

Technische Universität Dortmund

Fakultät Statistik

Wintersemester 2022/2023

Fallstudien I

# **Projekt 2: Multiple Lineare Regression**

Prof. Dr. Guido Knapp

M. Sc. Yassine Talleb

Bericht von: Louisa Poggel

Mitglieder der Gruppe 1:

Caroline Baer

Daniel Sipek

Julia Keiter

Louisa Poggel

17.11.2022

# Inhaltsverzeichnis

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Einleitung</b>   | <b>1</b>  |
| <b>2</b> | <b>Problemstellung</b>  | <b>1</b>  |
| <b>3</b> | <b>Statistische Methoden</b>                                    | <b>3</b>  |
| 3.1      | Modellbildung und Variablenselektion . . . . .                  | 3         |
| 3.2      | Modelldiagnostik . . . . .                                      | 6         |
| <b>4</b> | <b>Statistische Auswertung</b>                                  | <b>7</b>  |
| 4.1      | Deskriptive Beschreibung des Datensatzes . . . . .              | 7         |
| 4.2      | Modellbildung und Variablenselektion . . . . .                  | 9         |
| 4.3      | Modelldiagnostik . . . . .                                      | 11        |
| 4.4      | Interpretation der Koeffizienten des Parametervektors . . . . . | 13        |
| <b>5</b> | <b>Zusammenfassung</b>  | <b>14</b> |
|          | <b>Literaturverzeichnis</b>                                     | <b>16</b> |
|          | <b>Anhang</b>   | <b>17</b> |

# 1 Einleitung

Dieses Projekt beschäftigt sich mit der Ermittlung der ortsüblichen Vergleichsmiete auf Basis des Mietspiegels der Stadt München aus dem Jahr 2015. Unter Berücksichtigung von gesetzlichen Vorgaben wird dazu eine repräsentative Zufallsstichprobe aus allen Mietobjekten Münchens gezogen. Ziel ist es ein multiples Regressionsmodell zu erstellen, dass die *Nettomiete pro Monat* möglichst gut anhand zahlreicher Regressoren bezüglich der Ausstattung und Wohnlage der Mietobjekte beschreibt. Dabei soll das Modell bei der Findung einer korrekten Miete für zukünftige Immobilien helfen und ein neues Mietobjekt richtig einordnen können.

Es wird sich herausstellen, dass die *Nettomiete pro Monat* gut durch einen linearen Zusammenhang ohne polynomiale Koeffizienten beschrieben wird. Dabei werden sich die Variablen *Wohnfläche*, *Ausstattung Küche*, *Wohnlage*, *Baujahr*, *Warmwasserversorgung*, *gefliestes Bad*, *Zentralheizung* und *Ausstattung Bad* als geeignete Regressoren herausstellen. Im folgendem Kapitel 2 wird zunächst die Problemstellung, inklusive aller betrachteten Variablen, genauer erläutert. Darauf folgt in Kapitel 3 eine Darstellung der Statistischen Methoden unterteilt in die Modellbildung und Variablenselektion (3.1), sowie die Modelldiagnose (3.2). Im Kapitel 4, der statistischen Auswertung, wird nach kurzer deskriptiver Auswertung der Variablen (4.1) das Modell gebildet (4.2) und anschließend diagnostiziert und verbessert (4.3). Zuletzt erfolgt das Unterkapitel 4.4 zur Interpretation der Koeffizienten des Parametervektors, bevor alle zentralen Ergebnisse in Kapitel 5 zusammengefasst werden.

## 2 Problemstellung

Um eine sachliche Entscheidung über die Festlegung der Miete für eine bestimmte Immobilie zu erleichtern spielt die Betrachtung der Vergleichsmiete eine entscheidende Rolle. Diese wird aus einem in zahlreichen Städten erstellten Mietspiegel gewonnen und steht dabei unter Beachtung gesetzlicher Definitionen. In diesen wird eine feste Grundgesamtheit aus betrachteten Mieten festgelegt aus der eine repräsentative Zufallsstichprobe gezogen werden soll. Diese schließt beispielsweise gesetzlich festgelegte oder geförderte Mieten aus. Zudem resultiert aus dem Gesetzestext (BGB §558), dass die durchschnittli-

che Nettomiete der Regressand ist. Dieser soll durch mehrere Regressoren erklärt werden, welche sowohl die Wohnfläche, Ausstattung und Wohnlage einschließlich der Energieversorgung der Wohnung beachten. In Rahmen dieses Projektes wird die Vergleichsmiete im Raum München, wie obig beschrieben, anhand von Daten eines Ausschnittes des Mietspiegels aus dem Jahr 2015, mithilfe eines multiplen linearen Regressionsmodell geschätzt. Dabei steht vor allem die Einordnung bzw. Prognose von neuen Beobachtungen im Vordergrund. Der vorliegende Datensatz *mietspiegel2015* besteht dabei aus 13 Variablen und 3065 Beobachtungen. Eine ausführliche Beschreibung des Datenerhebungsprozesses ist weiterhin auf der Internetseite der Stadt München zu finden (Landeshauptstadt München - Sozialreferat (2015)). Die detaillierten Beschreibungen der Variablen sind in Tabelle 1 zu finden. Der Datensatz beinhaltet zunächst die in diesem

Tabelle 1: Variablen des Datensatzes

| Variablenname (kurz)                                | Ausprägung/Kodierung                 |                    | Skalenniveau      |
|---|--------------------------------------|--------------------|-------------------|
| <i>Nettomiete pro Monat (nm)</i>                    | reellwertig in Euro (€)              |                    | metrisch, stetig  |
| <i>Nettomiete pro Monat und Quadratmeter (nmqm)</i> | reellwertig in Euro (€)              |                    | metrisch, stetig  |
| <i>Wohnfläche (wfl)</i>                             | ganzzahlig in Quadratmeter ( $m^2$ ) |                    | metrisch, stetig  |
| <i>Anzahl Zimmer (räume)</i>                        | ganzzahlig als Anzahl                |                    | metrisch, diskret |
| <i>Baujahr (bj)</i>                                 | reellwertig als Zeitpunkt (Jahr)     |                    | metrisch, diskret |
| <i>gute Wohnlage (wohngut)</i>                      | 1<br>gute Lage                       | 0<br>andere Lage   | nominal, dichotom |
| <i>beste Wohnlage (wohnbest)</i>                    | 1<br>beste Lage                      | 0<br>andere Lage   | nominal, dichotom |
| <i>Warmwasserversorgung (ww)</i>                    | 1<br>nein                            | 0<br>ja            | nominal, dichotom |
| <i>Zentralheizung (zh)</i>                          | 1<br>nein                            | 0<br>ja            | nominal, dichotom |
| <i>gefliestes Bad (badkach)</i>                     | 1<br>nicht gefliest                  | 0<br>gefliest      | nominal, dichotom |
| <i>Ausstattung Bad (badextra)</i>                   | 1<br>gehoben                         | 0<br>nicht gehoben | nominal, dichotom |
| <i>Ausstattung Küche (küche)</i>                    | 1<br>gehoben                         | 0<br>nicht gehoben | nominal, dichotom |
| <i>Bezirkname (bez)</i>                             | Bezirkname in München                |                    | nominal           |

