

Modelling Heteroscedasticity for Fair Regression using Polynomial Models

Mana Douma¹ and Emma Beauxis-Aussalet¹[0000–0002–4657–892X]

Vrije Universiteit Amsterdam, De Boelelaan 1111, 1081 HV Amsterdam, The
Netherlands

{m.a.douma, e.m.a.l.beauxisaussalet}@vu.nl

Abstract. Ensuring an even distribution of AI errors across social groups is a key aspect of fairness: AI errors should not be systematically larger or more frequent for specific populations. For AI systems predicting numeric values (i.e., regression problems), residuals are the primary error metrics. Even or uneven distributions of residuals is a problem of homo- or heteroscedasticity: a regression system is fair if its residuals are randomly and homogeneously distributed across protected features (homoscedasticity). Thus modelling heteroscedasticity is important to identify fairness issues. State-of-the-art methods model residuals’ heteroscedasticity using linear models. We demonstrate key limitations of such approach, and how these limitations can be addressed by using polynomial models, and signed residuals for the most complex cases. Polynomial models can address complex cases of heteroscedasticity that would remained undetected using linear models. However, interpreting their results is more complex, and future work is needed to assess the impact of outliers and overfitting.

Keywords: Fair AI · Regression Problem · Heteroscedasticity.

1 Introduction

Algorithmic bias in artificial intelligence (AI) has become a prominent issue, as cases discrimination have been arising in many applications [7, 11, 13]. In regression problems, a key form of bias is the heteroscedasticity of residuals: if residuals are not uniformly distributed over the output or input features, or over additional sensitive features not used as input, heteroscedasticity is present. It means that the AI results can be more inaccurate for specific data points, which creates fairness issues and a risk of discrimination (e.g., as AI errors are larger for a specific social group).

Current methods use simple linear methods to detect heteroscedasticity [2–6, 8]. Given AI results $\hat{\mathbf{y}}$, and groundtruth values \mathbf{y} , state-of-the-art methods model the squared or absolute residuals as:

$$(\hat{\mathbf{y}} - \mathbf{y})^2 = \beta_0 + \beta_1 \mathbf{x} \quad (1)$$

$$|\hat{\mathbf{y}} - \mathbf{y}| = \beta_0 + \beta_1 \mathbf{x} \quad (2)$$

where β_0 is the intercept, β_1 the coefficient (also called slope), and \mathbf{x} can be a sensitive feature (e.g., gender), the AI results $\hat{\mathbf{y}}$, the groundtruth value \mathbf{y} , or any of the input features of the AI system. The presence of significant non-zero slope is an indicator of heteroscedasticity in the residuals' distribution over feature \mathbf{x} .

However, this approach can overlook complex cases of heteroscedasticity (e.g., Fig. 1) resulting in compromised validity of the heteroscedasticity analyses [12]. We thus investigate the use of polynomial models compared to linear models. Polynomial models are particularly suited for capturing complex, non-linear relationships between variables. We thus argue that we better model squared or absolute residuals as:

$$(\hat{\mathbf{y}} - \mathbf{y})^2 = \beta_0 + \beta_1 \mathbf{x} + \beta_2 \mathbf{x}^2 + \dots + \beta_n \mathbf{x}^n \quad (3)$$

$$|\hat{\mathbf{y}} - \mathbf{y}| = \beta_0 + \beta_1 \mathbf{x} + \beta_2 \mathbf{x}^2 + \dots + \beta_n \mathbf{x}^n \quad (4)$$

where β_0 is the intercept, β_1 the coefficient of the linear term, β_2, \dots, β_n the parameters β of the polynomial terms. With polynomial models, any non-zero coefficient indicate heteroscedasticity. To assume homoscedastic residuals (e.g., a fair AI model), all parameters β must be zero. It provides a stricter set of criteria that are sensitive to more complex patterns of variability in the residuals.

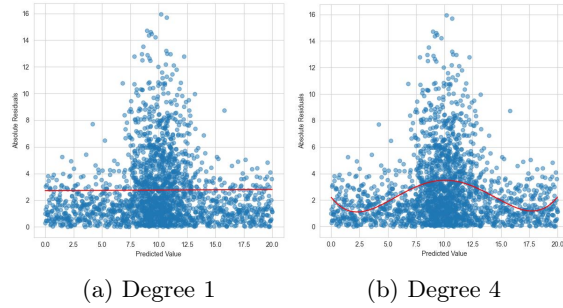


Fig. 1: Example of heteroscedasticity in (simulated) residuals that is not identified with linear models (a) but can be identified with polynomial models (b).

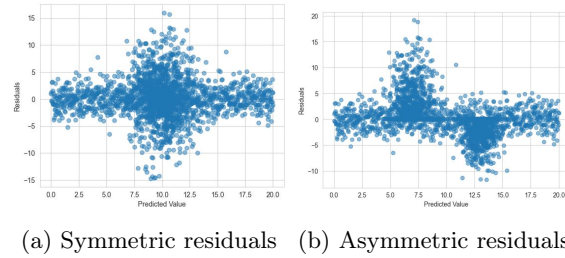


Fig. 2: Example of heteroscedasticity in (simulated) residual that can be modelled using squared or absolute residuals (a), or requires signed residuals (b).

Analysing complex cases of heteroscedasticity in AI errors with polynomial models allows both to better identify fairness issues, and to better estimate the uncertainty of AI results. Thus it provides insights into both the ethical implications of biased AI systems, and the robustness of AI systems.

While homoscedastic residuals might be reasonably expected (e.g., as a constraint that guides the convergence of robust AI algorithms), it may not be the case for all specific data subsets. Heteroscedastic residuals may occur for specific data subsets (e.g., for data points with specific sensitive features), while remaining generally homoscedastic for the whole dataset.

Complex patterns of heteroscedasticity may also remain undetected if squared or absolute residuals are used (e.g., Fig. 2). Signed residuals allow to detect heteroscedasticity when positive and negative residuals have different distributions. We thus model signed residuals as:

$$\hat{y} - y = \beta_0 + \beta_1 \mathbf{x} + \beta_2 \mathbf{x}^2 + \dots + \beta_n \mathbf{x}^n \quad (5)$$

Selecting an appropriate heteroscedasticity model is essential to guarantee the validity and completeness of AI bias analysis. To contribute to understanding the applicability bias analysis methods, we study the applicability of linear and polynomial models with the following research questions:

RQ1 How efficient are polynomial models for modelling heteroscedasticity?

RQ2 How to interpret polynomial models to identify fairness issues due to heteroscedasticity?

RQ3 How does using signed or absolute residuals impact the detection of heteroscedasticity?

