

## Lesson 2.4: Algorithmic Fairness and Bias – Quiz

Module 2: The Access Problem

Prof. Dr. Joerg Osterrieder

## Question 1

**A credit scoring model does not use race as an input feature. However, it uses zip code, which is highly correlated with race due to historical residential segregation. What type of discrimination does this represent?**

- A. Disparate treatment – the model explicitly uses a protected attribute
- B. Disparate treatment through proxy – which is the same as using race directly
- C. Disparate impact – a facially neutral feature produces disproportionate outcomes
- D. Neither – since race is not used, there is no discrimination

## Question 1

**A credit scoring model does not use race as an input feature. However, it uses zip code, which is highly correlated with race due to historical residential segregation. What type of discrimination does this represent?**

- A. Disparate treatment – the model explicitly uses a protected attribute
- B. Disparate treatment through proxy – which is the same as using race directly
- C. Disparate impact – a facially neutral feature produces disproportionate outcomes
- D. Neither – since race is not used, there is no discrimination

*[Answer hidden – compile with \solutionstrue to reveal]*

Zip code is facially neutral (not a protected attribute), but produces disproportionate outcomes because it correlates with race. This is classic disparate impact. Disparate treatment requires explicit use of the protected attribute.

## Question 2

**A lending model approves 60% of white applicants and 42% of Black applicants. What is the disparate impact ratio, and does it satisfy the four-fifths rule?**

- A.  $DI = 1.43$ , satisfies the four-fifths rule
- B.  $DI = 0.42$ , does not satisfy the four-fifths rule
- C.  $DI = 0.70$ , satisfies the four-fifths rule (above 0.60)
- D.  $DI = 0.70$ , does not satisfy the four-fifths rule (below 0.80)

## Question 2

**A lending model approves 60% of white applicants and 42% of Black applicants. What is the disparate impact ratio, and does it satisfy the four-fifths rule?**

- A.  $DI = 1.43$ , satisfies the four-fifths rule
- B.  $DI = 0.42$ , does not satisfy the four-fifths rule
- C.  $DI = 0.70$ , satisfies the four-fifths rule (above 0.60)
- D.  $DI = 0.70$ , does not satisfy the four-fifths rule (below 0.80)

*[Answer hidden – compile with \solutionstrue to reveal]*

$DI = 42\% / 60\% = 0.70$ . The four-fifths rule requires  $DI \geq 0.80$ . Since  $0.70 < 0.80$ , there is evidence of disparate impact.

## Question 3

**Which fairness metric requires that the approval rate be equal across all demographic groups, regardless of their qualifications?**

- A. Calibration
- B. Demographic parity
- C. Equalized odds
- D. Counterfactual fairness

## Question 3

**Which fairness metric requires that the approval rate be equal across all demographic groups, regardless of their qualifications?**

- A. Calibration
- B. Demographic parity
- C. Equalized odds
- D. Counterfactual fairness

*[Answer hidden – compile with \solutionstrue to reveal]*

Demographic parity requires  $P(\hat{Y} = 1|A = 0) = P(\hat{Y} = 1|A = 1)$ —equal positive prediction rates regardless of qualifications. Equalized odds conditions on the true label; calibration requires equal meaning of scores.

## Question 4

**The Chouldechova impossibility theorem states that it is impossible to simultaneously satisfy demographic parity, equalized odds, and calibration, except when:**

- A. The model is perfectly accurate ( $AUC = 1.0$ )
- B. Base rates (e.g., default rates) are equal across groups
- C. Post-processing fairness corrections are applied
- D. The model uses no proxy variables

## Question 4

**The Chouldechova impossibility theorem states that it is impossible to simultaneously satisfy demographic parity, equalized odds, and calibration, except when:**

- A. The model is perfectly accurate ( $AUC = 1.0$ )
- B. Base rates (e.g., default rates) are equal across groups
- C. Post-processing fairness corrections are applied
- D. The model uses no proxy variables

*[Answer hidden – compile with \solutionstrue to reveal]*

The impossibility theorem holds when base rates differ between groups. If Group A has a 5% default rate and Group B has a 12% default rate, no model can simultaneously achieve all three metrics. Only when base rates are equal is simultaneous satisfaction possible.