Projekt als Regressand verwendete metrisch, stetige Variable *Nettomiete pro Monat*. Diese setzt sich zusammen aus den, um Zuschläge bereinigten, an den Vermieter geleisteten Mietzahlungen aus dem Stichmonat des Januar 2014. Danach erfolgte noch der Abzug

der monatlichen Betriebskostenbeträge und die Addition einer eventuell auftretenden Mietminderung (Landeshauptstadt München - Sozialreferat (2015)). Die ebenfalls als Regressand geeignete Variable *Nettomiete pro Monat und Quadratmeter* wird in diesem Projekt nicht betrachtet. Alle weiteren genannten Variablen gehören zu den möglichen Regressoren. Dazu gehört das metrisch, stetige Merkmal *Wohnfläche*, welches auf ganze Zahlen gerundet in Quadratmeter vorliegt. Das *Baujahr* der Immobilie und die *Anzahl an Zimmern* sind hingegen metrisch, diskrete Variablen. Alle weiteren Merkmale sind nominal und beschreiben sowohl die Ausstattung (bzgl. Energieversorgung, Inventar des Bads und der Küche) als auch die Wohnlage der Immobilie. Dabei ist ein Großteil der Variablen dichotom mit den Ausprägungen „0“ und „1“. Lediglich der *Bezirksname* besteht aus den 25 Bezirken der Stadt München. Für die in den Variablen *gute Wohnlage*, textitbeste Wohnlage zu findende Bewertung der Wohnlage wurde außerdem ein Gutachter hinzugezogen. Die Bezeichnungen „ja“ und „nein“ bei den Variablen *Warmwasserversorgung* und *Zentralheizung* geben jeweils an ob dies von Vermieter gestellt ist. Durch eine vorherige Einteilung des *Baujahres* in Klassen, liegen durch die Klassenauflösung einige Werte reelwertig vor (beispielsweise 1957.5). Die Qualität der vorliegenden Daten ist sehr gut, da in keiner der Variablen fehlende Werte vorliegen.

## 3 Statistische Methoden

Alle folgenden statistischen Methoden werden in der Version 4.1.1 der Software R durchgeführt (R Core Team (2021)). Dabei wird bei Ergebnissen, wenn nicht anderes angegeben, auf zwei Nachkommastellen gerundet.

### 3.1 Modellbildung und Variablenselektion

Elementar für dieses Projekt ist das klassische allgemeine lineare Modell (Fahrmeir et al. (2009), S.62). Dieses besteht aus der Designmatrix  $X \in \mathbb{R}^{n \times k}$ , dem Parametervektor  $\beta \in \mathbb{R}^k$  und den Zufallsvektoren  $y, e \in \mathbb{R}^n$ . Dabei ist  $y$  der Vektor der Beobachtungen und  $e$  ein unbeobachtbarer Vektor der Fehler. Die Dimensionen sind durch  $n, k \in \mathbb{N}$  gegeben. Nun ergibt sich die Modellgleichung als  $y = X\beta + e$ . Zudem soll gelten, dass der Erwartungs-

wert und die Kovarianz von  $y$  existieren. Für den Erwartungswert gilt  $\mathbb{E}(y) = X\beta$  und somit  $\mathbb{E}(e) = 0$ . Zudem sollen die Fehler  $e_i$  mit homoskedastischer Varianz normalverteilt sein, d.h.  $e_i \sim N(0, \sigma^2)$ . Außerdem ist eine Unkorreliertheit der Fehler untereinander erwünscht, sodass  $Cov(e) = \sigma^2 I$  gilt. Hieraus ergibt sich für  $i = 1, \dots, n$  sowohl das in (1) definierte multiple lineare Regressionsmodell (Fahrmeir et al. (2009), S. 24) als auch die in (2) definierte polynomiale Regression (Fahrmeir et al. (2009), S. 153). Dabei werden Polynome vom Grad eins bis  $l \in \mathbb{N}$  angenommen.

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + e_i \quad (1)$$

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_n x_i^l + e_i \quad (2)$$

Eine Methode zur Bestimmung der Schätzung des Parametervektors  $\beta$  ist die Kleinste Quadrate Schätzung (KQ-Schätzer). Dieser Schätzer ist definiert als die Lösung des folgenden Minimierungsproblem, welches sich bei vollem Spaltenrang wie in (3) berechnen lässt (Fahrmeir et al. (2009), S. 90 bis 92).

$$\min_{\beta \in \mathbb{R}^k} \|y - X\beta\| \quad \text{d.h. falls } rg(X) = k \text{ ist } \hat{\beta} = (X^T X)^{-1} X^T y \quad (3)$$

Falls schwache Multikollinearität vorliegt, sind mindestens zwei Spalten der Designmatrix fast linear abhängig, was zu einer Erhöhung der Varianz des Schätzers  $\hat{\beta}$  und somit zu einer Unzuverlässigkeit des KQ-Schätzer führt (Fahrmeir et al. (2009), S. 102). Wie Multikollinearität diagnostiziert werden kann wird in Kapitel 3.2 beschrieben. In diesem Fall wird auf die Lasso (Least absolute shrinkage and selection operator) Schätzung zurückgegriffen. Diese ist im Falle eines multiplen Regressionsmodell wie in (4) definiert (James et al. (2021) S. 241). Analog kann der Schätzer auch für die polynomiale Regression verwendet werden.

$$\min_{\beta \in \mathbb{R}^k} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^k \beta_j x_{ij} \right)^2 + v \sum_{j=1}^k |\beta_j| \quad v \in \mathbb{R}, v > 0 \quad (4)$$

Der letzte Summand ist ein Strafterm, der die Koeffizienten von  $\hat{\beta}$  schrumpfen lässt. Dabei wird das  $v$  durch eine Kreuzvalidierung unter Verwendung der Funktion `cv.glmnet()` aus dem Paket `glmnet` (Friedman et al. (2010)) passend gewählt.

