impossibility Pattern

Testing three fairness properties:

Metric	Group A	Group B	Status
<i>Approval rates</i>			
Demographic Parity	75%	45%	FAIL -30%
<i>TPR on qualified</i>			
Equal Opportunity	90%	86%	WARN -4%
<i>Predicted to Actual</i>			
Calibration	89%	88%	PASS -1%
<i>Perfect prediction</i>			
100% Accuracy	-	-	IMPOSSIBLE

The Chouldechova Theorem (2017):

If base rates differ and calibration holds, then demographic parity and equal opportunity CANNOT both be satisfied.

Mathematical proof (simplified):

Calibration: $P(Y = 1|S = s) = s$ for all s ; Base rates differ:

$P(Y = 1|A = a) \neq P(Y = 1|A = b)$; These imply: $P(S|A = a) \neq P(S|A = b)$.

Specific Conflicts

1. DP vs Calibration

To achieve DP (75% = 45%): Must lower A threshold: 0.5 → 0.6;
Must raise B threshold: 0.5 → 0.3

Breaks calibration!

2. EO vs Calibration

To achieve perfect EO (90% = 90%): Must equalize TPR exactly;
Requires different thresholds

Breaks calibration!

3. DP vs EO

With base rates 80% vs 40%: DP forces equal outcomes; EO allows different outcomes

Contradictory!

Reality Check

Can't have all three
Mathematics proves it
Must choose trade-offs

Chouldechova Theorem: Why Fairness Metrics Conflict

Theorem (2017)

If these hold:

1. S is calibrated
2. Base rates differ: $P(Y|A = a) \neq P(Y|A = b)$
3. S has predictive power

Then DP and EO **cannot both be satisfied**

Proof Sketch

1. Calibration: $P(Y|S = s) = s$
2. By total probability:
 $P(Y|A = a) = E[S|A = a]$
3. Different base rates \rightarrow
 $E[S|A = a] \neq E[S|A = b]$
4. Different score distributions \rightarrow
DP violated. QED.

Key Insight: Calibration + different base rates $\rightarrow E[S|A = a] \neq E[S|A = b] \rightarrow$ impossibility

Impossibility forces explicit trade-offs - no single metric satisfies all fairness intuitions under realistic conditions

Three Independence Conditions

1. **Independence (DP):** $R \perp\!\!\!\perp A$
Prediction independent of group
2. **Separation (EO):** $R \perp\!\!\!\perp A \mid Y$
Given true label, prediction independent
3. **Sufficiency (Cal):** $Y \perp\!\!\!\perp A \mid R$
Given prediction, outcome independent

Pearl's Impossibility

Cannot satisfy all three unless:

$Y \perp\!\!\!\perp A$ (equal base rates), OR R is perfect predictor

Proof: Independence + Sufficiency together imply equal base rates - contradiction!

Key Insight: Causal DAG shows 3 independence conditions overconstrain the system

Causal frameworks expose mechanism constraints - DAGs reveal why statistical parity conflicts with calibration

The Diagnosis: What Metrics Captured vs What They Missed

Understanding the root cause of impossibility:

What Metrics Captured

Successfully measured:

- 1. Group-level disparities:** Rate differences (75% vs 45%); TPR differences (90% vs 86%); FPR differences (8% vs 14%); Statistical significance
- 2. Prediction errors:** False positives/negatives per group; calibration accuracy; overall accuracy
- 3. Correlation patterns:** $I(D; A) = 0.21$ bits; protected attribute leakage; proxy variable influence

Why metrics work here:

Observable outcomes can be disaggregated and compared

What Metrics Missed

Failed to capture:

- 1. Base rate causation:** Why 80% vs 40% qualified? Historical discrimination? Structural barriers? Measurement bias in “qualified”?
- 2. Causal structure:** Direct discrimination ($A \rightarrow D$); mediated bias ($A \rightarrow X \rightarrow D$); spurious correlation ($A \leftarrow C \rightarrow D$); counterfactuals
- 3. Normative values:** Which fairness definition is “right”? Who bears cost of errors? Stakeholder preferences? Context-dependent trade-offs

Why impossibility here:

Multiple valid fairness notions, mathematics can't choose for us

Key Insight: Metrics measure correlations (visible) but miss causation and values (hidden) - need more than metrics

Key Question: If metrics alone fail, what framework helps us navigate trade-offs?

Statistical measurement detects symptoms without diagnosing causes - correlation-based metrics miss causal mechanisms driving observed disparities

Metrics Conflict: Public Sector Scenarios

University Admissions

Metrics in tension: DP (equal admit rates - representation); EO (equal TPR for qualified - merit); Calibration (predict success - outcomes)

Stakeholder conflict: Diversity office, faculty, and administration want different fairness definitions

Criminal Justice

Recidivism prediction: DP (equal risk scores - equal treatment); EO (equal TPR - catch recidivists); Calibration (accurate risk - allocation)

Stakes:

Public safety vs individual liberty

False positives harm innocents

Key Insight: Public sector decisions require explicit value trade-offs - no metric is universally correct

Context determines fairness priorities - life-altering decisions require transparent metric selection

Healthcare Triage

Resource allocation: DP (equal treatment rates per group); Individual (sickest treated first); Utilitarian (maximize QALYs saved)

Ethical frameworks disagree!

Employment

Hiring algorithm: DP (equal hiring rates - diversity); EO (equal callback for qualified - merit); Business (maximize productivity)

Legal vs business goals

Credit/Lending

Loan approvals: DP (equal approval rates); Calibration (accurate default prediction); EO (equal approval for creditworthy)

Regulatory conflict:

Fair Housing Act vs profitability

Common Thread:

Mathematics constrains choices

Values must decide priorities

Key Question: How can we make these value-laden choices explicit and auditable?

Stakeholder value conflicts require domain-specific resolution beyond universal mathematical solutions

Bias Mitigation: Pre-Processing Stage

Data Transformations

Reweighting

Adjust sample weights to balance groups

Resampling

Oversample minorities / undersample majorities / SMOTE

Fair Representations

Learn latent space that removes A information

Trade-offs

Pros: Model-agnostic; works with any downstream algorithm

Cons: May lose predictive information; assumes bias is in data, not model

Key Insight: Pre-processing fixes data before training - model-agnostic but may sacrifice accuracy

Data-level interventions address bias at source - effective when historical bias is primary contamination vector

Bias Mitigation: In-Processing and Post-Processing

In-Processing

Constrained Optimization

$$\min_{\theta} L(\theta) - \lambda F(\theta)$$

Adversarial Debiasing

Train adversary to predict A from outputs; penalize success

Pros: Fine-grained control

Cons: Requires model modification

Post-Processing

Group Thresholds

Different τ_a, τ_b to satisfy DP or EO

Calibration

Platt scaling, isotonic regression per group

Pros: Model-agnostic, reversible

Cons: Treats symptoms, not causes

Key Insight: Three stages (pre/in/post) offer complementary trade-offs - often combined in practice