## Question 5

**A SHAP beeswarm plot shows that “zip code” is the third most important feature in a credit scoring model. The SHAP values for zip codes in predominantly minority neighborhoods are consistently negative (pushing toward denial). What should the fairness auditor conclude?**

- A. Zip code may be acting as a proxy for race, and further disparate impact analysis is needed
- B. The SHAP values are meaningless because SHAP cannot detect proxy variables
- C. The model is working correctly – zip code reflects legitimate risk differences
- D. Remove zip code and retrain – this conclusively proves racial bias

## Question 5

**A SHAP beeswarm plot shows that “zip code” is the third most important feature in a credit scoring model. The SHAP values for zip codes in predominantly minority neighborhoods are consistently negative (pushing toward denial). What should the fairness auditor conclude?**

- A. Zip code may be acting as a proxy for race, and further disparate impact analysis is needed
- B. The SHAP values are meaningless because SHAP cannot detect proxy variables
- C. The model is working correctly – zip code reflects legitimate risk differences
- D. Remove zip code and retrain – this conclusively proves racial bias

*[Answer hidden – compile with \solutionstrue to reveal]*

SHAP identifies *what* drives predictions, not *why*. Zip code may reflect legitimate risk (local economic conditions) or may proxy for race (redlining legacy). Further analysis (e.g., testing disparate impact with and without zip code) is needed before concluding bias.

## Question 6

**What is the key difference between SHAP and LIME for model explainability?**

- A. SHAP provides global explanations; LIME provides only local explanations
- B. LIME is more theoretically grounded than SHAP
- C. SHAP works only with tree-based models; LIME works with any model
- D. SHAP is based on cooperative game theory and provides both global and local explanations; LIME approximates the model locally using perturbations

**What is the key difference between SHAP and LIME for model explainability?**

- A. SHAP provides global explanations; LIME provides only local explanations
- B. LIME is more theoretically grounded than SHAP
- C. SHAP works only with tree-based models; LIME works with any model
- D. SHAP is based on cooperative game theory and provides both global and local explanations; LIME approximates the model locally using perturbations

*[Answer hidden – compile with \solutionstrue to reveal]*

SHAP uses Shapley values (cooperative game theory) and provides both global and local explanations. LIME perturbs inputs locally and fits a simple surrogate model. Both are model-agnostic, but SHAP has stronger theoretical guarantees.

## Question 7

**In the Apple Card case (2019), Goldman Sachs stated that gender was not used as a model input. Yet women consistently received lower credit limits. This failure illustrates which concept?**

- A. Fairness through unawareness – removing the protected attribute did not remove the bias
- B. Equalized odds – the model had equal error rates
- C. Disparate treatment – gender was explicitly used
- D. Calibration failure – the model's probabilities were wrong

## Question 7

**In the Apple Card case (2019), Goldman Sachs stated that gender was not used as a model input. Yet women consistently received lower credit limits. This failure illustrates which concept?**

- A. Fairness through unawareness – removing the protected attribute did not remove the bias
- B. Equalized odds – the model had equal error rates
- C. Disparate treatment – gender was explicitly used
- D. Calibration failure – the model's probabilities were wrong

*[Answer hidden – compile with \solutionstrue to reveal]*

“Fairness through unawareness” is the strategy of simply excluding protected attributes. Proxies (spending patterns, income sources) can reintroduce the attribute. *Important caveat:* NYDFS’s March 2021 investigation (*Report on Apple Card Investigation*) concluded Goldman’s underwriting did *not* violate New York fair-lending laws; the 2019 outcry reflected algorithmic opacity and customer-service failures more than demonstrated illegal discrimination. The question is about the pedagogical concept (proxy risk), not a finding of fact.

## Question 8

A model has the following performance across two groups:

Metric	Group A	Group B
Accuracy	92%	78%
False Positive Rate	5%	18%
False Negative Rate	3%	4%

Which fairness criterion does this model most clearly violate?

- A. Equalized odds – the false positive rate differs dramatically (5% vs. 18%)
- B. The model is fair because false negative rates are similar
- C. Calibration – predicted probabilities differ in meaning
- D. Demographic parity – the approval rates are unequal

## Question 8

A model has the following performance across two groups:

Metric	Group A	Group B
Accuracy	92%	78%
False Positive Rate	5%	18%
False Negative Rate	3%	4%

Which fairness criterion does this model most clearly violate?