Zur Variablenselektion werden zweiseitige t-Tests verwendet. Diese testen, wie in (5) beschrieben, ob ein geschätzter Koeffizient des Parametervektors  $\hat{\beta}$  signifikant von Null verschieden ist. Dies geschieht mit dem wie in (6) definierten Test  $\varphi$  (Fahrmeir et al. (2009),

S. 116). Dabei wird mit  $n_e = n - p$  der Fehlerfreiheitsgrad mit der Anzahl von Regressoren  $p \in \mathbb{N}$  bezeichnet.

$$H_0 : \hat{\beta}_i = 0 \text{ vs. } H_1 : \hat{\beta}_i \neq 0 \quad T = \frac{\hat{\beta}_i}{s_{\hat{\beta}_i}} \text{ mit } s_{\hat{\beta}_i} = \sqrt{\widehat{Var}(\hat{\beta}_i)} \quad (5)$$

$$\varphi := \begin{cases} 0 & \text{falls } |T| \leq t_{n_e, 1-\frac{\alpha}{2}} \\ 1 & \text{falls } |T| > t_{n_e, 1-\frac{\alpha}{2}} \end{cases} \quad (6)$$

Insbesondere ist der p-Wert als eine Überschreitungswahrscheinlichkeit des Tests  $\varphi$  von Interesse. Falls dieser kleiner als  $\alpha$  ist, lässt sich die Nullhypothese unter Einhaltung des Signifikanzniveaus  $\alpha$  ablehnen.

Bei Selektionsverfahren durch den p-Wert (James et al. (2021) S. 79) wird vorher ein kritischer Wert für  $\alpha$  festgelegt. Bei der Rückwärtselimination werden, bei Start des vollen Modells, schrittweise alle Variablen eliminiert, die einen p-Wert kleiner als den kritischen Wert haben. Bei der Vorwärtsselektion wird dieser Prozess umgedreht und es werden schrittweise die Variablen mit dem kleinsten p-Wert aufgenommen. Bei einer schrittweisen Selektion werden beide Verfahren kombiniert.

Ein weiteres Selektionskriterium zur Bewertung der Anpassung eines Modells ist das adjustierte Bestimmtheitsmaß  $R^2_{adj}$  (Fahrmeir et al. (2009), S. 160). Dieses setzt sich aus dem Bestimmtheitsmaß  $R^2$  in (7) zusammen, welches den Anteil der durch das Modell erklärten Streuung an der Gesamtstreuung angibt (Fahrmeir et al. (2009), S. 99).

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \in [0, 1] \quad (7)$$

Werte nahe eins sprechen für eine gute Anpassung des Modelles und Werte nahe Null für eine schlechte Anpassung. Da  $R^2$  aber auch bei Hinzunahme von unwichtigen Regressoren steigt, ist das in (8) definierte adjustierte Bestimmtheitsmaß besser zur Variablenselektion geeignet.

$$\tilde{R}^2 = 1 - \frac{n-1}{n-p}(1 - R^2) \leq 1 \quad (8)$$

Denn dieses bestraft die Hinzunahme von zusätzlichen Variablen mit einem Strafterm, der sich aus der Anzahl Beobachtungen ( $n$ ) und  $p \in \mathbb{N}$ , die Anzahl der Regressoren, zusammensetzt. Weiterhin ist das Maß nach oben beschränkt, sodass Werte nahe 1 weiterhin für eine gute Modellanpassung sprechen.

### 3.2 Modelldiagnostik

Zur Überprüfung der Modellannahmen werden die gewöhnlichen Residuen  $\hat{e}_i = y_i - \hat{y}_i$  betrachtet, die in einigen Anwendungen mithilfe des i-ten Diagonalelementes  $h_{ii}$  der Hat Matrix  $H = X(X^T X)^{-1} X^T$  und dem MSE (Mean Sum of Squares) wie in (9) standardisiert werden (Fahrmeir et al. (2009), S. 110).

$$e_i^* = \frac{\hat{e}_i}{\sqrt{MSE \cdot \sqrt{1 - h_{ii}}}} \quad MSE = \frac{\hat{e}^T \hat{e}}{n_e} \quad (9)$$

Bei einem einfachen Residualplot werden die Residuen  $\hat{e}$  gegen die angepassten Werte  $\hat{y} = X\hat{\beta}$  abgetragen (James et al. (2021), S. 93 bis 94). Außerdem wird ein „Scale-plot“ mit analogem Prinzip des Residualplots unter Verwendung der Transformation der Residuen  $\tilde{e}_i = \sqrt{|e_i^*|}$  genutzt. Diese Definition entspricht der in R implementierten Version aus der Funktion `plot.lm` unter Angabe des Arguments `which = 3` (R Core Team (2021)). Hier lässt sich, aufgrund der Standardisierung, gut Heteroskedastizität erkennen. Denn beim gewöhnlichen Residualplot kann der Effekt auftreten, dass mit größer werdendem  $\hat{y}$  die Residuen größer werden, obwohl die Varianz der Residuen nicht größer wird.

Zur Überprüfung der Normalverteilungsannahme der Fehler wird ein Quantile-Quantile-Plot benutzt (Hartung et al. (2009), S.847). Auch hier wird die in R implementierte Version der Funktion `plot.lm` verwendet, in der die empirischen Quantile der standardisierten Residuen  $e^*$  gegen die theoretischen Quantile der Normalverteilung abgetragen werden (R Core Team (2021)). Liegt ein Großteil der Punkte auf der Winkelhalbierenden, spricht dies für die Erfüllung der Normalverteilungsannahme.

Der Leverage (deutsch: Hebel, Einfluss) ist ein Abstandsmaß, dass bezüglich einer unabhängigen Variable den Abstand einer Beobachtung  $i$  zu den übrigen Beobachtungen angibt (Fahrmeir et al. (2009), S. 177 bis 178). Das Maß ist für die  $i$ -te Beobachtung ist definiert als das, zwischen 0 und 1 liegende,  $i$ -te Diagonalelement  $h_{ii}$  der Hatmatrix. Der durchschnittliche Leverage beträgt  $p/n$ . Beobachtungen mit einem Leverage score größer als  $2 \cdot (p/n)$  werden als Beobachtungen mit großer Hebelwirkung auf die Regressionsgerade bezeichnet. Haben diese zudem große Residuen, spricht man von einflussreichen Beobachtungen. Zu dessen Identifikation dient die Cook's Distance. Wie in (10) wird die Summe der quadrierten Änderungen, wenn die  $i$ -te Beobachtung entfernt wird, berechnet (Fahrmeir et al. (2009), S. 178). Eine Beobachtungen ab einem Wert von  $D_i > 0.5$



auffällig und sollte ab dem Wert  $D_i > 1$  auf jeden Fall näher untersucht werden .

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2}{p \cdot \text{MSE}} \quad (10)$$

Um zu ergründen ob Multikollinearität ein Problem darstellt, wird die Determinante von  $X^T X$  berechnet. Denn diese ist im Falle schwacher Multikollinearität nahezu Null (Toutenburg (2003), S. 114). Zudem eignet sich der Varianzinflationskoeffizient (VIF), definiert als  $VIF_i = 1/(1-R_i^2)$  als ein Indikator für Multikollinearität (Fahrmeir et al. (2009), S. 170 bis 171). Dabei bezeichnet  $R_i^2$  den multiplen Korrelationskoeffizient bei einer Regression wo  $x_i$  als abhängige Variable und alle weiteren Prädiktoren als unabhängige Variablen gesehen werden. Ein Wert des VIF größer als 10 spricht dabei für Multikollinearität. Die Umsetzung in R erfolgt über die Funktion `vif()` aus dem Paket `car` (Fox und Weisberg (2019)).