We will assess the performance of polynomial regression models across different levels of complexity (i.e., polynomials of degrees 1 to 6), applied to the dataset from Obermeyer et al. [11] where racial discrimination was identified. We will review a wide array of evaluation metrics for each heteroscedasticity model, and each of their coefficient (e.g., p-values, deviance). While our results clearly demonstrate the advantages of using polynomial models, we cannot assume that these models are sufficient to detect all forms of heteroscedasticity. Other approaches to modelling complex heteroscedasticity remain unexplored, e.g., parametric, non-parametric, or multivariate models with weighted sensitive features.

2 Related Work

Fairness metrics in regression problems focus either on statistical parity (which is not concerned with residuals) or loss-based (e.g., bounded group loss) which uses squared residuals [1]. These do not account for heteroscedasticity, i.e., fairness criteria may be satisfied while a protected group is more impacted by heteroscedasticity. Thus, with such approach an AI model may be considered fair while a protected group still suffers from systematically higher AI errors (e.g., for predicted values in [7.5, 12.5] in Fig. 1).

Heteroscedasticity can be identified with statistical tests of normality applied to the residuals. However, this approach does not model the magnitude of heteroscedasticity, nor how it can impact protected groups differently, e.g., for data points with specific predicted or actual values (\hat{y} or y).

Existing methods for modeling heteroscedasticity use linear models [2, 3, 8, 9]. A significant slope β_1 (1-2) indicates heteroscedasticity. These methods are not suited to identify all patterns of heteroscedasticity (e.g., Fig.1-2). These methods are also impacted by the choice of residual measurement, e.g., absolute, squared, or signed residuals (1-2). Squared residuals are frequently used, but are sensitive to outliers and distort small residuals, especially those between $[-1, 1]$. To address some of these limitations, other methods use absolute residuals as the dependent variable [4, 6]. These methods mitigate excessive outlier influence, and provide an alternative perspective on the residual distribution by focusing on the magnitude of differences between observed and predicted values. However, complex patterns of residuals may not be identifiable using neither absolute or squared residuals (Fig. 2).

Su et al. [14] apply local polynomial regressions, fitting multiple polynomial equations to different regions of the data. This method is non-parametric, thus does not require prior knowledge of the form of the heteroscedasticity function. Yet, it is applied exclusively to squared residuals, and only provide local models. In contrast, our approach provides global models of heteroscedasticity, and address asymmetric patterns by using absolute residuals. Our approach also decreases the computational complexity, requires fewer data points, and is easier to interpret.

Hsiao et al. [6] present a model to analyze complex financial data, incorporating various forms of heteroscedasticity such as linear, nonlinear, curvilinear (e.g., polynomial regressions), and composition functions. However, this method only uses squared residuals and its goal is to incorporating heteroscedasticity into the initial AI model. In contrast, our work also is also applied to signed residuals, and aims at identifying fairness issues with specific interpretations of polynomial parameters β .

3 Method

We used the dataset studied by Obermeyer et al. [11] because it contains racial bias creating fairness issues. It consists of a comprehensive collection of medical information, and AI-generated scores used throughout the USA to determine patient eligibility for a specialized healthcare program. Its data points represent patients from two racial groups: black and white. Each patient has a detailed set of medical information, including the variables *number of active chronic illness*, and *total medical expenditure* spent in hospital treatments. Their AI-generated risk scores range from 0 (low risk) to 100 (high risk), and are used to assess their health status. When an individual’s risk score exceeds the 55th percentile, they are referred for screening by a medical professional who determines their

eligibility for the healthcare program. Patients with risk scores above the 97th percentile are automatically admitted to the program.

Obermeyer et al. found significant biases when comparing signs of illnesses (*number of active chronic conditions*) at a given risk score, showing that black patients are considerably sicker than white patients at the same risk score. The AI bias is due to using *total medical expenses* as a predictor, as black patients often have less financial resources thus cannot afford to spend as much as white patients.

We used this dataset to create 2 experimental use cases: regression problems predicting (1) *number of active chronic conditions* and (2) *total medical expenses*, using risk score percentiles (Fig. 3). We analyzed the residuals from these regression problems, to assess potential heteroscedasticity. If heteroscedasticity is present, and impacts one race more than the other, a fairness issue is identified. We researched the effectiveness of polynomial models to identify such issues with heteroscedasticity.

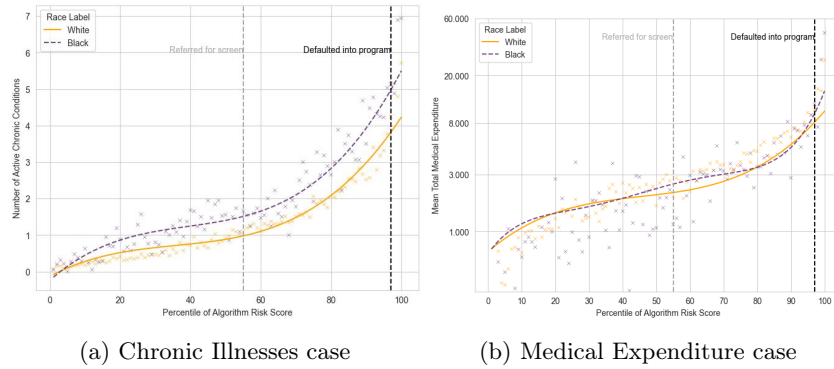


Fig. 3: Experimental data drawn from Obermeyer et al. (2019)

To ensure clarity in graphical representations, we opted for a data aggregation method akin to Obermeyer et al.’s approach, which entails the use of mean values per risk score percentile. It is imperative to clarify that the results presented in this paper do not directly mirror the levels of heteroscedasticity within the original dataset. Rather, they offer an overview of polynomial regression models’ performance in addressing heteroscedasticity.

Experimental Data We calculated the mean *active chronic illnesses* and the *total medical expenditure* per race and risk score percentile. This reduced the dataset from 48785 to 400 data points: 100 data points per race and variable (*active chronic illnesses* and the *total medical expenditure*). We then fit a polynomial model with degree 3 to the data points (Fig. 3).

We modelled the heteroscedascity in our two experimental use cases with polynomial regression of degrees 1 to 6, fitted with Ordinary Least Squares

(OLS). We assessed the heteroscedasticity modelling using well established metrics: mean squared error (MSE), p-values for each coefficient, deviance, and change in deviance between linear and polynomial models [10].

