

Modellierung der Kaltmiete

Vergleich Frankfurt am Main und Leipzig

Henrik Popp, Kai Herbst, Manuel Zeh

2024-02-16

Inhaltsverzeichnis

Aufgabenstellung	1
Einleitung	2
Datenerhebung	3
Explorative Datenanalyse	4
Parkplatz	9
Balkon	10
Zimmer	12
Etage	14
Modellierung	14
Modellierung über Durchschnittsmietpreis	14
Modellierung über die lineare Regression	16
Andere Modellierung	18
Modellierung über Durchschnittsmietpreis	18
Modellierung über die lineare Regression	19
Zusammenfassung	20
Quellen und Hilfsmittel	20

Aufgabenstellung

Abschnitt	Aufgabe	Reiner Textumfang	Erledigt
Einleitung	Auf inhaltliche Aufgabenstellung eingehen	0,5 - 1 Seiten	[]
Datenerhebung	Wie wurden die Daten erhoben? (Suchfilter, Sortierung)	1 - 3 Sätze	[]
Explorative Datenanalyse	Analyse + eventl. Datenvorverarbeitung	1 - 2 Seiten	[]
Modellierung	Modellierung + Interpretation	1 - 2 Seiten.	[]
Zusammenfassung	Gemeinsam kurz zentrale Ergebnisse zusammenfassen + Auf Grenzen der Analyse eingehen	0,5 - 1 Seiten.	[]

- Hier auch noch Literatur recherchieren:
 - <https://de.statista.com/statistik/daten/studie/258635/umfrage/bruttokaltmiete-bewohner-wohnungen-in-deutschland-nach-bundeslaendern/>
 - https://www.deutschlandatlas.bund.de/DE/Karten/Wie-wir-wohnen/040-Mieten.html#_6a54aw429
 - <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8053893/>
 - <https://de.statista.com/statistik/daten/studie/262508/umfrage/mietpreise-in-frankfurt-am-main/>
 - <https://de.statista.com/statistik/daten/studie/1312743/umfrage/mieten-in-leipzig-nach-dem-baualter-der-wohnung/>
 - <https://de.statista.com/statistik/daten/studie/535299/umfrage/mietpreise-auf-dem-wohnungsmarkt-in-leipzig/>
 - <https://de.statista.com/statistik/daten/studie/1312730/umfrage/entwicklung-der-angebotsmieten-in-leipzig/>
 - https://link.springer.com/chapter/10.1007/978-3-658-11757-3_4
 - https://www.ifo.de/DocDL/ifoDD_14-06_03-10.pdf

Einleitung

In dieser Fallstudie sollen die Kaltmieten der beiden Städte Frankfurt am Main und Leipzig miteinander verglichen und modelliert werden. Ziel ist es, die verschiedenen möglichen Einflussfaktoren auf die Kaltmiete in den jeweiligen Städte zu bestimmen und eine Modellierung der Kaltmiete zu erstellen.

Zu Beginn wird auf die Datenerhebung eingegangen. Hier soll erklärt werden, woher die verarbeiteten Daten stammen und unter welchen Bedingungen die Daten erhoben wurden. Mit der explorativen Datenanalyse sollen dann die erhobenen Daten beschrieben und veranschaulicht werden. Hierbei wird die Vorverarbeitung der Daten beschrieben, im Anschluss wird mithilfe von Grafiken und dazugehörigen Interpretationen eine Datenanalyse

erstellt. Dabei soll unter anderem herausgefunden werden, welche erhobenen Variablen den größten Einfluss auf die Kaltmiete einer Stadt haben oder wie hoch die eventuellen Unterschiede der Mieten in den beiden Städten sind. Den zentralen Teil des Dokuments stellt die Modellierung dar. Hier soll die Kaltmiete modelliert, also durch ein selbsterstelltes statistisches Modell dargestellt werden. Zudem wird das Modell interpretiert. Zum Abschluss werden die Ergebnisse in der Zusammenfassung aufgearbeitet und präsentiert.

Datenerhebung

Die Datenerhebung fand ausschließlich über den Online-Marktplatz für Wohnungen und Häuser [ImmobilienScout24](#) statt. Die untersuchten Objekte wurden dabei auf den Immobilientyp *Wohnung* beschränkt, was als Suchkriterium in der Suchleiste des Portals eingestellt werden kann. Weitere Suchkriterien haben sich auf den *Ort*, in diesem Fall Frankfurt am Main und Leipzig, und auf den Objekttyp, hier *Mieten*, beschränkt. Weitere Kriterien wie *Anzahl der Zimmer*, *Fläche* oder einem *maximalen Preis* wurden auf den Standardeinstellungen belassen. Anschließend wurden je Ort der Reihe nach bis zu 45 Objekte in der von ImmobilienScout24 generierten Reihenfolge überprüft und in eine Excel-Datei aufgenommen, die im Folgenden als Basis für die Auswertung dienen.

Aufgenommen in die Datenbasis wurden dabei die folgenden Variablen: der *Ort*, die *Kaltmiete* in Euro, die *Wohnfläche* in Quadratmetern, das Angebot eines *Parkplatz*, die *Etage*, Anzahl der *Zimmer*, Vorhandensein eines *Balkon*, das *Baujahr* des Objektes sowie der entsprechende Link zur Anzeige und dessen Abrufdatum.