## 4 Statistische Auswertung

### 4.1 Deskriptive Beschreibung des Datensatzes

Der Tabelle 2 sind die deskriptiven Kennzahlen aller metrischen Variablen zu entnehmen. Interessant zu betrachten ist der zukünftige Regressand *Nettomiete pro Monat (nm)*, der eine rechtsschiefe, spitze Verteilung mit einer deutlich größeren Standardabweichung als die *Nettomiete pro Monat und Quadratmeter(nmqm)* aufweist. Somit ähnelt die Vertei-

Tabelle 2: Deskriptive Kennzahlen der metrischen Variablen

|              | arithm. Mittel | Median  | Standardabweichung | IQR    | Schiefe | Wölbung |
|--------------|----------------|---------|--------------------|--------|---------|---------|
| <i>nm</i>    | 763.06         | 700.00  | 338.16             | 360.46 | 2.59    | 25.47   |
| <i>nmqm</i>  | 10.73          | 10.84   | 2.67               | 3.42   | 0.04    | 3.34    |
| <i>wfl</i>   | 71.98          | 70.00   | 25.74              | 30.00  | 1.35    | 8.33    |
| <i>räume</i> | 2.70           | 3.00    | 0.98               | 1.00   | 0.46    | 3.60    |
| <i>bj</i>    | 1964.21        | 1957.50 | 26.51              | 25.50  | -0.18   | 2.31    |

lung der *Nettomiete pro Quadratmeter* nicht wirklich einer Normalverteilung. Dies lässt sich auch gut grafisch in Abbildung 6, die im Anhang auf Seite 18 erkennen. Dabei wird

zur besseren Sichtbarkeit der Verteilung die Beobachtung 1975 mit einer *Nettomiete pro Monat* von 6000 Euro nicht dargestellt. Zur Berechnung der Schiefe- und Wölbungsmaße wurde das **moments** Paket (Komsta und Novomestky (2022)) verwendet.

In Tabelle 3 sind die relativen Häufigkeiten der dichotomen Variablen vorzufinden. Hier ist auffällig, dass bei den Merkmalen *beste Wohnlage*, *Warmwasserversorgung* und *Zentralheizung* die Ausprägung „1“ nur sehr selten vorkommt. Eine hier nicht aufgeführte,

Tabelle 3: Relative Häufigkeiten der dichotomen Variablen (n = 3065)

|                | <i>wohngut</i> | <i>wohnbest</i> | <i>ww</i> | <i>zh</i> | <i>badkach</i> | <i>badextra</i> | <i>kueche</i> |
|----------------|----------------|-----------------|-----------|-----------|----------------|-----------------|---------------|
| Ausprägung „0“ | 0.65           | 0.96            | 0.99      | 0.93      | 0.12           | 0.88            | 0.75          |
| Ausprägung „1“ | 0.35           | 0.04            | 0.01      | 0.07      | 0.88           | 0.12            | 0.25          |

vollständige Übersicht über die relativen Häufigkeiten der nominalen Variable *Bezirk* ist im Anhang auf Seite 17 in Tabelle 5 zu finden.

Die in Abbildung 1 dargestellten Korrelationen der metrischen Variablen, wurden mit dem Rangkorrelationskoeffizienten nach Spearman berechnet, um einen Vergleich zwischen diskreten und stetigen Variablen möglich zu machen. Dabei weist die Variable

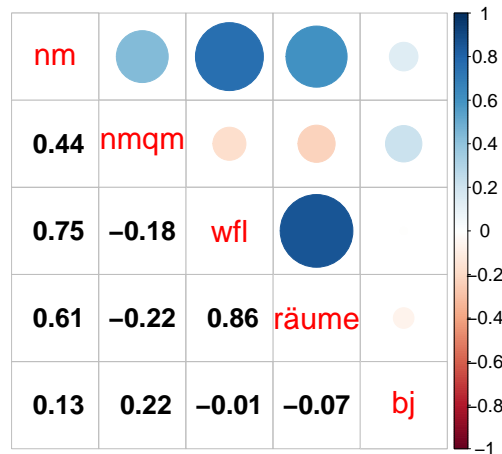


Abbildung 1: Korrelationen zwischen den metrischen Variablen

*Wohnfläche* mit 0.75 die höchste Korrelation mit der *Nettomiete pro Monat* auf, gefolgt vom Merkmal *Anzahl Zimmer* mit einer moderaten Korrelation von etwa 0.61. Auch die *Nettomiete pro Monat und Quadratmeter* ist mit 0.44 leicht mit dem Regressand korreliert. Dies sollte nicht weiter verwunderlich sein, das die *Nettomiete pro Monat und Quadratmeter* die *Nettomiete pro Monat* enthält. Aufgrund dieser Dopplung von Informationen wird die *Nettomiete pro Monat und Quadratmeter* als möglicher Regressor

ausgeschlossen. Die höchste Korrelation besteht mit 0.86 zwischen den Variablen *Wohnfläche* und *Anzahl der Zimmer*. Dies könnte ein Hinweis auf Multikollinearität sein, der später überprüft wird. Ansonsten liegen kaum nennenswerte sehr leichte negative und positive Korrelationen zwischen den Variablen vor.

## 4.2 Modellbildung und Variablenselektion

Zunächst werden vor der Modellbildung die beiden Variablen *beste Wohnlage* und *gute Wohnlage* in einem neuen Merkmal *Wohnlage* mit den drei Ausprägungen „beste“, „gute“ oder „andere“ Lagekategorie zusammengefasst.

Vielversprechend als erklärende Variable ist aufgrund der hohen Korrelation mit dem Regressand die Variable *Wohnfläche*. Führt man explorativ eine einfache lineare Regression mit der *Wohnfläche* aus, erhält man schon ein recht hohes adjustiertes Bestimmtheitsmaß von etwa  $R_{adj}^2 = 0.61$  und einen p-Wert für die *Wohnfläche*, der kleiner als  $2 \cdot 10^{-16}$  ist. Das heißt der Koeffizient  $\beta_1$  bezüglich der *Wohnfläche* ist in diesem Modell signifikant von Null verschieden. Auch der in Abbildung 2 erkennbare Verlauf der Regressionsgerade stimmt größtenteils mit dem Verlauf der Punktwolke überein. Aufgrund dem leicht