Deviance is calculated from the likelihood function, representing the probability of observing the residuals given the model's parameters $(\beta_0, \dots, \beta_n)$. Higher likelihood means that the observed residuals are more probable under that regression model. The likelihood function $L(\beta; r)$ is defined as the joint probability of the observed residuals given the parameters β :

$$L(\beta; r) = \prod_{i=1}^n f(r_i, \hat{r}_i, \beta)$$

where $f(r_i, \hat{r}_i, \beta)$ is the probability density function of the observed residuals r_i given the predicted residual \hat{r}_i and the model parameters β . We calculate deviance for eventually comparing models, thus its calculation can be simplified as:

$$D = -2L(\beta; r)$$

where $L(\beta; r)$ is the log-likelihood of the regression model. The negative twice log-likelihood ensures that larger deviance indicates poorer fit of the regression model. We also examined changes in deviance [10] between linear and polynomial models:

$$\Delta D = D_{polynomial} - D_{linear}$$

where D_{linear} is the deviance of the baseline linear model and $D_{polynomial}$ the deviance of a polynomial model. A negative ΔD suggests that the polynomial model is capturing the distribution of residuals more efficiently. This analysis helps select a model that strikes a balance between capturing the distribution of the residuals and avoiding unnecessary complexity. When the difference in deviance between two models is minimal, the simpler model may be preferable due to its potential for improved generalisability.

Table 1: Deviance D and difference with baseline ΔD for Black and White races.

| | | | Degree | | | | | |
|-------------------|--------------------|------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| | | | 1 | 2 | 3 | 4 | 5 | 6 |
| Chronic Illness | Signed Residuals | D | 1.2e ¹ | 8.5 | 8.1 | -4.2e ¹ | -7.8e ¹ | -9.5e ¹ |
| | | ΔD | 1.1e ² | 1.1e ² | 1.1e ² | 9.9e ¹ | 9.9e ¹ | 9.7e ¹ |
| | | | 0 | -3.5 | -3.9 | -5.4e ¹ | -9.0e ¹ | -1.1e ² |
| | Absolute Residuals | D | -4.4e ¹ | -6.2e ¹ | -1.0e ² | -1.2e ² | -1.3e ² | -1.5e ² |
| | | ΔD | 1.1e ¹ | 7.7 | 5.0 | 4.7 | 2.7 | -5.6 |
| | | | 0 | -1.8e ¹ | -5.7e ¹ | -7.5e ¹ | -8.7e ¹ | -1.1e ² |
| Total Expenditure | Signed Residuals | D | 1.8e ³ | 1.7e ³ | 1.7e ³ | 1.7e ³ | 1.7e ³ | 1.7e ³ |
| | | ΔD | 1.9e ³ | 1.9e ³ | 1.9e ³ | 1.9e ³ | 1.9e ³ | 1.8e ³ |
| | | | 0 | -1.5e ¹ | -3.0e ¹ | -4.9e ¹ | -7.3e ¹ | -9.7e ¹ |
| | Absolute Residuals | D | 1.8e ³ | 1.7e ³ | 1.7e ³ | 1.7e ³ | 1.7e ³ | 1.7e ³ |
| | | ΔD | 1.9e ³ | 1.9e ³ | 1.9e ³ | 1.8e ³ | 1.8e ³ | 1.8e ³ |
| | | | 0 | -2.0e ¹ | -3.4e ¹ | -5.5e ¹ | -7.7e ¹ | -1.1e ² |

4 Results

Table 1 shows the deviance values, and the difference in deviance (ΔD) with to the baseline model (Degree 1), for both the Black and White races. Deviance decreases when increasing the polynomial degree, indicating an improvement in goodness of fit. However, the decrease in deviance for the Black race is smaller than for the White race (especially for the Chronic Illness case). Thus the relationship between polynomial degree and deviance reduction varies between races. It indicates a variability in the predictive accuracy of the regression models across different racial groups. Regarding difference in deviance with baseline model (ΔD), it consistently decreases as the polynomial degree increases. This indicates that the polynomial models are more accurate for modelling the residuals and their heteroscedasticity.

Table 2: Modelling of signed residuals for the Chronic Illness case for **Black** and **White** races. Preferred model is in bold, parentheses show p-values.

| Degree | β_0 | β_1 | β_2 | β_3 | β_4 | β_5 | β_6 | MSE |
|--------|--|--|--|---|--|--|--|--------------------------------|
| 1 | $-1.3e^{-1}$ ($8.5e^{-3}$) $-1.1e^{-1}$ ($7.3e^{-2}$) | $2.4e^{-3}$ ($8.5e^{-3}$) $2.7e^{-3}$ ($7.3e^{-2}$) | | | | | | $6.7e^{-2}$ $1.8e^{-1}$ |
| 2 | $-2.0e^{-2}$ ($2.6e^{-1}$) $5.0e^{-2}$ ($2.7e^{-1}$) | $-4.0e^{-3}$ ($2.6e^{-1}$) $-6.6e^{-3}$ ($2.7e^{-1}$) | $6.4e^{-5}$ ($6.7e^{-2}$) $9.1e^{-5}$ ($1.1e^{-1}$) | | | | | $6.6e^{-2}$ $1.8e^{-1}$ |
| 3 | $-6.5e^{-2}$ ($9.0e^{-1}$) $-1.5e^{-2}$ ($9.5e^{-1}$) | $1.1e^{-3}$ ($9.0e^{-1}$) $9.9e^{-4}$ ($9.5e^{-1}$) | $-6.3e^{-5}$ ($7.6e^{-1}$) $-9.5e^{-5}$ ($7.8e^{-1}$) | $8.3e^{-7}$ ($5.4e^{-1}$) $1.0e^{-6}$ ($5.8e^{-1}$) | | | | $6.6e^{-2}$ $1.8e^{-1}$ |
| 4 | $4.6e^{-1}$ ($4.6e^{-5}$) $3.6e^{-1}$ ($9.8e^{-2}$) | $-9.8e^{-2}$ ($1.3e^{-9}$) $-7.0e^{-2}$ ($1.9e^{-2}$) | $4.3e^{-3}$ ($5.7e^{-11}$) $3.1e^{-3}$ ($1.1e^{-2}$) | $-6.6e^{-5}$ ($1.6e^{-11}$) $-4.7e^{-5}$ ($8.3e^{-3}$) | $3.3e^{-7}$ ($7.4e^{-12}$) $2.4e^{-7}$ ($6.4e^{-3}$) | | | $4.1e^{-2}$ $1.7e^{-1}$ |
| 5 | $4.2e^{-2}$ ($7.1e^{-1}$) $3.0e^{-1}$ ($2.6e^{-1}$) | $1.8e^{-2}$ ($4.2e^{-1}$) $-5.4e^{-2}$ ($3.1e^{-1}$) | $-3.6e^{-3}$ ($8.6e^{-3}$) $1.9e^{-3}$ ($5.5e^{-1}$) | $1.4e^{-4}$ ($5.6e^{-5}$) $-1.8e^{-5}$ ($8.3e^{-1}$) | $-2.0e^{-6}$ ($4.5e^{-7}$) $-8.9e^{-8}$ ($9.2e^{-1}$) | $9.1e^{-9}$ ($6.7e^{-9}$) $1.3e^{-9}$ ($7.1e^{-1}$) | | $2.9e^{-2}$ $1.7e^{-1}$ |
| 6 | $3.6e^{-1}$ ($8.0e^{-3}$) $5.8e^{-1}$ ($7.7e^{-2}$) | $-9.3e^{-2}$ ($5.8e^{-3}$) $-1.6e^{-1}$ ($7.1e^{-2}$) | $7.1e^{-3}$ ($1.3e^{-2}$) $1.2e^{-2}$ ($1.1e^{-1}$) | $-2.8e^{-4}$ ($8.4e^{-3}$) $-4.1e^{-4}$ ($1.3e^{-1}$) | $5.7e^{-6}$ ($2.5e^{-3}$) $7.1e^{-6}$ ($1.4e^{-1}$) | $-5.8e^{-8}$ ($4.3e^{-4}$) $-6.1e^{-8}$ ($1.4e^{-1}$) | $2.2e^{-10}$ ($5.4e^{-5}$) $2.1e^{-10}$ ($1.3e^{-1}$) | $2.4e^{-2}$ $1.7e^{-1}$ |