Für die nachfolgenden Auswertungen und Analysen lesen wir zunächst die Excel-Datei ein:

```
# Pfad zur Excel-Datei erstellen
pfad_mieten <- here("Mieten.xlsx")
# Daten einlesen
mieten <- read_excel(pfad_mieten)
```

Über die Ausgabe der ersten sechs Einträge erhalten wir einen Einblick in die Daten:

```
# Obere 6 Beobachtungen
head(mieten)
```

```
# A tibble: 6 x 12
```

	Ort	Kaltmiete	Wohnflaeche	Parkplatz	Etage	Zimmer	Balkon	Einbaukueche	Heizung
	<chr>	<dbl>	<dbl>	<chr>	<chr>	<dbl>	<chr>	<chr>	<chr>
1	Fran~	1800	70	ja	1	2	ja	ja	Fußbod~
2	Fran~	1500	60	ja	1	1	ja	ja	Zentra~
3	Fran~	2650	146.	ja	1	3	ja	ja	Fußbod~
4	Fran~	1500	72	nein	1	2	ja	ja	Fußbod~
5	Fran~	2000	113.	ja	3	4	ja	ja	Fußbod~
6	Fran~	1700	84.8	ja	3	3	ja	ja	Fußbod~

```
# i 3 more variables: Baujahr <dbl>, Link <chr>, Abrufdatum <dtm>
```

Explorative Datenanalyse

Zu Beginn der explorativen Datenanalyse muss geprüft werden, ob die in der Datenquelle enthaltenen Daten auf eine bestimmte Art und Weise vorverarbeitet oder angepasst werden müssen. Hierzu kann zunächst mit `str(mieten)` die Struktur des Datensatzes angezeigt werden. Zudem wird mit dem Befehl `subset` der Datensatz nach Stadt unterteilt, falls im weiteren Verlauf der Analyse eine stadtsspezifische Aussage getroffen werden soll.

```
str(mieten)
```

```
tibble [100 x 12] (S3: tbl_df/tbl/data.frame)
 $ Ort      : chr [1:100] "Frankfurt" "Frankfurt" "Frankfurt" "Frankfurt" ...
 $ Kaltmiete : num [1:100] 1800 1500 2650 1500 2000 1700 1480 2800 1080 2600 ...
 $ Wohnflaeche : num [1:100] 70 60 146 72 113 ...
 $ Parkplatz  : chr [1:100] "ja" "ja" "ja" "nein" ...
 $ Etage      : chr [1:100] "1" "1" "1" "1" ...
 $ Zimmer     : num [1:100] 2 1 3 2 4 3 2 3 2 4 ...
 $ Balkon     : chr [1:100] "ja" "ja" "ja" "ja" ...
 $ Einbaukueche: chr [1:100] "ja" "ja" "ja" "ja" ...
 $ Heizung    : chr [1:100] "Fußbodenheizung" "Zentralheizung" "Fußbodenheizung" "Fußbodenheizung" ...
 $ Baujahr    : num [1:100] 2022 1970 2017 2021 2015 ...
 $ Link       : chr [1:100] "https://www.immobilienscout24.de/expose/136299839?referrer=R...
 $ Abrufdatum : POSIXct[1:100], format: "2023-12-28" "2023-12-28" ...
```

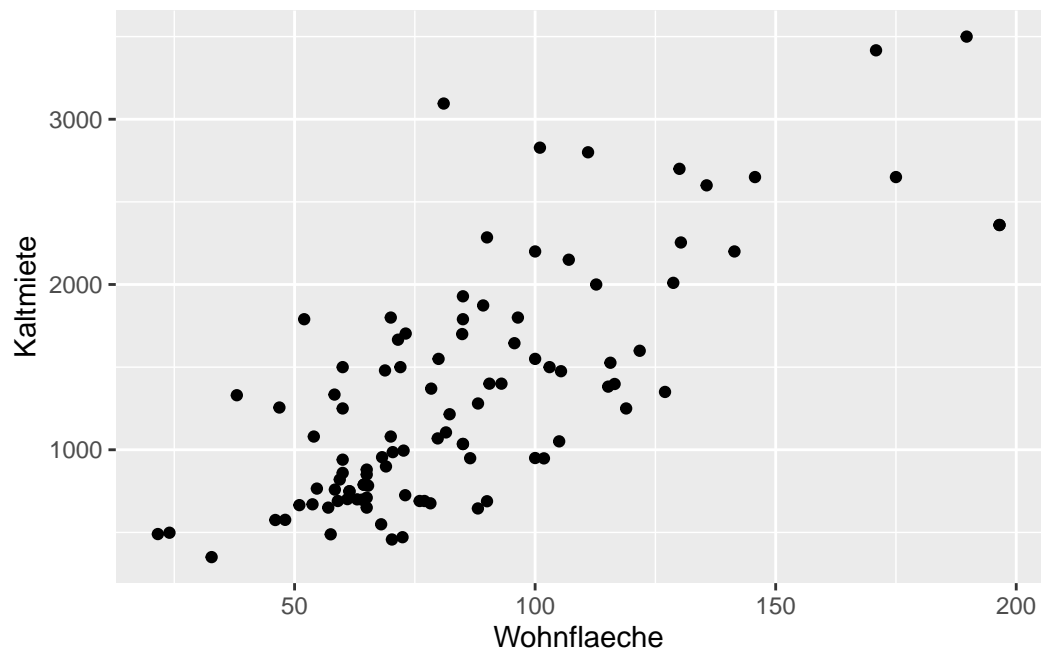
```
miete_ffm <- subset(mieten, Ort == "Frankfurt")
miete_lpz <- subset(mieten, Ort == "Leipzig")
```