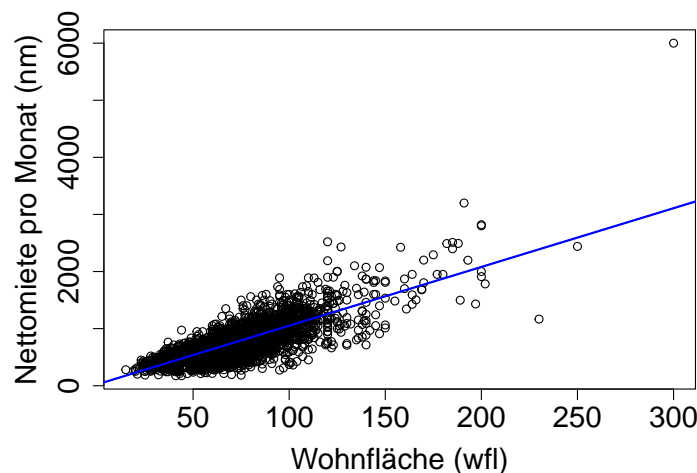


Abbildung 2: Einfache Regression durch die *Wohnfläche*

gekrümmten Verlauf der Punktwolke könnte auch eine polynomialer Ansatz sinnvoll sein (vgl. Abbildung 7 auf Seite 18). Jedoch ergeben sich bei einer polynomialen Regression

zweiten Grades, welche die *Nettomiete pro Monat* durch die *Wohnfläche* und die quadrierte *Wohnfläche* erklärt, Probleme mit der Multikollinearität (vgl. VIF von  $wfl \approx 17.70$ , VIF von  $wfl^2 \approx 11.27$ ). Bei einer Vorwärtsselektion mit einem kritischen Wert von 0.05 bleiben die beiden Terme der polynomialen Regression zwar enthalten. Wendet man aber aufgrund der Multikollinearität, nach einer Kreuzvalidierung zur Bestimmung des Parameters  $v \approx 0.43$  eine Lasso Regression an, wird der quadratische Term mit etwa 0.02 nahe Null geschätzt. Auch der KQ-Schätzer, der dem Lasso Schätzer aufgrund des kleinen  $v$  recht ähnlich ist schätzt den Einfluss von der quadrierten *Wohnfläche* als nur sehr gering ein. Somit wird dieser Ansatz im folgenden nicht weiter verfolgt.

Hingegen wird eine Vorwärtsselektion und Rückwärtselimination ohne quadratische Terme auf alle möglichen Regressoren angewendet. Mit dem vorher festgelegtem kritischem Wert für den p-Wert von  $\alpha = 0.05$ , werden bei der Vorwärtsselektion der Reihenfolge nach die Variablen *Wohnfläche*, *Ausstattung Küche*, *Wohnlage*, *Baujahr*, *Anzahl Zimmer*, *Warmwasserversorgung*, *gefliestes Bad*, *Zentralheizung* und *Ausstattung Bad* hinzugefügt. Zuletzt gilt zu entscheiden ob der *Bezirkname* aufgenommen werden soll. Da jedoch bei Betrachtung der Signifikanz der einzelnen Dummy-Variablen des Merkmales nur 5 von 25 Stadtteilen den kritischen Wert von 0.05 einhalten, wird auf die Hinzunahme des *Bezirks* zur Komplexitätsreduzierung des Modelles verzichtet.

Das resultierende Modell weist mit etwa 0.69 ein höheres adjustiertes Bestimmtheitsmaß als das einfache lineare Modell aus Abbildung 2 auf. Das heißt die hinzugenommenen Variablen tragen zur Erklärung des Regressanden bei und somit wurde das Modell der einfachen Regression verbessert. Bei der Rückwärtselimination erfolgt ebenfalls die Auswahl dieses Modells, da alle p-Werte im vollen Modell unter  $\alpha = 0.05$  liegen.

Noch näher zu betrachten ist die kritisch zu sehende Aufnahme der Variable *Anzahl der Zimmer*. Denn die in Abbildung 1 erkennbare hohe Korrelation mit der Variable *Wohnfläche* könnte zu Problemen führen. Zunächst wird untersucht, ob ein Problem durch Multikollinearität vorliegt. Die Determinante von  $X^T X$  liegt mit etwa  $1.67 \cdot 10^{35}$  sehr weit weg von der Null. Auch der VIF von etwa 3.56 der Variable *Wohnfläche* und der Variable *Anzahl Zimmer* ( $\approx 3.51$ ) ist nur etwas höher als bei den anderen Variablen. Erst ein VIF von 10 wird als kritisch in Bezug auf die Multikollinearität gesehen. Jedoch wird der Koeffizient der *Anzahl der Zimmer* mit einem fragwürdigem, negativem Wert von  $-55.54$  geschätzt. Das würde bedeuten, dass mehr Zimmer zu einer günstigeren Miete führen würde. Aufgrund dieses Ergebnisses, dass aus einer Wechselwirkung mit der *Wohnfläche* resultieren könnte, wird die Variable *Anzahl der Zimmer* entfernt.

Das nun resultierende Modell hat ein nur geringfügig kleineres adjustiertes Bestimm-

heitsmaß von etwa 0.68 und alle aufgenommenen Variablen haben einen p-Wert, der unter dem kritischen Wert von 0.05 liegt. Jegliche Werte des VIF liegen unter 1.2 und die Determinante von  $X^T X$  liegt mit etwa  $2 \cdot 10^{32}$  weit weg von der Null. Somit liegt keine Multikollinearität vor, sodass der KQ-Schätzer für die Schätzung von  $\hat{\beta}$  verwendet werden kann.

### 4.3 Modelldiagnostik

Zunächst wird untersucht, ob das Modell einflussreiche Beobachtungen enthält. In Abbildung 3 sind dazu die standardisierten Residuen gegen den Leverage abgetragen. Die orange vertikale Linie gibt den cut-off Wert von  $2 \cdot (8/3065) \approx 0.005$  an. Es ist erkennbar, dass nur ein standardisiertes Residuum der Beobachtung 1975 eine Cook's Distance größer als 0.5 und einen Leverage größer als 0.005 hat. Diese Beobachtung ist also einflussreich und hat eine große Hebelwirkung. Auch die Beobachtungen 1263 und 231 sind auffällig, da diese einen hohen Leverage und eine tendenziell höhere Cook's Distance als die anderen Werte aufweisen. Generell ist zu erkennen, dass viele Residuen einen hohen Leverage haben. Dies könnte daran liegen, dass für große  $\hat{y}$ -Werte nur wenig Beobachtungen vorliegen, welche nun einen großen Einfluss auf die Regression ausüben.