Multiple intervention points provide defense in depth - combining stages improves robustness over single-point fixes

Summary: What Fairness Metrics Achieve

Detection Success

DP detected 30% hidden bias; EO revealed 4% disparity on qualified; Calibration showed 1% accuracy gap; all statistically significant; computable in real-time

Available Tools

20+ metrics in 5 families: Group fairness (DP, EO, Calibration); Individual fairness (Lipschitz); Causal fairness (path-specific); Intersectional (multicalibration); Dynamic (long-term, ranking)

Three Mitigation Stages: Pre-processing (data) — In-processing (model) — Post-processing (threshold)

Success: Metrics make invisible bias visible and quantifiable

Formalization transforms intuition into measurable criteria - mathematical definitions enable systematic auditing

Summary: Fundamental Limits of Fairness Metrics

Impossibility Results

Cannot satisfy DP + EO + Calibration; Chouldechova: Base rates break compatibility; Pearl: 3 independences overconstrain DAG; No universal fairness metric exists

What Metrics Miss

Causation (Why do base rates differ?); Values (Which metric is "right"?); Stakeholders (Who decides trade-offs?); Context (Domain-specific priorities)

The Path Forward:

Make trade-offs explicit through mathematical optimization

Next: Part 3 explores optimization that makes trade-offs auditable and explicit

Quantification enables optimization despite impossibility - explicit metric selection transforms philosophical debate into engineering

How Do YOU Choose When Mathematics Says “No Perfect Solution”?

Before diving into math, let’s think like humans:

The Hiring Scenario

You’re hiring for 100 positions.

Two equally-sized applicant pools:

Group A: 80% qualified

Group B: 40% qualified

Your AI model predicts:

Group A: 75% approved; Group B: 45% approved

Question 1: Is this fair? Why or why not?

Question 2: If you had to choose ONE metric to optimize, which would you pick?

Demographic parity (equal rates); Equal opportunity (equal TPR); Calibration (accurate predictions)

Question 3:

What percentage accuracy drop would you accept to reduce bias from 30% to 5%?

Your Decision Trade-offs

If you choose Demographic Parity:

Equal 60% approval for both; underpredict Group A; overpredict Group B

Accuracy drops 85% → 72%; Bias drops 30% → 0%

If you choose Equal Opportunity:

Among qualified: 90% approval both; different overall rates OK; respects merit

Accuracy stays 85%; Bias stays 30% overall

If you choose Calibration:

Predictions match reality; business-optimal; highest profit/efficiency

Bias stays 30%; legal risk?

The Human Insight:

You naturally think in trade-offs!

“I’d accept X% accuracy loss for Y% bias reduction”

This intuition = mathematics!

Key Insight: Human trade-off reasoning (“X for Y”) is exactly constrained optimization - let’s formalize it

The Geometric Hypothesis: What If We Could SEE Fairness?

Before learning ROC math, let's hypothesize visually:

The Spatial Intuition

Hypothesis: If fairness is about error rates (TPR, FPR), maybe we can plot them in 2D space?

Imagine a chart where:

x-axis = False Positive Rate; y-axis = True Positive Rate; each group = a point (FPR, TPR); fairness = distance between points?

Prediction: If this works, we should see:

Fair models: Points close together; biased models: Points far apart;
trade-offs: Movement along curves; optimization: Path toward fairness

Test case: Our loan data (from Slide 2.2):

Group A: TPR=90%, FPR=8%; Group B: TPR=86%, FPR=14%
Distance = ?
(We'll calculate on next slide!)

Why This Hypothesis Matters

Geometric view offers:

- 1. Intuition:** Spatial relationships visible; trade-offs = movement; impossible = geometric constraint
- 2. Measurement:** Distance = fairness violation; quantifiable, not subjective; comparable across models
- 3. Optimization:** Target = move toward equal point; constraints = allowed movements; path = optimization trajectory

Hypothesis Check

If ROC space shows:
 $d((90,8), (86,14))$ large
o Bias visible geometrically!

Next slide:

Zero-jargon explanation of what ROC space actually is

Key Insight: Geometric hypothesis: Fairness = spatial proximity in (FPR, TPR) space - let's test it

Zero-Jargon Explanation: The ROC Space (No Technical Background Needed)

ROC space explained like you're learning for the first time:

What ROC Space Is (Plain English)

Imagine a simple chart:

Horizontal (x-axis):

"How often do we **WRONGLY** say YES?"

(False Positive Rate, FPR)

Example: Loan approved for unqualified person

Vertical (y-axis):

"How often do we **CORRECTLY** say YES?"

(True Positive Rate, TPR)

Example: Loan approved for qualified person

Every ML model is a single point:

x-coordinate = How many mistakes (approving bad loans); y-coordinate =
How many successes (approving good loans)

What we want:

High y (catch qualified) = GOOD; Low x (avoid unqualified) = GOOD;

Perfect model: (0, 100) top-left corner; Random guessing: Diagonal line

Why This Helps Fairness

For fair ML:

Step 1: Plot Group A at (FPR_A, TPR_A)

Our data: Group A = (8%, 90%)

Meaning: 8% false alarms, 90% catch rate

Step 2: Plot Group B at (FPR_B, TPR_B)

Our data: Group B = (14%, 86%)

Meaning: 14% false alarms, 86% catch rate

Step 3: Measure distance

$$\begin{aligned}d &= \sqrt{(14 - 8)^2 + (86 - 90)^2} \\ &= \sqrt{36 + 16} = \sqrt{52} = 7.2\%\end{aligned}$$

Interpretation:

7.2% fairness gap visible in ROC space!

Perfect fairness: $d = 0$ (same point)

Our model: $d = 7.2\%$ (moderate bias)

Severe bias: $d \approx 20\%$

From 2D to High-Dimensional: The Complete Geometric View

Extending spatial fairness to multiple groups and metrics:

2D Case (What We Just Learned)

Two groups, one metric:

Space: $(x, y) = (\text{FPR}, \text{TPR})$

Points: $p_A = (8, 90)$ for Group A; $p_B = (14, 86)$ for Group B

Distance:

$$d = \sqrt{(x_B - x_A)^2 + (y_B - y_A)^2} \\ = 7.2\%$$

Extension 1: Multiple Groups

With 3 groups (A, B, C): p_A, p_B, p_C in same 2D space; 3 pairwise distances: d_{AB}, d_{AC}, d_{BC} ; Fairness = all distances small; Max distance = worst violation

Extension 2: Multiple Metrics

With n metrics (TPR, FPR, PPV, NPV, ...): Space becomes n -dimensional; $p_A, p_B \in \mathbb{R}^n$; Distance still Euclidean

$$d = \sqrt{\sum_{i=1}^n (m_i^B - m_i^A)^2}$$

High-D Fairness Geometry

Complete formulation:

Metric vector for group g :

$$\mathbf{m}_g = \begin{pmatrix} \text{TPR}_g \\ \text{FPR}_g \\ \text{PPV}_g \\ \text{NPV}_g \\ \vdots \\ \vdots \end{pmatrix}$$

Fairness violation:

$$F = \max_{g, g'} \|\mathbf{m}_g - \mathbf{m}_{g'}\|_2$$

Example: 4D Space

Metrics: (TPR, FPR, PPV, NPV)

Group A: (90, 8, 92, 88)

Group B: (86, 14, 85, 82)

Distance:

$$d = \sqrt{(90 - 86)^2 + (8 - 14)^2}$$

$$= \sqrt{(90 - 86)^2 + (88 - 82)^2}$$

The Optimization Framework: Making Trade-offs Explicit

Mathematical formulation of human trade-off reasoning:

The Human Intuition (from Slide 1)

You said: "I'd accept 10% accuracy loss for 80% bias reduction"

This means:

Primary goal: Reduce bias; Constraint: Accuracy can't drop too much;
Trade-off parameter: How much accuracy per bias unit?