Es kann festgestellt werden, dass im Datensatz sowohl kategoriale nominale Variablen wie `Heizung` oder `Zimmer`, als auch metrische verhältnisskalierte Variablen wie `Kaltmiete` oder `Wohnflaeche` auftreten. Zunächst werden keine Variablen angepasst bzw. Werte ersetzt, da für die späteren Diagramme die kategorial nominalen Variablen als Achsenbeschriftung gut verwendet werden können. Um nicht nur den Gesamtpreis der Kaltmiete zu betrachten, wird der Quadratmeterpreis mit in den Datensatz aufgenommen:

```
mieten <- mieten |>
  mutate(ppqm = Kaltmiete / Wohnflaeche)
```

Zuerst soll auf den Zusammenhang von `Kaltmiete` und `Wohnflaeche` geschaut werden, bei zwei metrisch verhältnisskalierten Variablen bietet sich dafür ein Scatterplot an.

```
gf_point(Kaltmiete ~ Wohnflaeche, data = mieten)
```



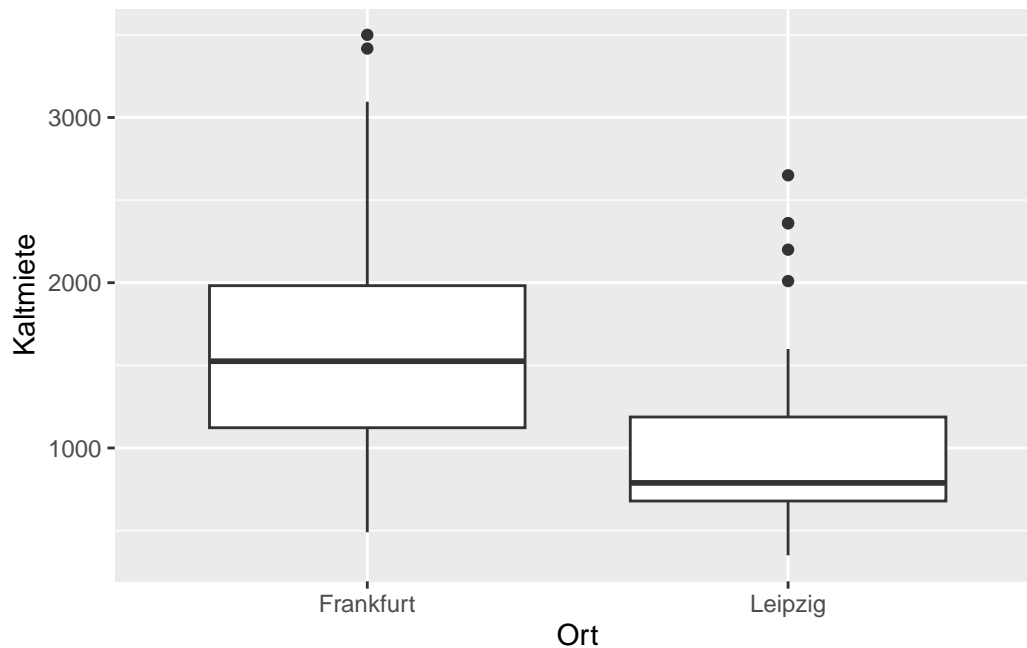
Grundsätzlich lässt sich ein positiver Zusammenhang zwischen **Kaltmiete** und **Wohnflaeche** erkennen, wobei die Streuung der Kaltmiete mit zunehmender Wohnfläche zunimmt. Nun muss dieses Diagramm jedoch um die Information des Ortes erweitert werden, um eine genauere Aussage treffen zu können. Dafür wird der Code um den Zusatz `color = ~ Ort` erweitert:

```
gf_point(Kaltmiete ~ Wohnflaeche, data = mieten, color = ~ Ort)
```



Hier lässt sich nun erkennen, dass die erfassten Mieten im Datensatz in Frankfurt tendenziell höher sind, als in Leipzig. Bei vergleichbarer Wohnfläche liegen die gefärbten Punkte für Frankfurt stets über den Punkten von Leipzig. Um den Eindruck der Mietunterschiede zu festigen, kann ein Boxplot verwendet werden.

```
gf_boxplot(Kaltmiete ~ Ort, data = mieten)
```



Das Boxplot zeigt, dass der Median für die Kaltmiete in Frankfurt deutlich über dem Median von Leipzig liegt. Zudem ist der 1,5-fache Interquartilsabstand bei Frankfurt größer als bei Leipzig, und die Whisker sind bei Frankfurt ebenfalls länger. Für Leipzig gibt es jedoch mehr Ausreißer, nämlich 4, verglichen mit den 2 Ausreißern von Frankfurt. Mithilfe der Funktion `mean` kann der Mittelwerte der Mieten in beiden Städten betrachtet werden.

```
mean(~ Kaltmiete, data = miete_ffm)
```

```
[1] 1652.821
```