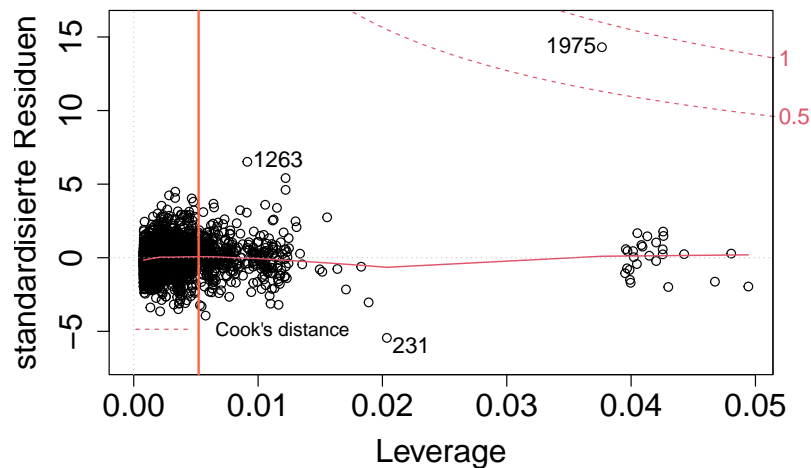


Abbildung 3: Leverage und Cook's Distance

Die obig genannten auffälligen Beobachtungen sind auch in der Abbildung 8 auf Seite 19 im Anhang erkennbar. Diese zeigt einen QQ-Plot, bei dem zu erkennen ist, dass im mittleren Teil ein Großteil der Punkte auf der Verbindungslinie zwischen dem ersten und dritten Quartil liegt, dessen Lage einer Winkelhalbierenden (in rot) gleicht. An den Rändern gibt es jedoch einige Punkte, die deutlich nach unten oder oben abweichen. Somit werden die Beobachtungen 1975, 1263 und 231 probenhalber entfernt, was zumindest zu einer Verbesserung der Ränder des Quantile-Quantile-Plot führt und somit der Erfüllung der Normalverteilungsannahme näher kommt. Zudem liegen die Residuen im Mittel bei einem Median von 2.11 vergleichsweise näher an Null als im vorherigen Modell mit einem Median von 3.34. Somit wird auch die Modellannahme, dass der  $\mathbb{E}(e) = 0$  ist, besser erfüllt. Dies ist auch gut in Abbildung 4 an der roten Durchschnittslinie erkennbar, die bis auf eine leichte anfängliche Schwankung bei Null liegt. Ansonsten sind keine Strukturen in den Residuen erkennbar, die auf autokorrelierte Fehler hinweisen. Jedoch stammt ein Großteil der Residuen aus dem Wertebereich zwischen 400 und 1000 der angepassten y-Werte und streut dort etwas weniger um die Null herum als in Wertebereich zwischen 1000 und 2500.

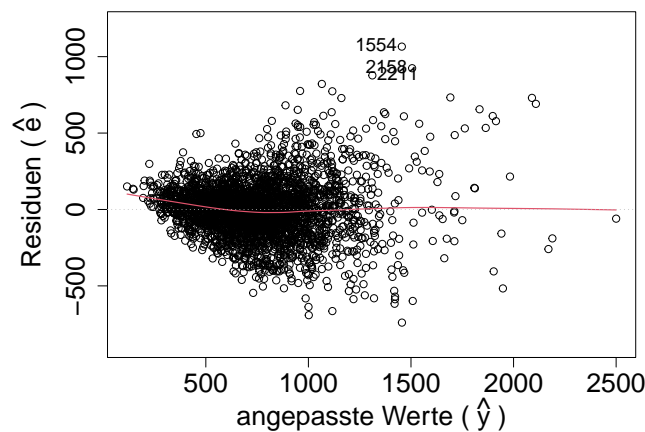


Abbildung 4: Residualplot ohne Beobachtungen 1975, 1263, 231

Um besser bewerten zu können, ob die in Abbildung 4 zu beobachtende größer werdende Streuung der Residuen an einer Verletzung der Homoskedastizität liegt, wird in Abbildung 5 ein Scale-plot unter Verwendung von standardisierten Residuen betrachtet. Dort ist zu erkennen, dass die rote Durchschnittsinie zunächst konstant bleibt und dann von etwa 0.5 auf 1.8 ansteigt. Auch die Form der Punktwolke ähnelt annähernd einer

Ellipse, dessen Achse der Ausrichtung der roten Linie entspricht. Beide Beobachtungen sprechen für Heteroskedastizität, da die Residuen  $\tilde{e}$  mit wachsenden  $\hat{y}$  immer größer werden.

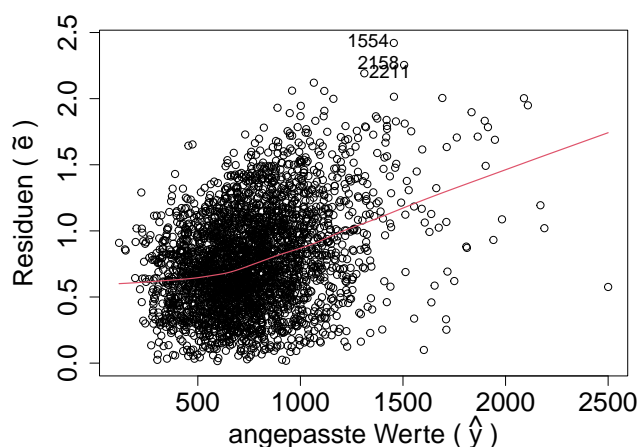


Abbildung 5: Scale-plot ohne Beobachtungen 1975, 1263, 231

Zusammenfassend lässt sich mit dem Modell also die Modellannahme, dass  $\mathbb{E}(y) = X\beta$  ist durch die Herausnahme der Beobachtungen 1975, 1263 und 231 verbessern und als erfüllt ansehen. Wie am Quantile-Quantile-Plot erkennbar ist, kann man auch die Normalverteilungsannahme der Fehler, unter Beobachtung von Abweichungen an den Rändern, akzeptieren. Kritischer zu sehen ist das Vorliegen der Heteroskedastizität, die schon deutlich erkennbar ist. Hinweise auf einen Verstoß der Unkorreliertheit der Fehler liegen nicht vor. Da das Eliminieren der Beobachtungen 1975, 1263 und 231 die Erfüllung Modellannahmen tendenziell zu begünstigen scheint, ist somit auch ein Verlust des Bestimmtheitsmaßes in dritter Nachkommastelle von etwa 0.679 auf etwa 0.677 zu rechtfertigen. Das Entfernen verbessert auch aus inhaltlicher Sicht die Prognoseeigenschaft des Modelles, da untypische Beobachtungen nicht die Einordnung neuer Mietobjekte verzerren.