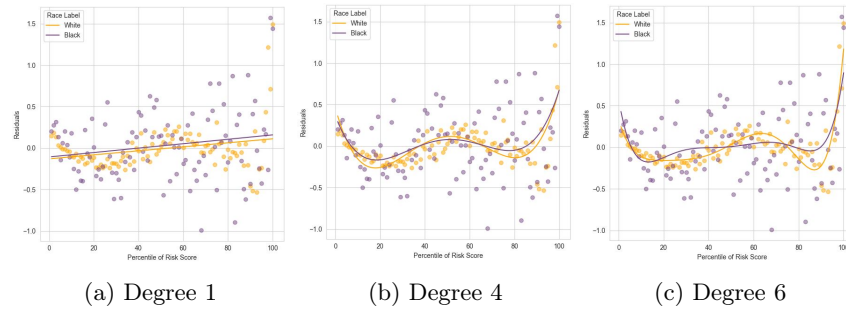


Fig. 4: Signed residuals vs risk score, for the Chronic Illness case.

4.1 Signed Residuals

Chronic Illness Case: Table 2 shows the regression parameters (with corresponding p-values) and MSE for polynomial models of degrees 1 to 6.

For the Black race, no model achieves statistical significance ($p < 0.05$) for all their parameters β . All models have equivalent mean squared error (MSE), although models with degrees 4 to 6 having slightly lower MSE. With statistical significance for all parameters β except the intercept, and the lowest MSE, the model of degree 4 seems preferable. The non-significant intercept indicates that the model is less accurate when risk scores are close to 0.

For the White race, MSE decreases as the models' degrees increase. It indicates that the polynomial regression models are more accurate for modelling the residuals. The model of degree 4 has parameters β with the lowest p-values. However, with degree 6 all the model's coefficient are also statistically significant, and MSE is lower by 41.5%. It is thus the preferred model for the White race.

These results indicate that modelling heteroscedasticity may require different models for different populations, which is crucial for assessing fairness. The varying MSE and significance of parameters β across different polynomial degrees, and across races, underscore the importance of tailoring models to specific demographic groups. Thus, in practice, a one-size-fits-all approach to heteroscedasticity modelling may not be appropriate for fairness assessments.

Figure 4 shows the distribution of signed residuals, the baseline model, and the preferred models. Both races have outliers when risk scores are around the maximum value (100), which slight skews the models. For the Black race, we observe a cone-shaped pattern of heteroscedasticity, where the range of residuals widens as risk scores increase. This pattern may be better modelled using absolute or squared residuals, since it is symmetrical around the zero line. For the White race, residuals are asymmetrical around the zero, showing a complex pattern that may not be identifiable using absolute or squared residuals. Thus modelling heteroscedasticity may require not only different model complexity (e.g., polynomial of degree 4 or 6), but also different types of residuals measurements (e.g., signed or absolute residuals).

Total Expenditure Case: Table 3 shows the regression parameters (with corresponding p-values) and MSE for polynomial models of degrees 1 to 6.

For both the Black and White races, MSE significantly decreases as the polynomial degree increases. For the White race, polynomial degrees 1, 3 and 5 exhibit parameters β that are all statistically significant. The model with degree 5 has the lowest MSE (50% lower than the baseline with degree 1). For the Black race, the only model where all parameters β are statistically significant is the baseline (degree 1). However, when disregarding the intercept, all the parameters β of models with degrees 2, 3, 5, and 6 are statistically significant. Since their MSEs are much lower than the baseline, polynomial models remain preferable. This suggests that estimating the intercept, e.g., estimating residuals for risk scores around 0, might be a recurrent issue with polynomial models.

Table 3: Modelling of signed residuals for the Medical Expenditure case for **Black** and **White** races. Preferred model is in bold, parentheses show p-values.