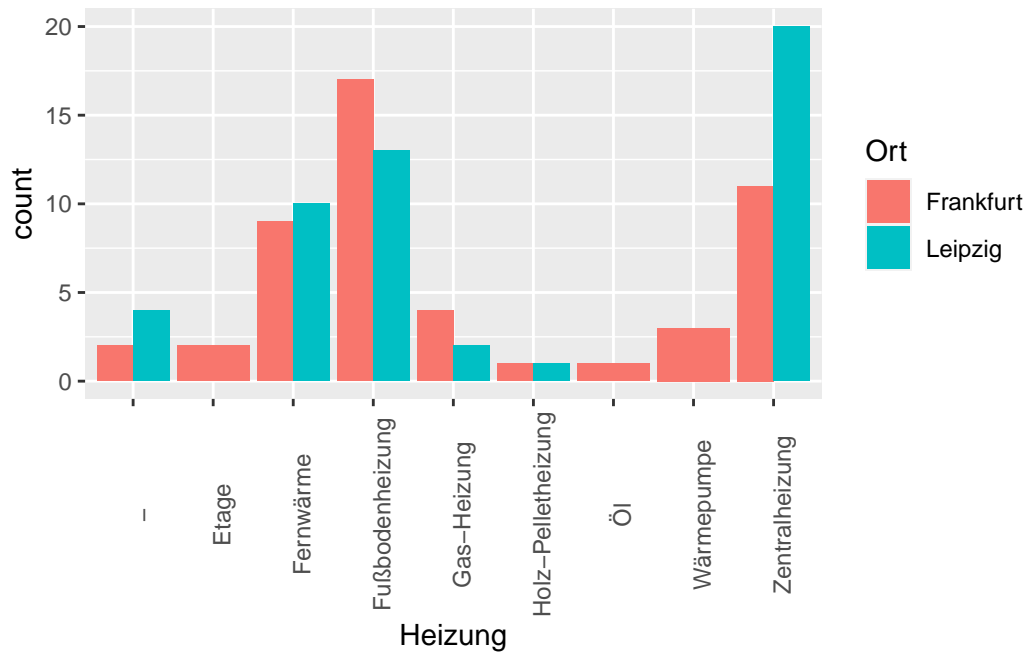
```
mean(~ Kaltmiete, data = miete_lpz)
```

```
[1] 1011.642
```

Der Mittelwert für die Kaltmiete liegt in Frankfurt bei 1.652,82€ und damit mehr als 600€ über dem Kaltmietendurchschnitt von Leipzig (1.011,64€).

Es sollen nun auch die weiteren Variablen untersucht werden, angefangen mit der Variablen **Heizung**. Um die verschiedenen Ausprägungen zu vergleichen und ihre absolute Häufigkeit darzustellen, eignet sich ein Säulendiagramm:

```
gf_bar( ~ Heizung, data = mieten, fill = ~ Ort, position = position_dodge()) + theme(ax...
```

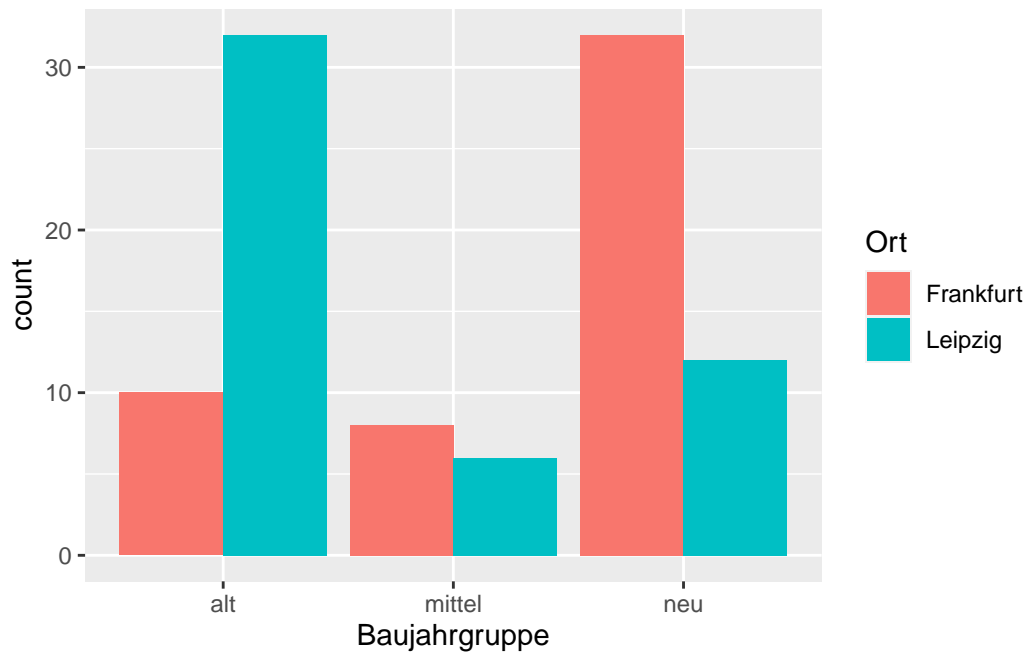


Das Säulendiagramm zeigt, dass die Fußbodenheizung in Frankfurt am weitesten verbreitet ist, gefolgt von der Zentralheizung und der Fernwärme. In Leipzig ist die Zentralheizung am weitesten verbreitet, gefolgt von der Fußbodenheizung und der Fernwärme. In einigen Beobachtungen und häufiger in Leipzig als in Frankfurt, wurde der Heizungstyp nicht angegeben.