#### 4.4 Interpretation der Koeffizienten des Parametervektors

Die Ergebnisse des nun resultierenden KQ-Schätzers sind in Tabelle 4 zu finden. Die zugehörigen p-Werte des Signifikanztests sind in Tabelle 6 auf Seite 17 im Anhang zu

finden. Der Einfluss der (Ausprägungen) der Variablen *Wohnfläche*, „gehobene“ *Ausstattung der Küche*, „beste“ *Wohnlage*, „gute“ *Wohnlage*, *Baujahr*, „nicht vorhanden sein“ eines *gefliesen Bades* und „gehobene“ *Ausstattung des Bades* wird positiv auf die Miete eingeschätzt. Das heißt in Bezug auf die numerischen Variablen *Wohnfläche* und *Bau-*

Tabelle 4: Geschätzte Koeffizienten von  $\hat{\beta}$

|                         |           |                                  |          |
|-------------------------|-----------|----------------------------------|----------|
| Intercept               | -2196.526 | <i>bj</i>                        | 1.093    |
| <i>wfl</i>              | 9.835     | <i>ww</i> (nicht vorhanden)      | -187.504 |
| <i>kueche</i> (gehoben) | 90.756    | <i>badkach</i> (nicht vorhanden) | 54.675   |
| <i>wohnlage</i> (beste) | 111.141   | <i>zh</i> (nicht vorhanden)      | -62.036  |
| <i>wohnlage</i> (gute)  | 87.188    | <i>badextra</i> (gehoben)        | 37.681   |

*jahr*, dass wenn alle anderen Variablen konstant im Modell vorliegen und die jeweilige Variable um eine Einheit steigt, dass die *Nettomiete pro Monat* um etwa 9.835 bzw. 1.093 Euro steigt. Bei den dichotomen Variablen steigt die *Nettomiete pro Monat* um den jeweiligen Koeffizienten im Vergleich zur Referenzkategorie. Das heißt beispielsweise, dass die *Nettomiete pro Monat* in „bester“ *Wohnlage* um etwa 111.141 Euro teurer ist als in einer „anderen“ *Wohnlage*. Mit negativem Koeffizienten werden der Intercept und das „nicht vorhanden sein“ einer *Warmwasserversorgung* und *Zentralheizung* geschätzt. Dabei ist der Intercept nicht sinnvoll interpretierbar, da die *Wohnfläche* nicht Null werden kann. Die anderen beiden negativen Einflussvariablen sind analog zu den positiven, aber nun in Form einer Mietverringerung zu interpretieren.

## 5 Zusammenfassung

Zur Schätzung der *Nettomiete pro Monat* als Vergleichsmiete, anhand eines Ausschnittes des Münchener Mietspiegels aus dem Jahr 2015, wurde eine möglichst gutes multiples lineares Regressionsmodell gesucht. Dieses soll die *Nettomiete pro Monat* mithilfe verschiedenster Regressoren, welche die Ausstattung, Energieversorgung und Wohnlage der Mietobjekte beschreiben, erklären.

Vor der Modellsuche wurde zunächst die Variable *Wohnfläche*, aufgrund der hohen Korrelation von 0.75 mit der Zielvariable, als vielversprechender Regressor ausgemacht. Ansätze einer polynomialen Regression mit der *Wohnfläche* und weiteren zusätzlichen Regressoren erwiesen sich dabei als nicht sinnvoll.



Stadtdessen wurde der Ansatz eines multiplen linearen Regressionsmodelles verfolgt. Als Konsequenz einer Vorwärtsselektion wird dabei die Variable *Bezirk* eliminiert, sodass das Modell die Variablen *Wohnfläche*, *Ausstattung Küche*, *Wohnlage*, *Baujahr*, *Anzahl Zimmer*, *Warmwasserversorgung*, *gefliestes Bad*, *Zentralheizung* und *Ausstattung Bad* enthält. Im weiteren Verlauf wurde die Variable *Anzahl der Zimmer* aufgrund logischer Überlegungen und möglichen Wechselwirkungen mit der *Wohnfläche* eliminiert. Im folgenden wurden dann die einflussreichen Beobachtungen 1975, 1263 und 231 entfernt um die Modellannahmen besser zu erfüllen und die Prognosefähigkeit des Modells zu verbessern. Eine detaillierte Darstellung der Koeffizienten von  $\hat{\beta}$  ist in Tabelle 4 zu finden. Alle im Modell enthaltenen Variablen, bis auf das „nicht vorhanden sein“ von *Warmwasserversorgung* oder *Zentralheizung*, haben dabei einen mieterhöhenden Einfluss. Leider stellte sich, anhand der Residuen, eine deutlich erkennbare Verletzung der Homoskedastizität der Varianzen der Fehler heraus, was dazu führt dass die Schätzung zwar erwartungstreu bleiben aber nicht mehr effizient sind (Auer und Rottmann (2010), S.518 bis 520). Das heißt es gibt unter Umständen einen besseren Schätzer für  $\hat{\beta}$ . Beispielsweise könnte man eine gewichtete KQ-Schätzung mit vorher geschätzten Gewichten  $\hat{w}_i = 1/\hat{\sigma}_i^2$  durchführen, die Beobachtungen mit größerer Streuung weniger gewichtet. Ein weiterer Effekt der Homoskedastizität ist, dass die Schätzung des Standardfehlers verzerrt ist, was zu verfälschten Testentscheidungen im t-Test geführt haben könnte (Auer und Rottmann (2010), S.518 bis 520). In Anbetracht dieser Information könnten es auch plausibel sein, die Merkmale *Warmwasserversorgung* bzw. *Zentralheizung* trotz hoher Signifikanz zu eliminieren. Denn das „nicht vorhanden sein“ des jeweiligen Merkmals kommt so selten vor, dass es aus inhaltlicher Sicht nicht unbedingt benötigt wird. Zudem könnte es sinnvoll sein Wechselwirkungsterme in das Modell zu integrieren (z.B. zwischen *Wohnfläche* und *Anzahl der Zimmer*). Zudem ist es zunächst verwunderlich, dass ein *gefliestes Bad* mietverringert wirkt. Auch hier könnte eine Wechselwirkung vorliegen. Denn *geflieste Bäder* sind häufiger in älteren Immobilien als in Neubauten zu finden, welche tendenziell eine höhere *Nettomiete pro Monat* haben. Jedoch könnte es auch sein, dass bei Mietobjekten tatsächlich Bäder ohne Fliesen bevorzugt werden. Dies könnte durch eine Plausibilitätsprüfung eines Fachmannes (z.B. Immobilienmakler) näher ergründet werden.

# Literatur

- Auer, B. und H. Rottmann (2010). *Statistik und Ökonometrie für Wirtschaftswissenschaftler: eine anwendungsorientierte Einführung*. 1. Auflage.
- Fahrmeir, L., T. Kneib und S. Lang (2009). *Regression: Modelle, Methoden und Anwendungen*. Statistik und ihre Anwendungen. 2. Auflage. Springer Berlin Heidelberg.
- Fox, J. und S. Weisberg (2019). *An R Companion to Applied Regression*. Third. Sage: Thousand Oaks CA. URL: <https://socialsciences.mcmaster.ca/jfox/Books/Companion/>.
- Friedman, J., T. Hastie und R. Tibshirani (2010). „Regularization Paths for Generalized Linear Models via Coordinate Descent“. In: *Journal of Statistical Software* 33(1), S. 1–22. URL: <https://www.jstatsoft.org/v33/i01/>.
- Hartung, J., B. Elpelt und K.-H. Klösener (2009). *Statistik Lehr- und Handbuch der angewandten Statistik*. 15.Auflage. Oldenbourg Verlag: München.
- James, G., D. Witten, T. Hastie und R. Tibshirani (2021). *An Introduction to Statistical Learning: with Applications in R*. Second Edition. Springer.
- Komsta, L. und F. Novomestky (2022). *moments: Moments, Cumulants, Skewness, Kurtosis and Related Tests*. R package version 0.14.1. URL: <https://CRAN.R-project.org/package=moments>.
- Landeshauptstadt München - Sozialreferat (2015). *Mietspiegel für München 2015 - Dokumentation*. URL: <https://archiv.mietspiegel-muenchen.de/2015/berechnungsprogramm/dokumentation.php> (besucht am 15.11.2022).
- R Core Team (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. URL: <https://www.R-project.org/>.
- Toutenburg, H. (2003). *Lineare Modelle: Theorie und Anwendung*. 2. Auflage. Springer-Verlag Berlin Heidelberg.