| degree | β_0 | β_1 | β_2 | β_3 | β_4 | β_5 | β_6 | MSE |
|--------|------------------------------|------------------------------|--------------------------------|---------------------------------|---------------------------------|---------------------------------|--------------------------------|----------|
| 1 | $-1.1e^3$ ($9.1e^{-4}$) | $3.7e^2$ ($6.75e^{-9}$) | | | | | | $2.8e^5$ |
| | $-1.5e^3$ ($3.4e^{-2}$) | $3.6e^1$ ($3.3e^{-3}$) | | | | | | $1.2e^7$ |
| 2 | $2.8e^2$ ($5.5e^{-1}$) | $-4.8e^1$ ($3.0e^{-2}$) | $8.3e^{-1}$ ($1.2e^{-4}$) | | | | | $2.4e^6$ |
| | $8.0e^2$ ($4.3e^{-1}$) | $-9.8e^1$ ($3.4e^{-2}$) | 1.3 ($3.3e^{-3}$) | | | | | $1.1e^7$ |
| 3 | $-1.3e^3$ ($3.3e^{-2}$) | $1.4e^2$ ($9.2e^{-3}$) | -3.7 ($2.2e^{-3}$) | $3.0e^{-2}$ ($1.8e^{-4}$) | | | | $2.1e^6$ |
| | $-2.0e^3$ ($1.3e^{-1}$) | $2.2e^2$ ($4.6e^{-2}$) | -6.6 ($1.1e^{-2}$) | $5.2e^{-2}$ ($2.1e^{-3}$) | | | | $9.8e^6$ |
| 4 | $6.4e^2$ ($3.6e^{-1}$) | $-2.3e^2$ ($1.7e^{-2}$) | $1.3e^1$ ($1.4e^{-3}$) | $-2.2e^{-1}$ ($2.0e^{-4}$) | $1.2e^{-3}$ ($2.5e^{-5}$) | | | $1.8e^6$ |
| | $1.0e^3$ ($5.3e^{-1}$) | $-3.4e^2$ ($1.2e^{-1}$) | $1.8e^1$ ($3.8e^{-2}$) | $-3.3e^{-1}$ ($1.2e^{-2}$) | $1.9e^{-3}$ ($3.6e^{-3}$) | | | $9.1e^6$ |
| 5 | $-1.7e^3$ ($3.5e^{-2}$) | $4.1e^2$ ($9.1e^{-3}$) | $-3.1e^1$ ($1.2e^{-3}$) | $9.2e^{-2}$ ($1.4e^{-4}$) | $-1.1e^{-2}$ ($1.7e^{-5}$) | $5.0e^{-5}$ ($2.0e^{-6}$) | | $1.4e^6$ |
| | $-3.0e^3$ ($1.1e^{-1}$) | $7.8e^2$ ($3.6e^{-2}$) | $-5.8e^1$ ($1.1e^{-2}$) | 1.7 ($3.3e^{-3}$) | $-2.0e^{-2}$ ($1.1e^{-3}$) | $8.8e^{-5}$ ($3.7e^{-4}$) | | $8.0e^6$ |
| 6 | $7.2e^2$ ($3.9e^{-1}$) | $-4.9e^2$ ($2.9e^{-2}$) | $5.5e^1$ ($4.1e^{-3}$) | -2.5 ($6.3e^{-4}$) | $5.1e^{-2}$ ($9.2e^{-5}$) | $-4.9e^{-4}$ ($1.3e^{-5}$) | $2.0e^{-6}$ ($2.0e^{-6}$) | $1.1e^6$ |
| | $2.4e^3$ ($2.4e^{-1}$) | $-1.3e^3$ ($2.0e^{-2}$) | $1.4e^2$ ($2.9e^{-3}$) | -6.1 ($5.2e^{-4}$) | $1.2e^{-1}$ ($1.1e^{-4}$) | $-1.2e^{-3}$ ($2.4e^{-5}$) | $4.0e^{-6}$ ($6.0e^{-6}$) | $6.5e^6$ |

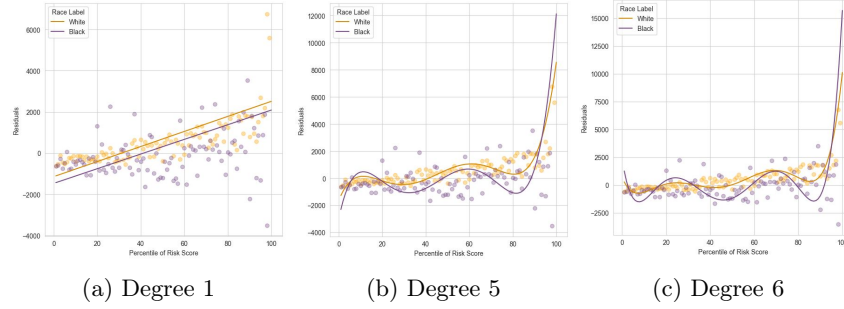


Fig. 5: Signed residuals vs risk score, for the Medical Expenditure case.

Figure 5 shows the distribution of signed residuals, the baseline model, and the preferred models. For visual clarity, 3 outliers with high risk scores are removed from these graphs¹ but were not removed when fitting the polynomial models. These outliers skew the models, especially the polynomial models.

For both races, we observe a cone-shaped pattern of heteroscedasticity, where the range of residuals widens as risk scores increase. However, with high risk scores, residuals tend to be negative for the Black race (i.e., risk scores are underestimated), and positive for the White race (i.e., risk scores are over-estimated). Furthermore, for all risk scores, the range of residuals is larger for the Black race. These discrepancies create fairness issues. Such cone-shaped patterns of heteroscedasticity might be better modelled using absolute residuals, which we investigate next.

¹ (1) Risk score = 99, residuals ≈ 1400 , race = Black. (2) Risk score = 100, residuals ≈ 1700 , race = White, (3) Risk score = 100, residuals ≈ 3000 , race = Black

4.2 Absolute Residuals

Chronic Illness Case: Table 4 shows the regression parameters (with corresponding p-values) and MSE for polynomial models of degrees 1 to 6.

For the White race models with degree 2 and 6 both have all their coefficient statistically significant. However, degree 6 has lower MSE (57.6% lower than degree 2). For the Black race, only the model with degree 1 has all its parameters β statistically significant. However, this model has the highest MSE. But other models offer only a slightly lower MSE (e.g., 10% lower at most with degree 6).

Figure 6 shows the distribution of signed residuals, the baseline model, and the preferred models. There are outliers around the maximum risk score, which skews the models. For the Black race, the absolute residuals have a wider range that consistently exceed those of the White race, and increases with the risk scores. In contrast, the White race shows less variability in the residuals.

Table 4: Modelling of absolute residuals for the Chronic Illness case for Black and White races. Preferred model is in bold, parentheses show p-values.