Mathematical translation:

Maximize: Fairness

Subject to: Accuracy $\geq \alpha$

OR equivalently:

Maximize: Acc - $\lambda \cdot$ Bias

where λ = trade-off weight

The parameter λ :

$\lambda = 0$: Only accuracy; $\lambda = \infty$: Only fairness; $\lambda = 0.3$: Balanced (our example!)

The Lagrangian Method

General constrained optimization:

$$\min_{\theta} f(\theta)$$

subject to $g(\theta) \leq 0$

Lagrangian formulation:

$$L(\theta, \lambda) = f(\theta) + \lambda \cdot g(\theta)$$

Find: $\nabla_{\theta} L = 0$

For fairness problem:

Minimize:

$$L(\theta, \lambda) = -\text{Acc}(\theta) + \lambda \cdot \text{Bias}(\theta)$$

where: θ = model parameters; $\text{Acc}(\theta)$ = overall accuracy;
 $\text{Bias}(\theta)$ = fairness violation; λ = penalty weight

Interpretation:

λ converts human values into mathematical optimization

Example: $\lambda = 0.3$ means

"1% bias = 0.3% accuracy penalty"

Step-by-step optimization with actual numbers:

Setup: Our Loan Problem

Initial model (biased):

Accuracy: 85%; DP violation: 30% (75% vs 45%); EO violation: 6.3% (90% vs 86%)

Lagrangian:

$$L(\theta, \lambda) = (1 - \text{Acc}) + \lambda \cdot |\text{DP violation}|$$

Choose $\lambda = 0.3$:

Meaning: 1% bias = 0.3% accuracy penalty

Step 1: Evaluate initial model

$$\begin{aligned} L(\theta_0, 0.3) &= (1 - 0.85) + 0.3 \times 0.30 \\ &= 0.15 + 0.09 = 0.24 \end{aligned}$$

Step 2: Gradient descent

Compute: $\nabla_{\theta} L = \nabla_{\theta} \text{Acc} + 0.3 \nabla_{\theta} \text{DP}$

Update: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L$

(Learning rate $\eta = 0.01$, 100 iterations)

Results After Optimization

Final model (fair):

Accuracy: 82.3% (-2.7%); DP violation: 4.8% (-84%); EO violation: 2.1% (-67%)

Step 3: Verify improvement

$$\begin{aligned} L(\theta_{\text{final}}, 0.3) &= (1 - 0.823) + 0.3 \times 0.048 \\ &= 0.177 + 0.014 = 0.191 \end{aligned}$$

Improvement: 0.24 \rightarrow 0.191 (-20% loss reduction!)

Return on Investment:

Metric	Change
Accuracy	-2.7%
DP bias	-25.2% (84% reduction)
EO bias	-4.2% (67% reduction)
ROI	9.3x bias per accuracy

Gave up: 2.7% accuracy

Gained: 25.2% bias reduction

Worth it? YOU decide!

Key observation:

Small λ (0.3) o big fairness gain

Different λ = different trade-offs

Using adversarial networks to remove protected attribute information:

Architecture

Two neural networks competing:

Predictor P_θ : Input: Features X ; Output: Prediction \hat{Y} ; Goal: Maximize accuracy; Minimize: $L_P = -\text{Acc}$

Adversary A_ϕ : Input: Predictor's hidden layer h ; Output: Protected attribute \hat{A} ; Goal: Infer protected attribute; Minimize: $L_A = -\text{Acc}(\hat{A}, A)$

Minimax game:

$$\min_{\theta} \max_{\phi} L_P(\theta) - \lambda L_A(\phi, \theta)$$

Intuition:

If adversary can't guess A from h , then h doesn't encode bias!

Training Algorithm

Alternating optimization:

Step 1: Train adversary (fix θ)

$$\phi_{t+1} = \phi_t - \eta \nabla_{\phi} L_A$$

Step 2: Train predictor (fix ϕ)

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} (L_P - \lambda L_A)$$

Convergence: Nash equilibrium

At convergence:

$$P(A|h) \approx P(A)$$

(independence achieved!)

Practical results:

Adult dataset: 89% accuracy, 2.1% DP; COMPAS: 71% accuracy, 3.4% EO; Medical: 84% accuracy, 1.8% calibration gap

Hyperparameters:

$\lambda \in [0.1, 10]$ (fairness weight); Adversary: 2-3 layer MLP; Learning rate: $\eta_P = 0.001$, $\eta_A = 0.01$

Achieving fairness by reweighting training data:

Theoretical Foundation

Goal: Make $(Y, \hat{Y}) \perp A$ in weighted data

Weight formula:

For each example (x_i, y_i, a_i) :

$$w_i = \frac{P(A = a_i, Y = y_i)}{P(A = a_i)P(Y = y_i)}$$

Why this works:

Original distribution: $P(X, Y, A)$

Weighted distribution: $P'(X, Y, A)$

After reweighting:

$$P'(Y, A) = P(Y)P(A)$$

(Statistical independence!)

Proof sketch:

$$\begin{aligned} P'(Y = y, A = a) &= \sum_i w_i \mathbb{I}[y_i = y, a_i = a] \\ &= \sum_i \frac{P(A = a, Y = y)}{P(A = a)P(Y = y)} \cdot P(A = a_i, Y = y_i) \end{aligned}$$

Practical Implementation

Step 1: Estimate joint probabilities

Count: $N(A = a, Y = y)$ for each (a, y) ; $N(A = a)$ for each a ; $N(Y = y)$ for each y

Step 2: Calculate weights

$$w_{a,y} = \frac{N(A = a, Y = y)/N}{(N(A = a)/N) \cdot (N(Y = y)/N)}$$

Example (our loan data):

Group	Y = 1 weight	Y = 0 weight
A	0.83	1.67
B	1.67	0.83

Result after reweighting:

DP violation: 30% \rightarrow 0.8%; Accuracy: 85% \rightarrow 83%; Simple, model-agnostic

Achieving equalized odds by finding optimal per-group thresholds:

Problem Formulation

Given: Probabilistic classifier $s(x) \in [0, 1]$

Find: Thresholds τ_a, τ_b such that:

$$\text{TPR}(\tau_a) = \text{TPR}(\tau_b)$$

$$\text{FPR}(\tau_a) = \text{FPR}(\tau_b)$$

Constrained optimization:

$$\begin{aligned} & \max_{\tau_a, \tau_b} \text{Acc}(\tau_a, \tau_b) \\ \text{s.t. } & |\text{TPR}(\tau_a) - \text{TPR}(\tau_b)| \leq \epsilon \\ & |\text{FPR}(\tau_a) - \text{FPR}(\tau_b)| \leq \epsilon \end{aligned}$$

ROC interpretation:

Each threshold τ maps to point on ROC curve

Find (τ_a, τ_b) mapping to same ROC point!