Bei der Variablen **Baujahr** handelt es sich hier um eine diskrete Variable. Aufgrund der Vielzahl an verschiedenen Jahren im Datensatz eignet sich jedoch die Verwendung des tatsächlichen Baujahres nicht, da die Darstellungen sonst sehr unübersichtlich werden. Stattdessen soll eine Klassifizierung in "alt - mittel - neu" vorgenommen werden, um die Baujahre zusammenzufassen.

```
mieten <- mieten %>%
  mutate(Baujahrgruppe = case_when(
    is.na(Baujahr) ~ "NA",
    as.integer(Baujahr) < 1970 ~ "alt",
    between(as.integer(Baujahr), 1970, 2000) ~ "mittel",
    TRUE ~ "neu"
  ))

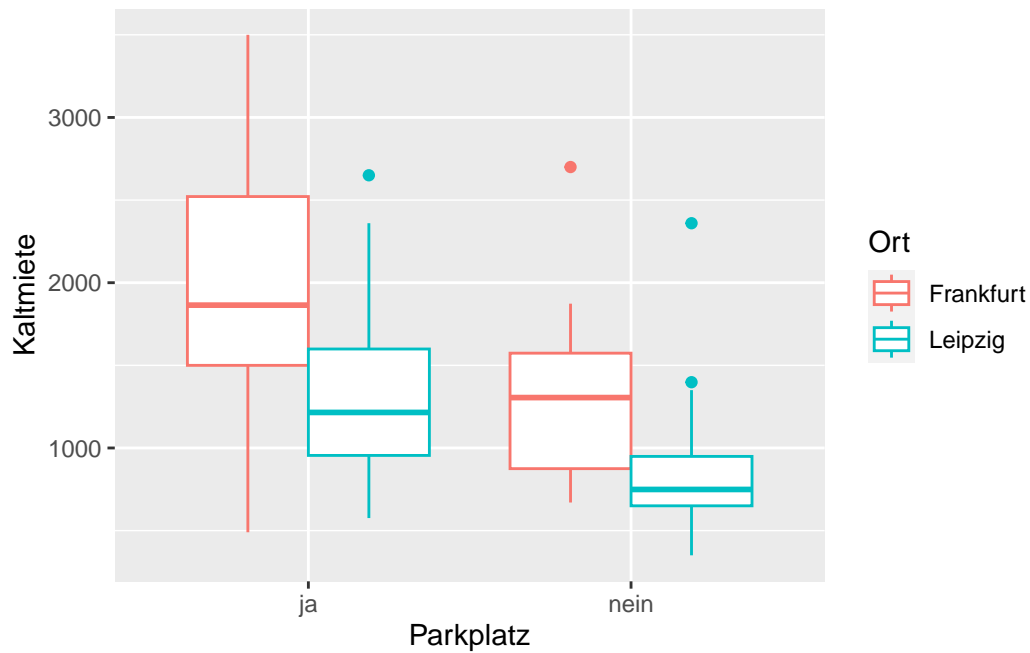
gf_bar(~ Baujahrgruppe, data = mieten, fill = ~ Ort, position = position_dodge())
```

Im Säulendiagramm ist erkennbar, dass die meisten Beobachtungen in Frankfurt in die Kategorie **neu** (Baujahr > 2000) fallen, gefolgt von **alt** (Baujahr < 1970) und **mittel**. In Leipzig dominieren die Beobachtungen mit **alt**, dann kommen **neue** Baujahre. Es treten in beiden Städten nur wenige Beobachtungen mit Baujahren zwischen 1970 und 2000 auf.

Parkplatz

```
gf_boxplot(Kaltmiete ~ Parkplatz, data = mieten, color = ~ Ort)
```