| degree | β_0 | β_1 | β_2 | β_3 | β_4 | β_5 | β_6 | MSE |
|--------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|-------------|
| 1 | $4.3e^{-2}$ ($2.8e^{-1}$) | $2.5e^{-3}$ ($4.9e^{-4}$) | | | | | | $3.9e^{-2}$ |
| | $1.6e^{-1}$ ($3.2e^{-3}$) | $3.4e^{-3}$ ($2.8e^{-4}$) | | | | | | $6.7e^{-2}$ |
| 2 | $2.3e^{-1}$ ($8.8e^{-5}$) | $-8.3e^{-3}$ ($1.4e^{-3}$) | $1.1e^{-4}$ ($2.8e^{-5}$) | | | | | $3.3e^{-2}$ |
| | $2.7e^{-1}$ ($9.8e^{-4}$) | $-3.0e^{-3}$ ($4.0e^{-1}$) | $6.3e^{-5}$ ($6.9e^{-2}$) | | | | | $6.5e^{-2}$ |
| 3 | $-5.7e^{-2}$ ($3.6e^{-1}$) | $2.4e^{-2}$ ($1.2e^{-5}$) | $-7.0e^{-4}$ ($9.8e^{-8}$) | $5.0e^{-6}$ ($1.2e^{-9}$) | | | | $2.2e^{-2}$ |
| | $1.5e^{-1}$ ($1.6e^{-1}$) | $1.0e^{-2}$ ($2.5e^{-1}$) | $-2.7e^{-4}$ ($2.0e^{-1}$) | $2.0e^{-6}$ ($1.1e^{-1}$) | | | | $6.4e^{-2}$ |
| 4 | $1.4e^{-1}$ ($5.4e^{-2}$) | $-1.3e^{-2}$ ($1.9e^{-1}$) | $9.6e^{-4}$ ($1.7e^{-2}$) | $-2.0e^{-5}$ ($9.0e^{-4}$) | $1.3e^{-7}$ ($3.2e^{-5}$) | | | $1.9e^{-2}$ |
| | $2.0e^{-1}$ ($1.4e^{-1}$) | $9.5e^{-4}$ ($9.6e^{-1}$) | $1.5e^{-4}$ ($8.4e^{-1}$) | $-4.2e^{-6}$ ($7.0e^{-1}$) | $3.1e^{-8}$ ($5.6e^{-1}$) | | | $6.5e^{-2}$ |
| 5 | $-2.7e^{-2}$ ($7.6e^{-1}$) | $3.4e^{-2}$ ($4.9e^{-2}$) | $-2.2e^{-3}$ ($3.2e^{-2}$) | $6.3e^{-5}$ ($1.5e^{-2}$) | $-8.0e^{-7}$ ($4.8e^{-3}$) | $3.7e^{-9}$ ($1.1e^{-3}$) | | $1.7e^{-2}$ |
| | $6.5e^{-2}$ ($6.9e^{-1}$) | $3.8e^{-2}$ ($2.4e^{-1}$) | $-2.4e^{-3}$ ($2.3e^{-1}$) | $6.2e^{-5}$ ($2.1e^{-1}$) | $-7.1e^{-7}$ ($1.9e^{-1}$) | $2.9e^{-9}$ ($1.7e^{-1}$) | | $6.4e^{-2}$ |
| 6 | $2.1e^{-1}$ ($3.0e^{-2}$) | $-5.5e^{-2}$ ($3.2e^{-2}$) | $6.3e^{-3}$ ($4.1e^{-3}$) | $-2.7e^{-4}$ ($9.3e^{-4}$) | $5.3e^{-6}$ ($2.7e^{-4}$) | $-5.0e^{-8}$ ($8.7e^{-5}$) | $1.8e^{-10}$ ($2.7e^{-5}$) | $1.4e^{-2}$ |
| | $3.7e^{-1}$ ($5.7e^{-2}$) | $-7.8e^{-2}$ ($1.4e^{-1}$) | $8.8e^{-3}$ ($4.8e^{-2}$) | $-3.7e^{-4}$ ($2.2e^{-2}$) | $7.4e^{-6}$ ($1.3e^{-2}$) | $-6.7e^{-8}$ ($8.1e^{-3}$) | $2.3e^{-10}$ ($5.6e^{-3}$) | $6.0e^{-2}$ |

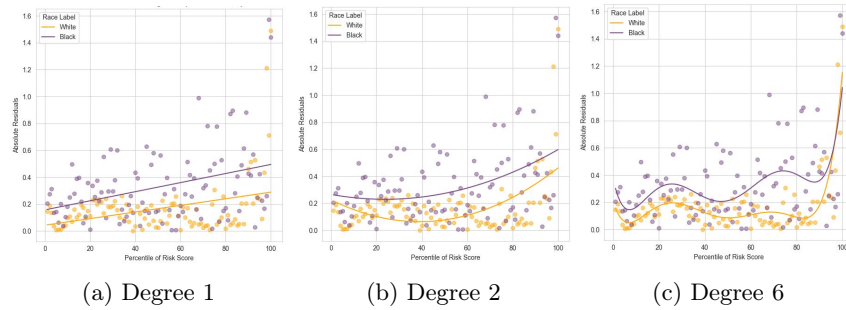


Fig. 6: Absolute residuals vs risk score, for the Chronic Illness case.

Medical Expenditure case: Table 5 shows the regression parameters (with corresponding p-values) and MSE for polynomial models of degrees 1 to 6.

The parameters β for degrees 2, 4, and 6 are all statistically significant for the White race, with degree 6's MSE being considerably lower than degree 1's MSE (44.4% lower). For the Black race, only degree 2 has all its parameters β statistically significant. For degree 6, only the intercept β_0 is not statistically significant. The MSE of degree 6 is much lower than the MSE of degree 2 (51.6% lower), thus it is also a well-performing model.

Figure 7 shows the distribution of signed residuals, the baseline model, and the preferred models. For the Black race, the spread of residuals is notably wider than for the White race. For both races, residuals increase with the risk scores, and this increase is more pronounced around the maximum risk scores. There are outliers for high risk scores, which significantly skew the models, especially for polynomial models, and for the Black race.

Table 5: Modelling of absolute residuals for the Medical Expenditure case for **Black** and **White** races. Preferred model is in bold, parentheses show p-values.