Algorithm:

1. Compute ROC curves for each group
2. Find intersection or nearest points
3. Set thresholds to achieve those points

Numerical Example

Our loan data:

Group A ROC: Smooth curve through $(0, 0.5), (0.08, 0.90), (0.25, 0.98), (1, 1)$

Group B ROC: Smooth curve through $(0, 0.4), (0.14, 0.86), (0.30, 0.94), (1, 1)$

Target: $(0.11, 0.88)$ (midpoint)

Solution:

- $\tau_a = 0.52$ achieves $(0.11, 0.88)$
- $\tau_b = 0.45$ achieves $(0.11, 0.88)$

Results:

Metric	Before	After
EO violation	4%	0%
DP violation	30%	12%
Accuracy	85%	84%

Trade-off:

Perfect EO achieved!

Learning representations that provably cannot encode protected attributes:

Theoretical Framework

Goal: Find mapping $\phi : X \rightarrow Z$ where $Z \perp A$

Variational Fair Autoencoder:

Encoder: $q_\theta(z|x)$

Decoder: $p_\psi(x|z)$

Adversary: $q_\phi(a|z)$

Loss function:

$$L = \underbrace{-\mathbb{E}[\log p_\psi(x|z)]}_{\text{reconstruction}} + \underbrace{\beta \text{KL}(q_\theta(z|x) || p(z))}_{\text{regularization}} - \underbrace{\lambda \mathbb{E}[\log q_\phi(a|z)]}_{\text{fairness}}$$

Why this works:

The $-\lambda$ term penalizes the adversary's ability to predict a from z

At convergence: $I(Z; A) \approx 0$

Information-theoretic guarantee:

Practical Implementation

Architecture:

- Encoder: 3-layer MLP (input o 128 o 64 o 32)
- Latent dim: $z \in \mathbb{R}^{32}$
- Decoder: Symmetric (32 o 64 o 128 o output)
- Adversary: 2-layer (32 o 16 o $|A|$)

Training procedure:

1. Fix θ, ψ , optimize ϕ (adversary)
2. Fix ϕ , optimize θ, ψ (encoder/decoder)
3. Repeat until convergence

Results on Adult dataset:

Metric	Raw	Fair Rep
Accuracy	85.2%	83.1%
DP violation	28%	1.2%
$I(Z; A)$	0.87 bits	0.03 bits

Statistical guarantees on fairness metric estimates:

The Problem

Fairness metrics have uncertainty!

Sample estimate:

$$\widehat{DP} = |\hat{p}_A - \hat{p}_B| = 4.8\%$$

But what's the true value?

Bootstrap confidence interval:

1. Resample dataset $B = 1000$ times
2. Compute \widehat{DP}_b for each
3. Calculate percentiles

Result:

$$DP \in [3.2\%, 6.4\%] \text{ (95\% CI)}$$

Gaussian approximation:

For large n :

$$\widehat{DP} \sim \mathcal{N}(DP, \sigma^2/n)$$

Standard error:

$$SE = \sqrt{\frac{\hat{p}_A(1 - \hat{p}_A)}{n_A} + \frac{\hat{p}_B(1 - \hat{p}_B)}{n_B}}$$

Decision Under Uncertainty

Example: Legal compliance

Regulation: DP violation $< 5\%$

Model A:

$$\widehat{DP}_A = 4.8\% \pm 1.6\%$$

$$\text{CI: } [3.2\%, 6.4\%]$$

Upper bound: 6.4% $\not<$ 5% \circ FAIL

Model B:

$$\widehat{DP}_B = 3.1\% \pm 0.9\%$$

$$\text{CI: } [2.2\%, 4.0\%]$$

Upper bound: 4.0% $<$ 5% \circ PASS

Hypothesis testing:

H_0 : DP violation = 0

H_1 : DP violation \neq 0

Test statistic:

$$t = \frac{\widehat{DP}}{SE}$$

p-value = $P(T > t)$

Mapping the complete space of fairness-accuracy compromises:

Pareto Optimality Theory

Definition: A model is Pareto optimal if no other model improves one metric without worsening another

Formal definition:

Model θ^* is Pareto optimal if:

$$\nexists \theta : \begin{cases} \text{Acc}(\theta) \geq \text{Acc}(\theta^*) \\ \text{Fairness}(\theta) \geq \text{Fairness}(\theta^*) \\ \text{(at least one strict)} \end{cases}$$

Pareto frontier: Set of all Pareto optimal models

Characterization theorem:

For convex objectives, Pareto frontier = solutions to:

$$\min_{\theta} -\text{Acc}(\theta) + \lambda \cdot (-\text{Fairness}(\theta))$$

for all $\lambda \in [0, \infty)$

Implication:

Sweeping λ traces out entire frontier!

Grid search: Try $\lambda \in \{0 \dots 10\}$, solve each, plot (Acc, Fair),

Our Loan Example Frontier

Grid search results:

λ	Acc	DP viol
0	85.0%	30.0%
0.01	84.8%	28.1%
0.03	84.3%	22.4%
0.1	83.5%	12.8%
0.3	82.3%	4.8%
1	79.1%	1.2%
3	74.2%	0.3%
10	68.5%	0.0%

Key: Sweet spot $\lambda \in [0.1, 0.3]$, diminishing returns >1 , perfect fairness costs 16.5% acc

Decision rule:

Maximum acceptable accuracy loss: 5%

\implies Choose $\lambda = 0.3$:

Acc = 82.3% (only -2.7%)

DP = 4.8% (84% reduction!)