- A. Equalized odds – the false positive rate differs dramatically (5% vs. 18%)
- B. The model is fair because false negative rates are similar
- C. Calibration – predicted probabilities differ in meaning
- D. Demographic parity – the approval rates are unequal

*[Answer hidden – compile with \solutionstrue to reveal]*

Equalized odds requires equal TPR and FPR across groups. The FPR difference (5% vs. 18%) is dramatic – Group B members are 3.6x more likely to be falsely flagged. This is exactly the pattern ProPublica found in COMPAS.

## Question 9

**Which bias mitigation approach modifies the model's decision thresholds per group after training to equalize outcomes?**

- A. Pre-processing: disparate impact remover
- B. Pre-processing: reweighting
- C. Post-processing: threshold adjustment
- D. In-processing: adversarial debiasing

## Question 9

**Which bias mitigation approach modifies the model's decision thresholds per group after training to equalize outcomes?**

- A. Pre-processing: disparate impact remover
- B. Pre-processing: reweighting
- C. Post-processing: threshold adjustment
- D. In-processing: adversarial debiasing

*[Answer hidden – compile with \solutionstrue to reveal]*

Post-processing threshold adjustment sets different classification thresholds per group (e.g., approve if score  $> 0.55$  for Group A,  $> 0.48$  for Group B) to equalize outcomes. It is model-agnostic but may raise legal concerns about explicit differential treatment.

## Question 10

**Counterfactual fairness asks: “Would the decision have been different if the individual belonged to a different group?” What does this approach fundamentally require?**

- A. A causal model specifying how protected attributes influence other features
- B. A perfectly accurate model with  $AUC = 1.0$
- C. A large test dataset with balanced group representation
- D. Protected attributes to be included as direct model inputs

## Question 10

**Counterfactual fairness asks: “Would the decision have been different if the individual belonged to a different group?” What does this approach fundamentally require?**

- A. A causal model specifying how protected attributes influence other features
- B. A perfectly accurate model with  $AUC = 1.0$
- C. A large test dataset with balanced group representation
- D. Protected attributes to be included as direct model inputs

*[Answer hidden – compile with \solutionstrue to reveal]*

Counterfactual fairness requires a causal graph: flipping race means adjusting all causally downstream features (income, education, neighborhood). Without specifying these causal relationships, the counterfactual is undefined.

## Question 11

**Under the EU AI Act, credit scoring algorithms are classified as:**

- A. Limited risk – transparency obligations only
- B. High risk – requiring risk management, data governance, transparency, and human oversight
- C. Unacceptable risk – prohibited entirely
- D. Minimal risk – no requirements beyond general transparency

## Question 11

**Under the EU AI Act, credit scoring algorithms are classified as:**

- A. Limited risk – transparency obligations only
- B. High risk – requiring risk management, data governance, transparency, and human oversight
- C. Unacceptable risk – prohibited entirely
- D. Minimal risk – no requirements beyond general transparency

*[Answer hidden – compile with \solutionstrue to reveal]*

The EU AI Act classifies AI systems used for creditworthiness assessment and credit scoring as high-risk, requiring conformity assessments, technical documentation, data governance, transparency, human oversight, and accuracy/robustness standards.

## Question 12

**A bank's model development team reports that their new credit model has "no bias" because overall accuracy is 91%. What is the fundamental flaw in this claim?**

- A. 91% accuracy is too low for production deployment
- B. The team should have used AUC-ROC instead of accuracy
- C. Accuracy is not a valid metric for credit scoring
- D. Overall accuracy can mask significant disparities in error rates across demographic groups

## Question 12

**A bank's model development team reports that their new credit model has "no bias" because overall accuracy is 91%. What is the fundamental flaw in this claim?**

- A. 91% accuracy is too low for production deployment
- B. The team should have used AUC-ROC instead of accuracy
- C. Accuracy is not a valid metric for credit scoring
- D. Overall accuracy can mask significant disparities in error rates across demographic groups

*[Answer hidden – compile with \solutionstrue to reveal]*

A model with 91% overall accuracy could have 96% accuracy for Group A and 72% for Group B. Disaggregated metrics are essential. Fairness cannot be assessed from aggregate performance alone.