| degree | β_0 | β_1 | β_2 | β_3 | β_4 | β_5 | β_6 | MSE |
|--------|---|--|---|---|--|---|--|----------------------------|
| 1 | $-5.4e^2$ ($1.2e^{-1}$) | $2.8e^1$ ($3.0e^{-6}$) | | | | | | $2.8e^6$ |
| | $-4.0e^2$ ($5.4e^{-1}$) | $3.3e^1$ ($4.0e^{-3}$) | | | | | | $1.0e^7$ |
| 2 | $1.1e^3$ ($2.2e^{-2}$) | $-6.7e^1$ ($1.2e^{-5}$) | $9.5e^{-1}$ ($2.3e^{-3}$) | | | | | $2.4e^6$ |
| | $2.0e^3$ ($4.0e^{-2}$) | $-1.1e^2$ ($1.6e^{-2}$) | 1.4 ($1.2e^{-3}$) | | | | | $9.3e^6$ |
| 3 | $-4.7e^2$ ($4.3e^{-1}$) | $1.1e^2$ ($3.3e^{-3}$) | -3.5 ($2.7e^{-2}$) | $2.9e^{-2}$ ($1.9e^{-4}$) | | | | $2.1e^6$ |
| | $-1.4e^3$ ($2.5e^{-1}$) | $2.8e^2$ ($6.7e^{-3}$) | -8.1 ($6.9e^{-4}$) | $6.2e^{-2}$ ($6.9e^{-5}$) | | | | $7.9e^6$ |
| 4 | $1.5e^3$ ($3.0e^{-2}$) | $-2.6e^2$ ($7.4e^{-4}$) | $1.3e^1$ ($6.1e^{-3}$) | $-2.3e^{-1}$ ($1.1e^{-4}$) | $1.3e^{-3}$ ($1.3e^{-5}$) | | | $1.7e^6$ |
| | $2.5e^3$ ($7.7e^{-2}$) | $-4.4e^2$ ($1.9e^{-2}$) | $2.4e^1$ ($2.0e^{-3}$) | $-4.3e^{-1}$ ($2.1e^{-4}$) | $2.4e^{-3}$ ($2.4e^{-5}$) | | | $6.7e^6$ |
| 5 | $-6.6e^2$ ($3.9e^{-1}$) | $3.4e^2$ ($2.9e^{-3}$) | $-2.8e^1$ ($2.6e^{-2}$) | $8.5e^{-1}$ ($3.3e^{-4}$) | $-1.1e^{-2}$ ($4.2e^{-5}$) | $4.7e^{-5}$ ($5.0e^{-6}$) | | $1.4e^6$ |
| | $-1.6e^3$ ($2.9e^{-1}$) | $6.9e^2$ ($2.6e^{-2}$) | $-5.4e^1$ ($4.6e^{-3}$) | 1.6 ($7.7e^{-4}$) | $-2.0e^{-2}$ ($1.2e^{-4}$) | $8.9e^{-5}$ ($1.8e^{-5}$) | | $5.5e^6$ |
| 6 | $1.9e^3$ ($2.5e^{-2}$) | $-6.1e^2$ ($6.3e^{-3}$) | $6.3e^1$ ($8.9e^{-4}$) | -2.7 ($1.2e^{-4}$) | $5.5e^{-2}$ ($1.7e^{-5}$) | $-5.3e^{-4}$ ($2.3e^{-6}$) | $2.0e^{-6}$ ($3.2e^{-7}$) | $1.0e^6$ |
| | $2.8e^3$ ($9.5e^{-2}$) | $-1.0e^3$ ($2.9e^{-2}$) | $1.1e^2$ ($5.2e^{-3}$) | -4.8 ($1.0e^{-3}$) | $9.7e^{-2}$ ($2.1e^{-4}$) | $-9.3e^{-4}$ ($4.0e^{-5}$) | $3.2e^{-6}$ ($8.0e^{-6}$) | $4.5e^6$ |

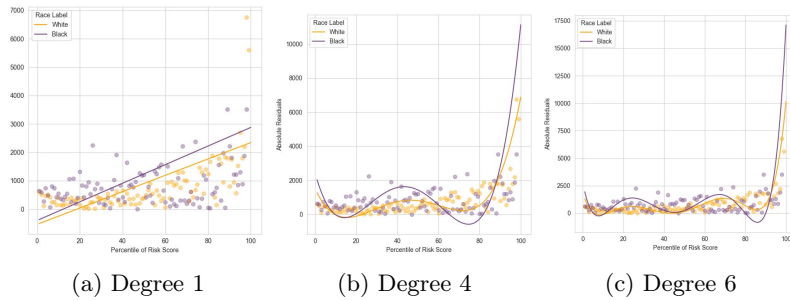


Fig. 7: Absolute residuals vs risk score, for the Medical Expenditure case.

4.3 Preferred Model

Table 6 recaps the evaluation of the preferred models for all use cases. We can observe that the models' performance is worse for the Black race compared to the White race, except for the Chronic Illness case with signed residuals. Thus modelling the residuals is generally more uncertain for the Black race.

Signed vs absolute residuals: For the Chronic Illness case, with the Black race, the linear model using absolute residuals is statistically more reliable than the one using signed residuals (deviance, p-values of β parameters). However, its MSE is much higher. This is consistent with our assumption that the Black race's cone-shaped pattern of heteroscedasticity, horizontally centered around the zero line (Fig. 4), can be modelled with simple linear models and absolute residuals.

For the Medical Expenditure case, cone-shaped pattern of heteroscedasticity were also observed (Fig. 5). For these, the preferred models using absolute residuals are less complex than those using signed residuals (e.g., degree 2 vs 6 for the Black race). The deviance is similar with signed or absolute residuals, but MSE is lower with signed residuals. These cone-shaped patterns were not well-centered around the horizontal zero line, which explains why using signed residuals is preferable.

We conclude that using signed residuals is preferable to absolute residuals if the pattern of heteroscedasticity is not symmetrical around the horizontal zero line. However, with signed residuals the intercept parameter β_0 might be less reliable (e.g., not statistically significant), e.g., impeding the modelling of residuals and heteroscedasticity for risk scores around zero.

Linear vs polynomial models: A linear model is preferable only for one case out of eight. For the remaining cases, compared to linear models, MSE is significantly lower with polynomial models, and the reduction of deviance ΔD is noticeable (especially for the White race). The MSE consistently decreased as model complexity increased from degree 1 to 6. We conclude that using polynomial regression enhances the ability to capture the distribution of residuals and the heteroscedasticity.

Table 6: Preferred models for each use case and Black and White races.

| | | Degree | MSE | Deviance | ΔD | $p_\beta < 0.05$ |
|-------------------|--------------------|--------|-------------|-----------|------------|-------------------|
| Chronic Illness | Signed Residuals | 4 | $4.1e^{-2}$ | $-4.2e^1$ | $-5.4e^1$ | Yes |
| | | 4 | $1.7e^{-1}$ | $9.9e^1$ | $-1.1e^1$ | All but β_0 |
| | Absolute Residuals | 6 | $1.4e^{-2}$ | $-1.5e^2$ | $-1.1e^2$ | Yes |
| | | 1 | $6.7e^{-2}$ | $1.1e^1$ | 0 | Yes |
| Total Expenditure | Signed Residuals | 5 | $1.4e^6$ | $1.7e^3$ | $-7.3e^1$ | Yes |
| | | 6 | $6.5e^6$ | $1.8e^3$ | $-6.4e^1$ | All but β_0 |
| | Absolute Residuals | 4 | $1.7e^6$ | $1.7e^3$ | $-5.5e^1$ | Yes |
| | | 2 | $9.3e^6$ | $1.9e^3$ | $-1.1e^1$ | Yes |

5 Discussion

What is the effectiveness of polynomial regression in modelling complex cases of heteroscedasticity? Polynomial models offers an important advantage over linear models in capturing curvature within residuals. This curvature allows for more flexibility when modelling heteroscedasticity where the distribution of residuals is not linear or constant, an ability that linear regression, limited to straight lines, does not have. Our finding shows that relying solely on linear regression for modelling heteroscedasticity may be insufficient. This approach requires to tune the model complexity (i.e., the polynomial degrees). Beware that even if non-significant parameters are found for a given degree, and significant parameters may yet be found for a higher degree higher (e.g., Table 4). However, the model complexity may not be increased indefinitely as this may lead to overfitting.