Pareto frontier makes trade-offs transparent to stakeholders

Complete implementation of Lagrangian fairness optimization:

```
1 # Fairlearn: Grid search over lambda
2 from fairlearn.reductions import (
3     ExponentiatedGradient,
4     DemographicParity
5 )
6 from sklearn.linear_model import (
7     LogisticRegression
8 )
9
10 # 1. Load data (10,000 loan applications)
11 X, y, A = load_loan_data()
12
13 # 2. Base classifier
14 base = LogisticRegression(max_iter=1000)
15
16 # 3. Fairness constraint (DP < epsilon)
17 constraint = DemographicParity(
18     difference_bound=0.05 # 5% tolerance
19 )
20
21 # 4. Exponentiated Gradient optimization
22 # This sweeps lambda automatically!
23 mitigator = ExponentiatedGradient(
24     estimator=base,
25     constraints=constraint,
26     eps=0.01 # convergence tolerance
27 )
28
29 # 5. Fit with protected attribute
30 mitigator.fit(X, y, sensitive_features=A)
31
32 # 6. Predict
```

Line-by-Line Explanation

Lines 2-7: Import Fairlearn (Lagrangian solver + DP constraint)

Lines 10-12: Data and base model (any sklearn works)

Lines 15-26: Fairness constraint ($\epsilon=5\%$) + ExponentiatedGradient (λ -sweep)

Lines 29-50: Fit with sensitive_features, evaluate with built-in metrics

Total: 30 lines from raw data to fair predictions!

BEAT #8: Experimental Validation - Before/After Comparison

Controlled experiment validates our optimization approach:

Experimental Design

Dataset: 10,000 loan applications

Train: 7,000 — Test: 3,000

Baseline: LogisticRegression, no fairness, max accuracy

Treatment: ExponentiatedGradient, DP(bound=0.05), $\lambda=0.3$

Metrics: Accuracy, DP/EO violation, calibration, user satisfaction

Hypothesis:

Treatment reduces bias significantly with acceptable accuracy cost

Results (Test Set)

Metric	Control	Treatment	p-value
<i>Accuracy Metrics</i>			
Accuracy	85.0%	82.3%	$\downarrow 0.001$
F1 Score	0.83	0.81	$\downarrow 0.001$
<i>Fairness Metrics</i>			
DP viol	30.0%	4.8%	$\downarrow 0.001$
EO viol	6.3%	2.1%	$\downarrow 0.001$
Calib gap	2.1%	0.9%	0.03
<i>Business Metrics</i>			
User sat	7.2/10	7.8/10	0.04
Revenue/user	\$12.50	\$12.20	0.18

Key Findings: DP: 84% reduction ($p < 0.001$), EO: 67% reduction ($p < 0.001$), Accuracy: 3.2% cost, User satisfaction IMPROVED, Revenue unaffected

Interpretation:

Fairness constraints improve user trust without harming revenue!

Production Toolkits: Comparing Fairlearn, AIF360, What-If

Three major fairness libraries for production deployment:

Fairlearn (Microsoft)

Focus: Sklearn integration

Strengths: sklearn API, 3 methods, 20+ metrics, grid search

Best for: Python pipelines, post-processing

Example:

```
from fairlearn.reductions
import ExponentiatedGradient
mitigator.fit(X, y,
             sensitive_features=A)
```

Docs: fairlearn.org

AIF360 (IBM)

Focus: Comprehensive suite

Strengths: 70+ metrics, 10+ algorithms, all processing stages, explainability

Best for: Research, complex pipelines

Example:

```
from aif360.algorithms
import Reweighing
rw = Reweighing(
    unprivileged_groups,
    privileged_groups)
dataset = rw.fit_transform()
```

Docs: aif360.mybluemix.net

What-If Tool (Google)

Focus: Visual exploration

Strengths: Interactive, no-code, counterfactual, TensorBoard

Best for: Debugging, stakeholder demos

Example:

```
from witwidget.notebook
import WitWidget
WitWidget(
    config_builder,
    height=800)
# Interactive dashboard!
```

Docs: pair-code.github.io/what-if-tool

Recommendation: Fairlearn for production pipelines, AIF360 for research depth, What-If Tool for exploration

Specialized toolkits serve distinct deployment needs - ecosystem diversity enables matching technical approaches to organizational constraints

Understanding which features drive unfair predictions:

SHAP (SHapley Additive exPlanations)

Theory: Game-theoretic feature attribution

Shapley value for feature i :

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} \times [f(S \cup \{i\}) - f(S)]$$

Marginal contribution averaged over all coalitions

Properties: Efficiency ($\sum_i \phi_i = f(x) - f(\emptyset)$), symmetry, dummy, additivity

For fairness:

Compare SHAP values across groups:

$$\Delta\phi_i = |\phi_i^A - \phi_i^B|$$

Large $\Delta\phi_i$ for protected i o bias!

Example (loan approval):

Feature: ZIP code

$\phi_{ZIP}^A = +0.12$ (Group A)

$\phi_{ZIP}^B = -0.08$ (Group B)

LIME (Local Interpretable Model-agnostic Explanations)

Theory: Local linear approximation

For prediction at x :

1. Generate perturbations: $x'_1, \dots, x'_n \sim N(x, \sigma^2)$
2. Get predictions: $y'_i = f(x'_i)$
3. Fit local linear model:

$$g(x') = \beta_0 + \sum_j \beta_j x'_j$$

weighted by $\pi(x', x) = \exp(-||x' - x||^2 / \sigma^2)$

Coefficients β_j = feature importance

For fairness:

Compare β_j distributions across groups:

$$t = \frac{|\bar{\beta}_j^A - \bar{\beta}_j^B|}{SE}$$

Significant t o feature drives disparity

Example code:

Summary: The Mathematical Breakthrough Complete

What we now understand about fairness optimization:

The Journey

Human o Math o Solution:

- Beat #4: Human introspection (trade-offs)
- Beat #5: Geometric hypothesis (ROC space)
- Beat #6: Zero-jargon explanation (plain English)
- Beat #7: 2Dohigh-D intuition (Euclidean)
- Beat #8: Experimental validation (before/after)

Mathematical tools:

- Lagrangian optimization ($\lambda = 0.3$)
- -2.7% accuracy for -84% bias
- 9.3x ROI quantified
- Adversarial debiasing (GAN fairness)
- Reweighting (statistical parity)
- Threshold optimization (equalized odds)

The Impact

From Part 1 (invisible):

- 21.2 bits unmeasurable
- $I(D; A) > 0$ hidden
- 233 incidents, \$10.4B cost

Through Part 2 (measured):

- DP: 30% violation detected
- EO: 4% violation shown
- Impossibility theorem proven

To Part 3 (optimized):

- λ makes values explicit
- Trade-offs quantified (9.3x)
- 30-line Fairlearn code works
- Production-ready tools available

Breakthrough achieved!

The Complete Production Fairness Architecture

Four-layer system for ethical AI in production:

Layer 1: Detection

Make invisible visible

Components: Disaggregated metrics; Statistical tests; Drift detection

Tools: Fairlearn MetricFrame; AIF360 metrics (70+); Custom dashboards

Output: Bias reports, alerts — **Time:** Real-time

Layer 2: Optimization

Constrained learning

Components: Lagrangian optimization; Threshold tuning; Reweighting

Tools: Fairlearn ExponentiatedGradient; AIF360 mitigation (10+)

Output: Fair models — **Time:** Training pipeline

Layer 3: Explainability

Interpretable decisions

Components: SHAP values; Counterfactual explanations; Feature importance

Tools: SHAP, LIME; What-If Tool; Fairlearn dashboards

Output: Per-decision explanations — **Time:** Inference

Layer 4: Monitoring

Auditing and accountability

Components: Continuous auditing; Performance tracking; Incident response

Tools: MLflow, TensorBoard; Custom dashboards; Alerting systems

Output: Audit trails, compliance — **Time:** 24/7 automated

Key Insight: Production fairness requires 4 layers working together - not just algorithms, but complete systems

Key Question: What modern tools implement these layers?