## Question 13

**In the COMPAS case, ProPublica showed that Black defendants had a higher false positive rate, while Northpointe showed the model was calibrated. Both were correct. How is this possible?**

- A. One of them made a statistical error
- B. The false positive rates were actually equal after adjusting for income
- C. COMPAS was only used for sentencing, not prediction
- D. They measured different aspects of fairness, and the impossibility theorem proves both cannot be satisfied simultaneously when base rates differ

## Question 13

**In the COMPAS case, ProPublica showed that Black defendants had a higher false positive rate, while Northpointe showed the model was calibrated. Both were correct. How is this possible?**

- A. One of them made a statistical error
- B. The false positive rates were actually equal after adjusting for income
- C. COMPAS was only used for sentencing, not prediction
- D. They measured different aspects of fairness, and the impossibility theorem proves both cannot be satisfied simultaneously when base rates differ

*[Answer hidden – compile with \solutionstrue to reveal]*

ProPublica measured equalized odds (equal error rates); Northpointe measured calibration (equal score meaning). The impossibility theorem proves both cannot hold when base rates differ – and recidivism base rates differed. This is the theorem's most famous real-world demonstration.

## Question 14

**Which of the following is a proxy variable for race in US credit scoring?**

- A. Annual income
- B. Employer name
- C. All of the above
- D. Zip code

## Question 14

**Which of the following is a proxy variable for race in US credit scoring?**

- A. Annual income
- B. Employer name
- C. All of the above
- D. Zip code

*[Answer hidden – compile with \solutionstrue to reveal]*

All three can act as proxies: zip code correlates with race due to residential segregation; income reflects historical wage gaps; employer name correlates with industry concentration by race. Any feature correlated with a protected attribute is a potential proxy.

## Question 15

**A credit model's fairness audit reveals that “years at current address” disproportionately penalizes minority applicants. Which bias source does this represent?**

- A. Label bias – the target variable is incorrectly defined
- B. Measurement bias – the feature is measured incorrectly for minority applicants
- C. Historical bias – systemic housing instability for minorities is encoded in this feature
- D. Aggregation bias – the model treats all groups identically when it should not

## Question 15

**A credit model's fairness audit reveals that “years at current address” disproportionately penalizes minority applicants. Which bias source does this represent?**

- A. Label bias – the target variable is incorrectly defined
- B. Measurement bias – the feature is measured incorrectly for minority applicants
- C. Historical bias – systemic housing instability for minorities is encoded in this feature
- D. Aggregation bias – the model treats all groups identically when it should not

*[Answer hidden – compile with \solutionstrue to reveal]*

Historical discrimination in housing (redlining, discriminatory lending, exclusion from homeownership) means minority individuals are more likely to be renters and to move more frequently. The feature encodes this historical inequality, creating systematic disadvantage.

**Adversarial debiasing is an in-processing bias mitigation technique. How does it work?**

- A. It adjusts decision thresholds per group after training
- B. It augments the training data with synthetic samples from underrepresented groups
- C. It removes all features correlated with protected attributes before training
- D. It trains a secondary model to predict group membership from the main model's output, then penalizes the main model when group membership is predictable

**Adversarial debiasing is an in-processing bias mitigation technique. How does it work?**

- A. It adjusts decision thresholds per group after training
- B. It augments the training data with synthetic samples from underrepresented groups
- C. It removes all features correlated with protected attributes before training
- D. It trains a secondary model to predict group membership from the main model's output, then penalizes the main model when group membership is predictable

*[Answer hidden – compile with \solutionstrue to reveal]*

Adversarial debiasing trains a predictor and an adversary simultaneously. The adversary tries to predict the protected attribute from the predictor's output. The predictor is penalized when the adversary succeeds, forcing it to learn representations that do not encode group membership.