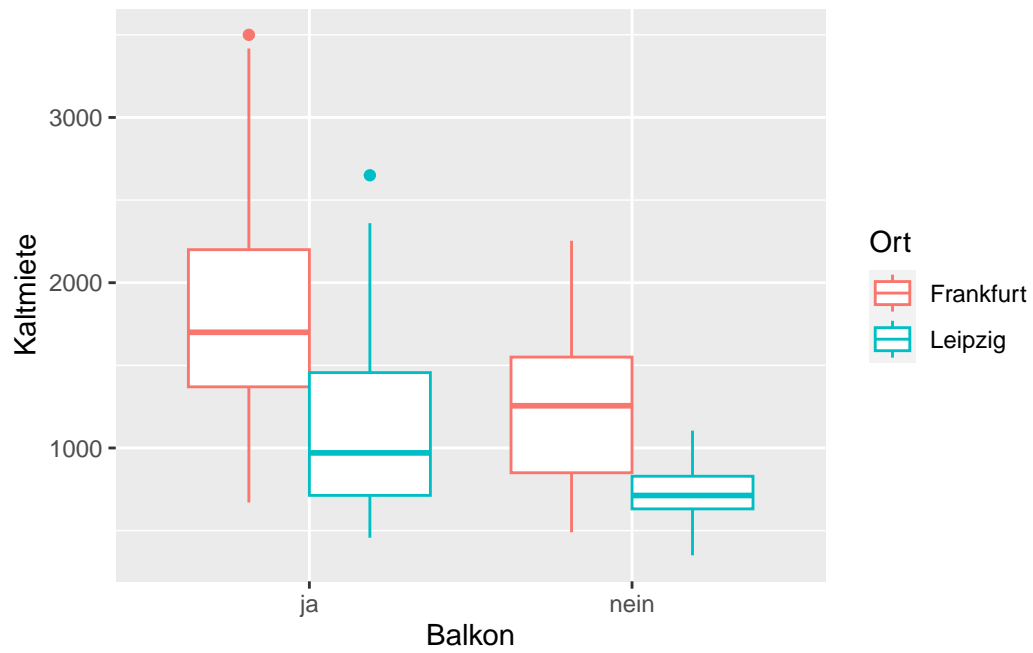
```
tally(Ort ~ Parkplatz, data = mieten)
```

	Parkplatz	
Ort	ja	nein
Frankfurt	26	24
Leipzig	17	33

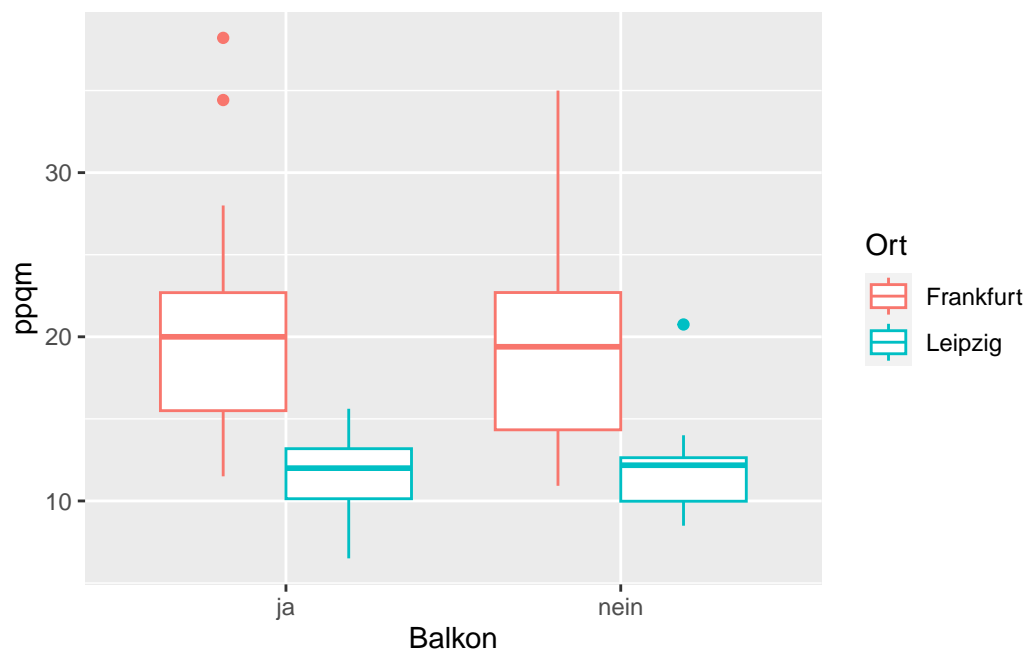
Klare Tendenz in beiden Städten Wohnungen mit Parkplatz sind im Schnitt teurer ? Parkplatz in der Kaltmiete enthalten? ? Eventuell teurere Wohnung hat eher eine Tiefgarage oder einen sonstigen Stellplatz auf dem Grundstück?

Balkon

```
gf_boxplot(Kaltmiete ~ Balkon, data = mieten, color = ~ Ort)
```



```
gf_boxplot(ppqm ~ Balkon, data = mieten, color = ~ Ort)
```