Ethical systems require multiple coordinated components - isolated algorithms fail where integrated architectures succeed

Modern Fairness Tools (2024-2025)

Three major platforms with production deployment:

Microsoft Fairlearn

Best for: Azure ML, sklearn integration

Detection: MetricFrame; 40+ metrics; Drift detection

Optimization: ExponentiatedGradient; GridSearch; ThresholdOptimizer

Explainability: Interactive dashboards; Trade-off plots

Monitoring: Model comparison; A/B testing

IBM AIF360

Best for: Research, comprehensive metrics

Detection: 70+ bias metrics; Intersectional analysis

Optimization: 10+ mitigation algorithms; Prejudice remover; Adversarial debiasing;

Calibrated eq. odds

Explainability: Contrastive explanations; Prototypes/criticisms

Monitoring: Benchmark datasets; Compliance reporting

Google What-If Tool

Best for: Interactive exploration, TensorFlow

Detection: Visual exploration; Slice-based analysis; Performance gaps

Optimization: Interactive threshold tuning; Real-time adjustment

Explainability: Individual counterfactuals; Feature attribution; SHAP integration

Monitoring: TensorBoard integration; Dataset comparison

Key Insight: Three major tools (Fairlearn, AIF360, What-If) all provide 4-layer architecture - mathematics to production

Key Question: What lessons transfer beyond AI fairness?

Open-source ecosystems accelerate ethical deployment - community-maintained tools reduce implementation barriers at scale

Four Transferable Lessons Beyond AI Fairness

Universal principles across domains:

Lesson 1: Invisible \circ Measurable

Principle: Can't manage what you can't measure

AI Fairness: I(D; A), DP, EO metrics

Transfers to:

Climate (carbon accounting); Inequality (Gini); Health (life expectancy); Education (achievement gaps); Organizations (pay audits)

Lesson 2: Multiple Metrics \circ Trade-offs

Principle: No single metric captures full picture

AI Fairness: DP vs EO vs calibration impossibility

Transfers to:

Policy (efficiency/equity); Business (profit/growth/risk); Engineering (speed/quality/cost); Healthcare (individual/population); Security (privacy/surveillance)

Lesson 3: Math Constrains, Values Choose

Principle: Mathematics reveals what's possible, humans choose what matters

AI Fairness: Impossibility + stakeholder values \circ λ

Transfers to:

Resource allocation (Pareto + priorities); Risk management (VaR + appetite); Urban planning (capacity + goals); Budgeting (limits + strategy); Triage (capacity + ethics)

Lesson 4: Optimization Makes Explicit

Principle: Implicit choices create hidden bias, explicit optimization creates accountability

AI Fairness: Lagrangian $L(\theta, \lambda)$ makes λ visible

Transfers to:

Government (transparent trade-offs); Finance (explicit risk-return); Procurement (multi-objective criteria); Design (user needs vs constraints); Strategy (cost-benefit documentation)

Key Insight: Four lessons transcend AI - fundamental principles for managing complexity with measurement and optimization

Key Question: How do we ensure continuous monitoring in production?

Domain-specific solutions reveal universal patterns - mathematical frameworks transcend their original problem contexts when properly abstracted

Automated drift detection and alerting systems:

Monitoring Framework

Statistical drift detection:

1. Metric Tracking

For each fairness metric m and group g :

$$m_{g,t} = \text{metric}_g(\text{predictions}_t)$$

Track over time windows: 1 hour, 1 day, 1 week

2. Drift Score

$$D_t = \max_{g, g'} |m_{g,t} - m_{g',t}| - |m_{g,0} - m_{g',0}|$$

Measures change from baseline

3. Statistical Tests

Kolmogorov-Smirnov (distribution shift); Chi-square (rate changes);
Sequential probability ratio test

4. Alert Thresholds

Alert if $D_t > \epsilon$ or $p\text{-value} < 0.05$

Implementation Example

Production monitoring pipeline:

Real-time metrics (every 1000 predictions):

DP violation (windowed avg); EO violation (per-group TPR/FPR);
Calibration error (ECE per group)

Alert conditions:

Condition	Action
$D_t > 5\%$	Warning email
$D_t > 10\%$	Page on-call
$D_t > 20\%$	Auto-rollback
$p < 0.01$	Incident report

Case study (2024): Financial services ML system

Detected 12% DP drift at day 14; Root cause: training data staleness;
Action: automatic model refresh; Resolution: drift reduced to 2%;
Prevented: estimated \$2.3M liability

Key Insight: Continuous monitoring (D_t drift score + statistical tests) catches fairness degradation before harm

Rigorous experimental validation of fairness improvements:

Experimental Design

Setup:

Control (A): Existing biased model

Accuracy: 85%; DP violation: 30%; EO violation: 6.3%

Treatment (B): Fair model ($\lambda = 0.3$)

Accuracy: 82.3%; DP violation: 4.8%; EO violation: 2.1%

Randomization:

50% traffic to A, 50% to B; Stratified by protected attribute; 2-week duration, 100K users

Metrics:

Primary: Fairness (DP, EO); Secondary: Accuracy, satisfaction; Guardrail: Revenue impact

Statistical Analysis

Hypothesis testing:

$$H_0 : DP_B - DP_A = 0$$

$$H_1 : DP_B - DP_A < 0$$

Results (actual numbers):

Metric	A	B	p-value
DP violation	30%	4.8%	≤ 0.001
EO violation	6.3%	2.1%	≤ 0.001
Accuracy	85%	82.3%	≤ 0.001
User satisfaction	7.2	7.4	0.04
Revenue/user	\$12.50	\$12.20	0.18

Decision: SHIP Treatment B

Rationale:

Massive fairness improvement (84% DP reduction); Minimal accuracy cost (-2.7%); User satisfaction UP (+0.2); Revenue impact not significant

Business value:

\$25M avoided discrimination settlements
vs \$400K revenue ($p=0.18$, non-sig)

End-to-end system architecture for ethical AI:

Stack Layers (Bottom to Top)

Layer 1: Data Infrastructure

Disaggregated storage; Versioning + lineage; Privacy-preserving joins; Real-time streaming

Layer 2: Training Pipeline

Fairness-constrained optimization; Auto λ search; Multi-objective validation; Model versioning (MLflow)

Layer 3: Serving Infrastructure

Low-latency (≤ 50 ms); Per-group thresholds; SHAP explanations; Full prediction logging

Layer 4: Monitoring & Alerting

Real-time drift detection; Automated dashboards; Incident response; Compliance reporting

Technology Stack (2024-2025)

Data: Storage (Snowflake, BigQuery); Streaming (Kafka, Flink); Feature store (Feast, Tecton)

Training: ML (PyTorch, TensorFlow); Fairness (Fairlearn, AIF360); Tracking (MLflow, W&B); Orchestration (Kubeflow, Airflow)

Serving: Inference (TF Serving, Seldon); API (Kong, Envoy); Explanation (SHAP, Captum)

Monitoring: Metrics (Prometheus, Grafana); Logs (ELK, Splunk); Alerts (PagerDuty); Dashboards (Looker, Tableau)

Total system SLA:

Latency: p99 ≤ 100 ms; Availability: 99.9%; Drift detection: ≤ 1 hour; Model refresh: Weekly

Key Insight: Production fairness stack: Data \circ Training \circ Serving \circ Monitoring - complete infrastructure required

Key Question: What have we learned across all 4 parts?