We do claim that polynomial models are necessarily the best approach to modelling heteroscedasticity: other approaches remain unexplored, such as non-parametric models, and multivariate models to account for feature interactions, or intersectionality [15] (e.g., by weighing sensitive features).

How to interpret polynomial models to determine the fairness issues due to heteroscedasticity? Heteroscedasticity can be identified when at least one parameters β_1, \dots, β_n is non-zero and statistically significant (regardless of intercept $beat_0$), with a model using either absolute or signed residuals (and without overfitting). A model of heteroscedasticity must be fitted to each protected group. If all protected groups are impacted by heteroscedasticity, it remains challenging the assess whether one group is more impacted than the other.

To address this issue, the β parameters could be compared, e.g., fairness issues arise if β_n is higher for a protected group compared to another group. However, this approach is difficult to apply if the heteroscedasticity models have different complexity for each group, or if parameters β_n are not significant for each groups. Other approaches to explore in future work include (1) studying the difference between residuals models for each groups (e.g., $f_{Black}(\hat{y} - y) - f_{White}(\hat{y} - y)$), or (2) identifying the range of input values x for which residuals is higher, and the magnitude of this heteroscedastic increase in residuals (e.g., fairness issues arise if within the same range of x values, residuals are higher for one protected group).

What is the difference in effectiveness between using signed residuals and absolute residuals for detecting heteroscedasticity? The difference in effectiveness between using signed residuals and absolute residuals for detecting heteroscedasticity is complex to determine. Using absolute residuals can reduce the cumulative prediction errors of polynomial models, and give a more accurate indication of heteroscedasticity (e.g., with lower MSEs and β parameters' p-values). This especially applies to patterns of heteroscedasticity that are symmetrical around the horizontal zero line, e.g., when using signed residuals.

However, the models using signed residuals performed systematically better in Medical Expenditure case, and offered much lower MSE for the Black race in the Chronic Illness case. Signed residuals preserve information about over- or under-estimation. The identification of systematic biases in AI models can depend on this information. Signed residuals also offer crucial information to assess the practical impacts of heteroscedasticity, as over- or under-estimations can have opposite consequences.

Finally, the impact of outliers may be different with signed or absolute residuals which requires future work.

6 Conclusion

This study highlights the advantages and limitations of polynomial regression models in detecting and modelling heteroscedasticity, compared to traditional linear methods. We demonstrate the polynomial model’s ability to capture patterns with curvature. This is a distinct advantage for modelling complex distributions of residuals compared to linear models, which are limited to straight lines. This study also underscores the importance of using multiple evaluation methods, such as β parameters’ p-values, deviance, mean squared error (MSE), and data visualization, to support model selection. This study also underlines the importance of considering diverse modelling approaches tailored to each protected group, as there is no one-size-fits-all approach. The potential of polynomial regression to offer a deeper understanding of complex distributions of residuals is evident, yet challenges remain with comparing the levels of heteroscedasticity across protected groups, and with the impact of outliers. Future work must address these issues, but also explore other approaches to modelling heteroscedasticity with multivariate and non-parametric models.

References

1. Agarwal, A., Dudík, M., Wu, Z.S.: Fair regression: Quantitative definitions and reduction-based algorithms. In: International Conference on Machine Learning. pp. 120–129. PMLR (2019)
2. Breusch, T.S., Pagan, A.R.: A simple test for heteroscedasticity and random coefficient variation. *Econometrica* **47**(5), 1287 (September 1979). <https://doi.org/10.2307/1911963>
3. Cook, R.D., Weisberg, S.: and university of minnesota school of statistics department of applied statistics, “diagnostics for heteroscedasticity in regression. Technical Report 405,” (May 1982)
4. Glejser, H.: A new test for heteroskedasticity. *Journal of the American Statistical Association* **64**(325), 316–323 (March 1969). <https://doi.org/10.1080/01621459.1969.10500976>
5. Goldfeld, S.M., Quandt, R.E.: Some tests for homoscedasticity. *Journal of the American Statistical Association* **60**(310), 539–547 (June 1965). <https://doi.org/10.1080/01621459.1965.10480811>

6. Hsiao, C.W., Chan, Y.C., Lee, M.Y., Lu, H.P.: Heteroscedasticity and precise estimation model approach for complex financial time-series data: An example of taiwan stock index futures before and during covid-19. *Mathematics* **9**(21), 2719 (October 2021). <https://doi.org/10.3390/math9212719>
7. J. L. A. Mattu Lauren Kirchner, S.: How we analyzed the compas recidivism algorithm (2020), <https://www.propublica.org/Article/How-We-Analyzed-the-Compas-Recidivism-Algorithm>, proPublica, Feb. 29
8. Klein, A.G.: 1. C. Gerhard, R. D. Büchner, S. Diestel, and K. Schermelleh-Engel, “The detection of heteroscedasticity in regression models for psychological data,” *Psychological Test and Assessment Modeling* (December 2016)
9. MacKinnon, J.G., White, H.: Some heteroskedasticity consistent covariance matrix estimators with improved finite sample properties (1983), <http://hdl.handle.net/10419/189084>
10. McCullagh, P., Nelder, J.A.: *Generalized linear models* (2019)
11. Obermeyer, Z., Powers, B., Vogeli, C., Mullainathan, S.: Dissecting racial bias in an algorithm used to manage the health of populations. *Science* **366**(6472), 447–453 (October 2019), <https://gitlab.com/labsysmed/dissecting-bias>
12. Rosopa, P.J., University, C., Schaffer, M.M., Walmart, A.N.S., University, W.K.: *Managing heteroscedasticity in general linear models* (2013)
13. Skeem, J.: University of california (June 2016), <https://ssrn.com/abstract=2687339>, berkeley, C. T. Lowenkamp, and Administrative Office, U. S. Courts, “Risk, Race, & Recidivism: Predictive bias and Disparate impact,” *journal-article*
14. Su, L., Zhao, Y., Yan, T.: Two-stage method based on local polynomial fitting for a linear heteroscedastic regression model and its application in economics. *Discrete Dynamics in Nature and Society* **2012**, 1–17 (January 2012). <https://doi.org/10.1155/2012/696927>
15. Wang, A., Ramaswamy, V.V., Russakovsky, O.: Towards intersectionality in machine learning: Including more identities, handling underrepresentation, and performing evaluation. In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. pp. 336–349 (2022)