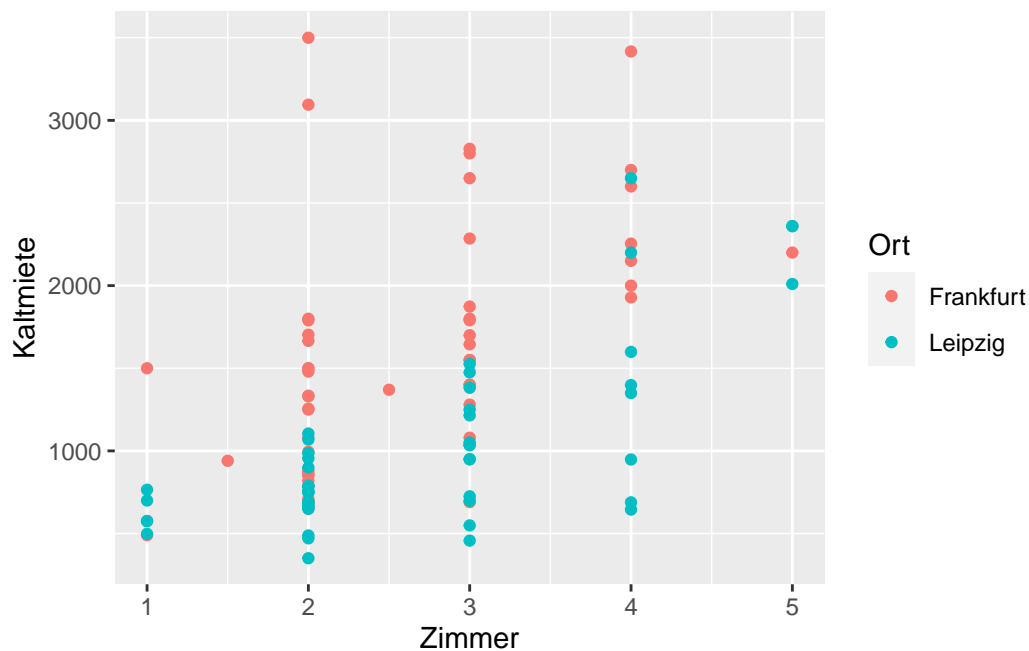
```
tally(Ort ~ Balkon, data = mieten)
```

Ort	Balkon	
	ja	nein
Frankfurt	37	13
Leipzig	30	20

Erstes Boxplot zeigt Wohnung mit Balkon sind meist teurer Da der Balkon auch zu 25% in die gesamte Wohnflaeche eingerechnet ist, wird im zweiten Boxplot auch der Quadratmeterpreis betrachtet. Hier ist auch ein Zusammenhang zwischen Balkon “Ja” und einen höheren Kaltmiete zu erkennen. ?Interpretation: Balkon steigert Kaltmiete unabhängig von der Wohnfläche. Interessant wäre zusätzlich eine Analyse, wenn Wohnfläche um die Fläche des Balkons bereinigt wäre.

Zimmer

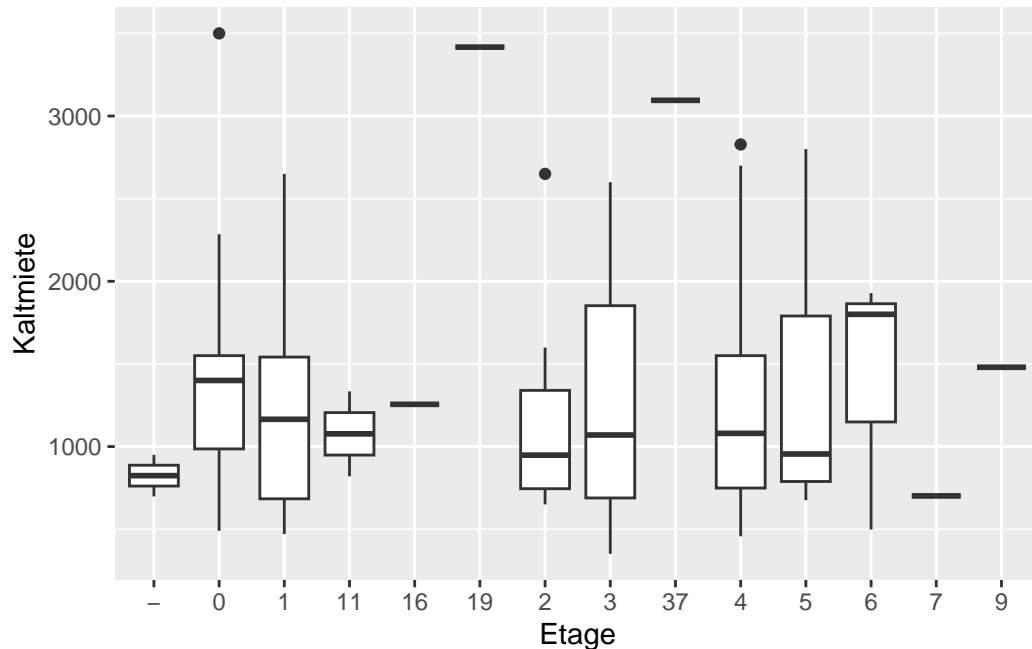
```
gf_point(Kaltmiete ~ Zimmer, data = mieten, color = ~ Ort)
```



```
gf_point(Wohnflaeche ~ Zimmer, data = mieten, color = ~ Ort)
```


Etage

```
gf_boxplot(Kaltemiete ~ Etage, data = mieten)
```



Schlecht bis gar nicht beschreibbar/interpretierbar

Modellierung

Aus den zahlreichen Diagrammen des vorherigen Abschnitts der explorativen Datenanalyse konnten sich bereits diverse Zusammenhänge erkennen lassen. Dieser letzte Teil der Untersuchung der gegebenen Daten beschäftigt sich abschließend mit der Modellierung der Kaltmiete unter Verwendung der zur Verfügung stehenden Variablen wie der Wohnfläche, der Art der Heizung oder dem Vorhandensein eines Balkons. Ziel ist hierbei die Erstellung eines Modells, durch das die Variable **Kaltemiete** bestmöglich erklärt werden kann.