## Question 17

**Under the US Equal Credit Opportunity Act (ECOA), what must a lender provide when denying a credit application?**

- A. Nothing – lenders are not required to explain denials
- B. Specific reasons for the denial (adverse action notice) that the applicant can understand and act upon
- C. A general statement that the applicant did not meet lending criteria
- D. The full model source code and training data

## Question 17

**Under the US Equal Credit Opportunity Act (ECOA), what must a lender provide when denying a credit application?**

- A. Nothing – lenders are not required to explain denials
- B. Specific reasons for the denial (adverse action notice) that the applicant can understand and act upon
- C. A general statement that the applicant did not meet lending criteria
- D. The full model source code and training data

*[Answer hidden – compile with \solutionstrue to reveal]*

ECOA requires specific, actionable adverse action reasons (e.g., “high debt-to-income ratio,” “insufficient credit history”). The CFPB has clarified (2022) that using complex ML models does not excuse vague explanations. SHAP/LIME can generate these reasons.

## Question 18

**A fairness audit finds that a lending model's overall disparate impact ratio is 0.83 (above the 0.80 threshold). However, for applicants aged 55+, the DI ratio is 0.62. What concept does this illustrate?**

- A. Simpson's paradox – aggregate fairness masks subgroup unfairness
- B. Calibration failure – the model is poorly calibrated for older applicants
- C. Intersectionality – unfairness appears at the intersection of protected attributes
- D. The impossibility theorem – you cannot satisfy multiple fairness criteria

## Question 18

**A fairness audit finds that a lending model's overall disparate impact ratio is 0.83 (above the 0.80 threshold). However, for applicants aged 55+, the DI ratio is 0.62. What concept does this illustrate?**

- A. Simpson's paradox – aggregate fairness masks subgroup unfairness
- B. Calibration failure – the model is poorly calibrated for older applicants
- C. Intersectionality – unfairness appears at the intersection of protected attributes
- D. The impossibility theorem – you cannot satisfy multiple fairness criteria

*[Answer hidden – compile with \solutionstrue to reveal]*

Intersectionality means that unfairness may be hidden when examining single protected attributes in isolation. The model may be “fair” for race overall, but unfair for older minority applicants specifically. Audits must test intersectional subgroups.

**Which statement about the fairness–accuracy tradeoff is most accurate?**

- A. Fairness constraints never affect accuracy – they only change the threshold
- B. Moderate fairness improvements often require small accuracy sacrifices; the tradeoff follows a Pareto frontier
- C. The tradeoff only exists for linear models, not for deep learning
- D. Enforcing fairness always destroys model accuracy – fairness and accuracy are incompatible

**Which statement about the fairness–accuracy tradeoff is most accurate?**

- A. Fairness constraints never affect accuracy – they only change the threshold
- B. Moderate fairness improvements often require small accuracy sacrifices; the tradeoff follows a Pareto frontier
- C. The tradeoff only exists for linear models, not for deep learning
- D. Enforcing fairness always destroys model accuracy – fairness and accuracy are incompatible

*[Answer hidden – compile with \solutionstrue to reveal]*

Research shows the fairness–accuracy Pareto frontier often has a “knee” where significant fairness gains come at modest accuracy cost. The tradeoff is real but often smaller than feared, and it applies to all model types.

## Question 20

**A financial institution is deploying a credit scoring model in both the EU and the US. Which combination of regulatory requirements must it satisfy?**

- A. ECOA and Fair Lending (US) plus EU AI Act (EU) – both jurisdictions' requirements apply to their respective markets
- B. Neither – credit scoring is exempt from AI regulation
- C. US regulations only – the EU AI Act has not yet been enforced
- D. EU AI Act only – it supersedes all national regulations

## Question 20

**A financial institution is deploying a credit scoring model in both the EU and the US. Which combination of regulatory requirements must it satisfy?**

- A. ECOA and Fair Lending (US) plus EU AI Act (EU) – both jurisdictions' requirements apply to their respective markets
- B. Neither – credit scoring is exempt from AI regulation
- C. US regulations only – the EU AI Act has not yet been enforced
- D. EU AI Act only – it supersedes all national regulations

*[Answer hidden – compile with \solutionstrue to reveal]*

Multinational institutions must satisfy each jurisdiction's requirements for operations in that jurisdiction. The US focuses on outcome-based testing (disparate impact), while the EU mandates process-based compliance (risk management system, documentation, human oversight).