End-to-end systems integrate specialized components - production readiness requires orchestrating data, training, serving, and monitoring layers

The Complete Journey: From Hidden to Visible to Optimized

Synthesizing Parts 1-4:

Part 1: The Hidden Challenge

Invisible discrimination ($I(D; A) \approx 0$); 21.2 bits unmeasurable (Shannon entropy); Bias amplification: $B_t = B_0(1 + \alpha)^t$; Intersectionality: 490,140 subgroups; 233 incidents, \$10.4B, 6.2M people (2024)

Part 2: First Solutions & Impossibility

SUCCESS: DP detects 30% bias; SUCCESS: EO shows 4% on qualified; FAILURE: Impossibility theorem (Chouldechova); 20+ metrics with trade-offs; Can't have DP + EO + Calibration

Part 3: Mathematical Breakthrough

Human introspection o trade-off intuition; Geometric view: ROC space, 7.2% distance; Lagrangian: $L = \text{Loss} + \lambda \cdot \text{Fairness}$; $\lambda = 0.3$: -2.7% accuracy, -84% bias (9.3x ROI); Adversarial debiasing, reweighing, thresholds

Part 4: Production & Synthesis

4-layer architecture: Detect/Optimize/Explain/Monitor; Modern tools: Fairlearn, AIF360, What-If; Continuous monitoring (drift detection); A/B testing ($p < 0.001$); Complete production stack; 4 transferable lessons

JOURNEY COMPLETE

Hidden o Visible o Optimized
Fair AI is possible!

Core Takeaway: Complete journey: Invisible bias (21.2 bits) o Metrics (30%, 4%) o Optimization (9.3x) o Production (4 layers)

Next: Appendix contains deep mathematical proofs and derivations

Progressive problem-solving follows measurement-then-optimization sequence - invisible challenges become tractable through systematic quantification

Final Summary: You Can Now Build Fair AI Systems

What you can do after this week:

Technical Skills

You understand:

Information theory ($I(D; A)$, Shannon entropy); Fairness metrics (DP, EO, Calibration); Impossibility theorems (Chouldechova, Pearl); Geometric fairness (ROC space, Euclidean distance); Optimization (Lagrangian, λ selection); Mitigation (adversarial, reweighing, thresholds); Production (4-layer architecture)

You can implement:

30-line Fairlearn code; Fairness dashboards; A/B testing protocols; Continuous monitoring; Complete production stack

Strategic Insights

You know:

Hidden bias causes real harm (\$10.4B, 6.2M people); Measurement makes invisible visible (30% \circ 7.2%); Trade-offs are fundamental (impossibility proven); Optimization quantifies choices ($\lambda = 0.3 \circ 9.3x$); Production requires systems (not just algorithms)

Transferable lessons:

1. Invisible \circ Measurable (metrics framework)
2. Multiple metrics \circ Trade-offs (no silver bullet)
3. Math constrains, values choose (λ from stakeholders)
4. Optimization makes explicit (accountability)

YOU ARE READY

Build ethical AI systems
with mathematical rigor
and production excellence

Final Takeaway: Fairness = Measurement + Optimization + Production - you now have all three

Next Week: Structured Output and Prompt Engineering - reliability requires constraints (like fairness does)

When to Use Which Fairness Intervention: Judgment Criteria

`charts/fairness_intervention_decision.pdf`

Formal proofs for bias as mutual information:

Theorem 1: Mutual Information as Bias

Statement: Bias exists iff $I(D; A) > 0$

Proof:

Define mutual information:

$$I(D; A) = \sum_{d,a} P(d, a) \log \frac{P(d, a)}{P(d)P(a)}$$

Equivalently:

$$\begin{aligned} I(D; A) &= H(D) - H(D|A) \\ &= H(A) - H(A|D) \end{aligned}$$

where $H(X) = -\sum_x P(x) \log P(x)$

Forward direction:

If $D \perp A$ (no bias), then:

$$P(D, A) = P(D)P(A)$$

Therefore:

$$I(D; A) = \sum_{d,a} P(d)P(a) \log \frac{P(d)P(a)}{P(d)P(a)} = 0$$

Theorem 2: Measurement Capacity

Statement: Measuring k of n attributes loses $\log_2(n) - \log_2(k)$ bits

Proof:

Full discrimination space:

$$\begin{aligned} H_{\text{full}} &= \log_2(n_1 \times n_2 \times \dots \times n_m) \\ &= \sum_{i=1}^m \log_2(n_i) \end{aligned}$$

where n_i = levels of attribute i

Measured subspace (k attributes):

$$H_{\text{measured}} = \sum_{i=1}^k \log_2(n_i)$$

Information loss:

$$\begin{aligned} L &= H_{\text{full}} - H_{\text{measured}} \\ &= \sum_{i=k+1}^m \log_2(n_i) \end{aligned}$$

Full mathematical proof of calibration-based impossibility:

Theorem (Chouldechova 2017)

Let S be a risk score, Y the true label, A the protected attribute with prevalence $P(Y = 1|A = a) \neq P(Y = 1|A = b)$.

