



D-BIAS Analysis Report

cached_dataset.csv

Generated on 11/18/2025

Executive Summary

Fairness Score

31/100

Bias Risk

Critical

Fairness Label

Critical

Reliability

Moderate

Dataset Information

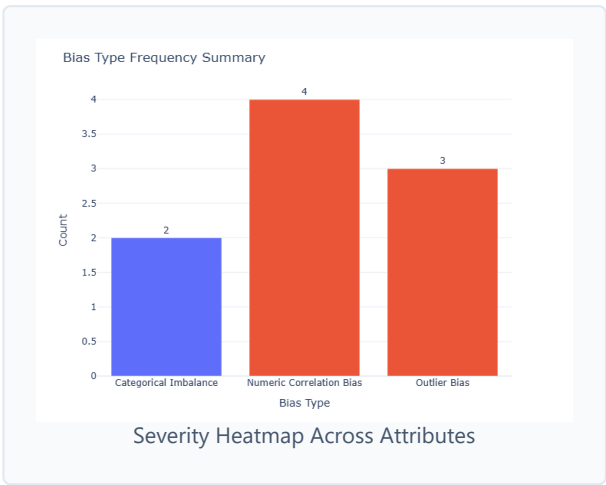
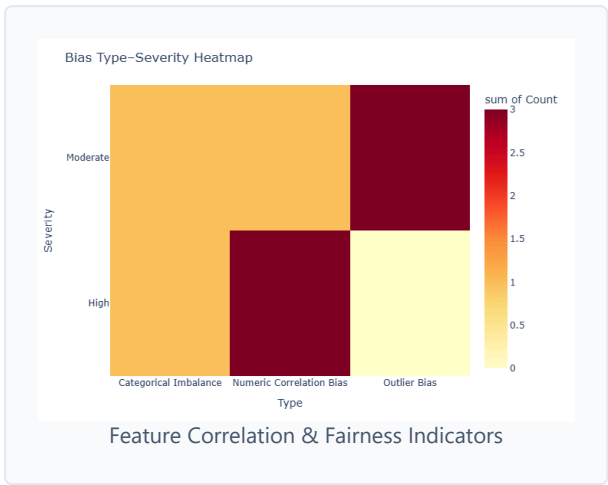
Rows: 660

Columns: 27

Mean: 31.328

Median: 0

Visualizations



Detected Biases

Categorical Imbalance

Moderate

Column: job_simp

Description: 'data scientist' dominates 67.7% of 'job_simp' values (entropy=1.61).

AI Explanation

Type: Categorical Imbalance; Feature: job_simp; 'data scientist' dominates 67.7% of 'job_simp' values (entropy=1.61).; Severity: Moderate

Definition: Categorical Imbalance

Categorical Imbalance

High

Column: seniority

Description: 'na' dominates 85.9% of 'seniority' values (entropy=0.60).

AI Explanation

Type: Categorical Imbalance; Feature: seniority; 'na' dominates 85.9% of 'seniority' values (entropy=0.60).; Severity: High

Definition: Categorical Imbalance

Numeric Correlation Bias

High

Column: min_salary ↔ max_salary

Description: Strong correlation $r=0.905$.

AI Explanation

Type: Numeric Correlation Bias; Feature: min_salary ↔ max_salary; Strong correlation $r=0.905$.; Severity: High

Definition: Numeric Correlation Bias

Numeric Correlation Bias

High

Column: min_salary ↔ avg_salary

Description: Strong correlation $r=0.966$.

AI Explanation

Type: Numeric Correlation Bias; Feature: min_salary ↔ avg_salary; Strong correlation $r=0.966$.; Severity: High

Definition: Numeric Correlation Bias

Numeric Correlation Bias

High

Column: max_salary ↔ avg_salary

Description: Strong correlation $r=0.984$.

AI Explanation

Type: Numeric Correlation Bias; Feature: max_salary ↔ avg_salary; Strong correlation $r=0.984$;; Severity: High

Definition: Numeric Correlation Bias

Numeric Correlation Bias

Moderate

Column: hadoop ↔ spark

Description: Strong correlation $r=0.532$.

AI Explanation

Type: Numeric Correlation Bias; Feature: hadoop ↔ spark; Strong correlation $r=0.532$;; Severity: Moderate

Definition: Numeric Correlation Bias

Outlier Bias

Moderate

Column: rating

Description: 7.6% of 'rating' values are outliers (left-skewed).

AI Explanation

Type: Outlier Bias; Feature: rating; 7.6% of 'rating' values are outliers (left-skewed).; Severity: Moderate

Definition: Outlier Bias

Outlier Bias

Moderate

Column: avg_salary

Description: 6.2% of 'avg_salary' values are outliers (right-skewed).

AI Explanation

Type: Outlier Bias; Feature: avg_salary; 6.2% of 'avg_salary' values are outliers (right-skewed).; Severity: Moderate

Definition: Outlier Bias

Outlier Bias

Moderate

Column: company_age

Description: 9.4% of 'company_age' values are outliers (right-skewed).

AI Explanation

Type: Outlier Bias; Feature: company_age; 9.4% of 'company_age' values are outliers (right-skewed).; Severity: Moderate

Definition: Outlier Bias

Recommendations

- **Resolve Multicollinearity (Highest Priority):** Immediately drop `min_salary` and `max_salary` from the feature set, retaining only `avg_salary` to represent compensation. This is essential for model stability.
- **Address Missing `seniority` Data:** Due to the extreme level of missingness, the `seniority` column should be removed from any modeling efforts unless the missing values can be reliably filled from an external source.
- **Counteract `job_simp` Imbalance:** Implement class weighting during model training to force the model to pay more attention to underrepresented job titles. This is crucial for improving fairness and utility across different roles.
- **Mitigate Outlier Effects:** Apply a log transformation to skewed numeric features like `avg_salary` and `company_age` to reduce the disproportionate influence of extreme values and create more stable models.
- **Strategic Data Augmentation:** To build a truly representative model, the dataset must be supplemented with more examples of non-'data scientist' roles and job listings that include seniority information.

Conclusion

: The dataset's overall "fairness health score" is **Needs Significant Improvement**. The combination of high-severity data quality issues, strong imbalances, and statistical redundancies means that using this data as-is would likely produce a biased, unreliable, and potentially unfair model. Foundational issues must be resolved before it can be trusted for any serious application. *