## Anhang

| Bezirk                      | relative Häufigkeit |
|-----------------------------|---------------------|
| Allach-Untermenzing         | 0.01                |
| Altstadt-Lehel              | 0.02                |
| Au-Haidhausen               | 0.05                |
| Aubing                      | 0.02                |
| Berg am Laim                | 0.03                |
| Bogenhausen                 | 0.05                |
| Fledmoching-Hasenbergel     | 0.03                |
| Hadern                      | 0.03                |
| Laim                        | 0.03                |
| Ludwigvorstadt-Isarvorstadt | 0.05                |
| Maxvorstadt                 | 0.05                |
| Milbersthoften-Am Hart      | 0.04                |
| Moosach                     | 0.03                |
| Neuhausen-Nymphenburg       | 0.08                |
| Obergiesing                 | 0.05                |
| Pasing-Obermenzing          | 0.04                |
| Ramersdorf-Perlach          | 0.06                |
| Schwabing-Freimann          | 0.05                |
| Schwabing West              | 0.05                |
| Schwanthalerhoehe           | 0.03                |
| Sendling                    | 0.04                |
| Sendling-Westpark           | 0.04                |
| Thalkirchen                 | 0.06                |
| Trudering-Riem              | 0.03                |
| Untergiesing                | 0.04                |

Tabelle 5: Relative Häufigkeiten - Bezirke Münchens (n = 3065)

|                         |                       |                                  |                       |
|-------------------------|-----------------------|----------------------------------|-----------------------|
| Intercept               | $< 2 \cdot 10^{-16}$  | <i>bj</i>                        | $9.27 \cdot 10^{-16}$ |
| <i>wfl</i>              | $< 2 \cdot 10^{-16}$  | <i>ww</i> (nicht vorhanden)      | $5.30 \cdot 10^{-07}$ |
| <i>kueche</i> (gehoben) | $< 2 \cdot 10^{-16}$  | <i>badkach</i> (nicht vorhanden) | $1.28 \cdot 10^{-07}$ |
| <i>wohnlage</i> (beste) | $1.32 \cdot 10^{-09}$ | <i>zh</i> (nicht vorhanden)      | $1.23 \cdot 10^{-05}$ |
| <i>wohnlage</i> (gute)  | $< 2 \cdot 10^{-16}$  | <i>badextra</i> (gehoben)        | 0.000432              |

Tabelle 6: Koeffizienten des finalen Modells - p-Werte des t-Tests

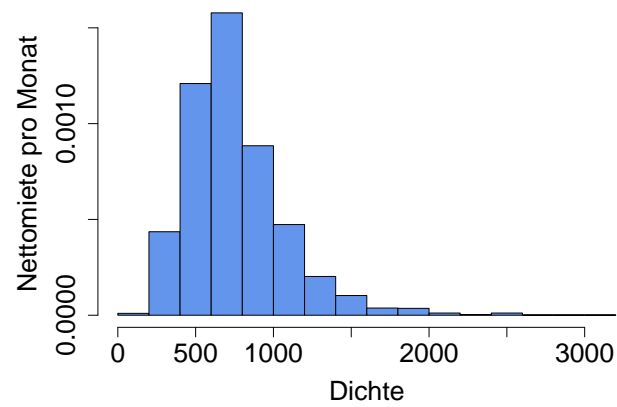


Abbildung 6: Verteilung der *Nettomiete pro Monat* ohne Beobachtung 1957

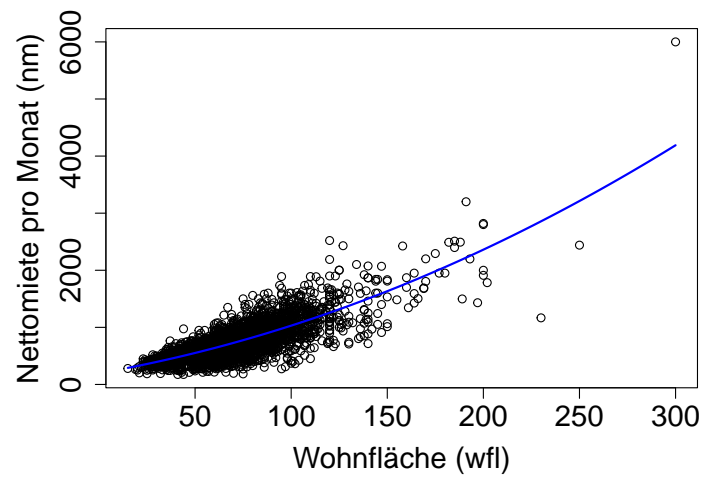


Abbildung 7: Polynomiale Regression zweiten Grades durch die *Wohnfläche*

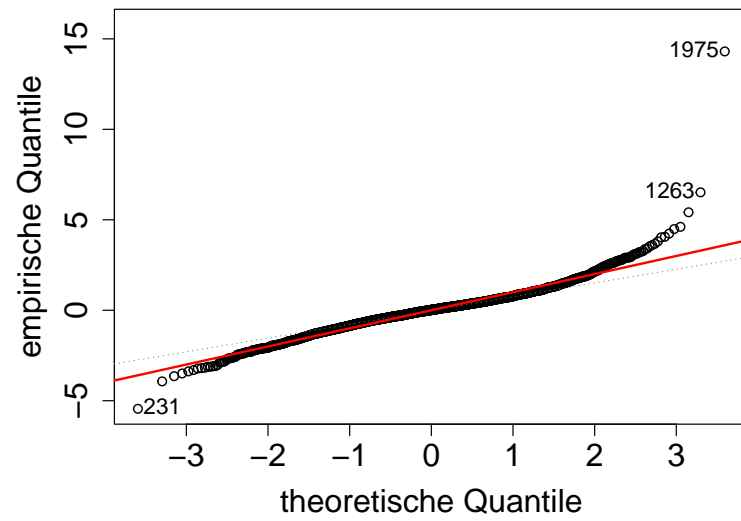


Abbildung 8: Quantile-Quantile-Plot