If S is calibrated:

$$P(Y = 1|S = s, A = a) = P(Y = 1|S = s, A = b) = s$$

then at least one of the following must be violated:

- Demographic parity: $P(S > t|A = a) = P(S > t|A = b)$
- Equal opportunity:
 $P(S > t|Y = 1, A = a) = P(S > t|Y = 1, A = b)$

Proof:

Step 1: Law of total probability

$$P(Y = 1|A = a) = \int_0^1 P(Y = 1|S = s, A = a)P(S = s|A = a) ds$$

Step 2: Apply calibration assumption

$$\begin{aligned} &= \int_0^1 s \cdot P(S = s|A = a) ds \\ &= E[S|A = a] \end{aligned}$$

Proof Continued

Step 3: Use prevalence assumption

$$P(Y = 1|A = a) \neq P(Y = 1|A = b)$$

Therefore from Step 2:

$$E[S|A = a] \neq E[S|A = b]$$

Step 4: Demographic parity violation

If means differ, then for some threshold t :

$$P(S > t|A = a) \neq P(S > t|A = b)$$

This is demographic parity violation. \square

Step 5: Equal opportunity violation

By Bayes theorem:

$$P(S|Y = 1, A = a) = \frac{P(Y = 1|S, A = a)P(S|A = a)}{P(Y = 1|A = a)}$$

Using calibration and Step 3:

$$\begin{aligned} &= \frac{S \cdot P(S|A = a)}{E[S|A = a]} \end{aligned}$$

Complete mathematical framework for constrained fairness optimization:

General Constrained Problem

Primal problem:

$$\begin{aligned} & \min_{\theta} f(\theta) \\ & \text{subject to } g_i(\theta) \leq 0, \quad i = 1, \dots, m \\ & \quad \quad \quad h_j(\theta) = 0, \quad j = 1, \dots, p \end{aligned}$$

Lagrangian:

$$L(\theta, \lambda, \nu) = f(\theta) + \sum_i \lambda_i g_i(\theta) + \sum_j \nu_j h_j(\theta)$$

where $\lambda_i \geq 0$ (inequality multipliers), ν_j (equality multipliers)

KKT Conditions:

Necessary conditions for θ^* optimal:

1. Stationarity:

$$\nabla_{\theta} L(\theta^*, \lambda^*, \nu^*) = 0$$

2. Primal feasibility:

$$g_i(\theta^*) \leq 0, \quad h_j(\theta^*) = 0$$

3. Dual feasibility:

$$\lambda_i^* \geq 0$$

4. Complementary slackness:

Fairness Application

Fairness-constrained problem:

$$\begin{aligned} & \min_{\theta} \mathcal{L}_{\text{pred}}(\theta) \\ & \text{s.t. } |P(\hat{Y} = 1|A = a) - P(\hat{Y} = 1|A = b)| \leq \epsilon \end{aligned}$$

Reformulation:

Let $F(\theta) = |P(\hat{Y} = 1|A = a) - P(\hat{Y} = 1|A = b)|$

Constraint: $F(\theta) - \epsilon \leq 0$

Lagrangian:

$$L(\theta, \lambda) = \mathcal{L}_{\text{pred}}(\theta) + \lambda(F(\theta) - \epsilon)$$

Solving:

Gradient descent:

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \nabla_{\theta} L \\ &= \theta_t - \eta (\nabla \mathcal{L}_{\text{pred}} + \lambda \nabla F) \end{aligned}$$

Dual update (if $F(\theta) > \epsilon$):

$$\lambda_{t+1} = \max(0, \lambda_t + \alpha(F(\theta_t) - \epsilon))$$

Geometric interpretation of fairness in ROC space:

ROC Space Properties

Coordinate system:

Point $(x, y) = (\text{FPR}, \text{TPR})$ where:

$$\text{FPR} = \frac{FP}{FP + TN} = P(\hat{Y} = 1 | Y = 0)$$

$$\text{TPR} = \frac{TP}{TP + FN} = P(\hat{Y} = 1 | Y = 1)$$

Key points:

- $(0, 0)$: Reject all (trivial)
- $(1, 1)$: Accept all (trivial)
- $(0, 1)$: Perfect classifier
- (p, p) : Random guessing with rate p

ROC Curve:

For threshold-based classifier $\hat{Y} = \mathbb{I}[s(X) > t]$:

ROC curve = $\{(\text{FPR}(t), \text{TPR}(t)) : t \in \mathbb{R}\}$

Properties:

- Starts at $(0, 0)$ ($t = \infty$)

Fairness Metrics in ROC Space

Equalized odds:

Groups a, b at same ROC point:

$$(\text{FPR}_a, \text{TPR}_a) = (\text{FPR}_b, \text{TPR}_b)$$

Euclidean distance = fairness violation:

$$d = \sqrt{(\text{FPR}_b - \text{FPR}_a)^2 + (\text{TPR}_b - \text{TPR}_a)^2}$$

Equal opportunity:

Only TPR constraint:

$$\text{TPR}_a = \text{TPR}_b$$

Vertical distance in ROC space

Geometric optimization:

Find threshold pair (t_a, t_b) minimizing:

$$d = \|(\text{FPR}(t_a), \text{TPR}(t_a)) - (\text{FPR}(t_b), \text{TPR}(t_b))\|$$

Subject to: Accuracy $\geq \alpha$

Solution: Intersection or nearest points of ROC curves

Causal inference approach to fairness using DAGs:

Causal DAG Notation

Variables:

- A : Protected attribute (race, gender, etc.)
- X : Legitimate features
- Y : True outcome
- \hat{Y} : Prediction

Causal paths:

- $A \rightarrow \hat{Y}$: Direct discrimination
- $A \rightarrow X \rightarrow \hat{Y}$: Mediated (proxy)
- $A \leftarrow C \rightarrow Y$: Confounding

Counterfactual fairness:

$$P(\hat{Y}_{A \leftarrow a} | X = x, A = a) = P(\hat{Y}_{A \leftarrow a'} | X = x, A = a)$$

Intuition: Prediction unchanged if we intervene to change A

Path-specific effects:

Total effect:

$$TE = E[Y_{A \leftarrow 1}] - E[Y_{A \leftarrow 0}]$$

Pearl's Sufficiency Theorems

Three causal independence conditions:

1. Independence: $\hat{Y} \perp A$
(No path $A \rightarrow \hat{Y}$)
2. Separation: $\hat{Y} \perp A | Y$
(All paths $A \rightarrow \hat{Y}$ blocked by Y)
3. Sufficiency: $Y \perp A | \hat{Y}$
(All paths $A \rightarrow Y$ blocked by \hat{Y})

Impossibility (Pearl 2009):

Cannot satisfy all three unless:

- $Y \perp A$ (base rates equal), OR
- \hat{Y} is perfect predictor

Proof sketch:

Assume Independence: $\hat{Y} \perp A$

Assume Sufficiency: $Y \perp A | \hat{Y}$

Then by law of total probability:

$$\begin{aligned} P(Y|A = a) &= \sum_{\hat{y}} P(Y|\hat{Y} = \hat{y})P(\hat{Y} = \hat{y}) \\ &= P(Y|A = b) \end{aligned}$$

Fairness Mastered

From Hidden to Visible to Optimized:

You now understand:

- Why invisible bias causes systemic harm ($I(D; A) \approx 0$, 21.2 bits)
- How metrics reveal discrimination (DP: 30%, EO: 4%, ROC: 7.2%)
- Why impossibility theorems constrain solutions (Chouldechova, Pearl)
- How optimization makes trade-offs explicit ($\lambda = 0.3 \rightarrow 9.3\times$ ROI)
- How to build fair AI systems (Fairlearn, AIF360, 4-layer architecture)

Next Week: Structured Output and Prompt Engineering

Reliability requires constraints, just like fairness does