Modellierung über Durchschnittsmietpreis

Das erste Diagramm der explorativen Datenanalyse, in dem die **Kaltemiete** der Inserate zusammen mit deren **Wohnfläche** im Streudiagramm dargestellt wurden, lässt einen positiven Zusammenhang der **Kaltemiete** zur **Wohnfläche** vermuten. In einem ersten einfachen Modell, mit dem dieser Zusammenhang modelliert werden soll, kann aus den Daten beispielsweise der Durchschnitt des Quadratmeterpreises berechnet werden.

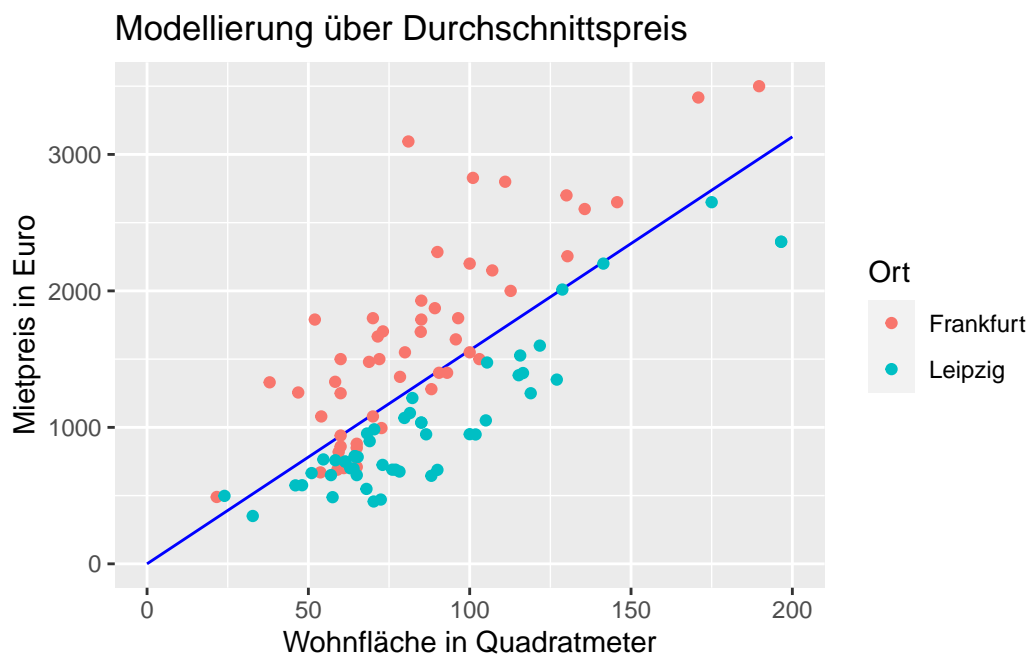
```
sum_wohnflaeche <- sum(~ Wohnflaeche, data = mieten)
sum_kaltemiete <- sum(~ Kaltmiete, data = mieten)
price_per_squaremeter <- sum_kaltemiete / sum_wohnflaeche
```

Damit ergibt sich als erstes einfaches Modell für die Kaltmiete unter Verwendung der Wohnfläche als unabhängige Variable folgende Gleichung:

$$\text{Kaltmiete} = \text{Wohnfläche} \cdot 15.65\text{€}$$

Wir betrachten das Modell, indem die berechnete Gerade in das Streudiagramm der explorativen Datenanalyse eingezeichnet wird:

```
pps_x = c(0, 200)
pps_y = c(0, price_per_squaremeter * 200)
gf_line(pps_y ~ pps_x, color = "blue") |>
  gf_point(Kaltmiete ~ Wohnflaeche, data = mieten, color = ~ Ort) |>
  gf_labs(x = "Wohnfläche in Quadratmeter",
    y = "Mietpreis in Euro",
    title = "Modellierung über Durchschnittspreis")
```



Dazu kann noch der Korrelationskoeffizient bestimmt werden.

```
cor_miete_flaeche = cor(Wohnflaeche ~ Kaltmiete, data = mieten)
```

Dabei wird festgestellt, dass die `Kaltemiete` zu 74.23% mit der Wohnfläche korreliert. Das Modell kann demnach bereits für einen groben Richtwert verwendet werden. Ziel ist jedoch eine noch genauere Modellierung der Kaltemiete unter Berücksichtigung der weiteren Daten.

Modellierung über die lineare Regression

Um die weiteren Variablen miteinzubeziehen, verwenden wir die lineare Regression, die mit der `lm()`-Funktion auf die Daten angewendet werden kann. Wir starten zunächst erneut mit der `Kaltemiete` und der `Wohnflaeche`.

```
km.lm1 <- lm(Kaltemiete ~ Wohnflaeche, data = mieten)
summary(km.lm1)
```

Call:

```
lm(formula = Kaltemiete ~ Wohnflaeche, data = mieten)
```

Residuals:

Min	1Q	Median	3Q	Max
-732.9	-294.9	-117.0	322.4	1826.9

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	16.882	129.179	0.131	0.896
Wohnflaeche	15.447	1.408	10.968	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 480 on 98 degrees of freedom

Multiple R-squared: 0.5511, Adjusted R-squared: 0.5465

F-statistic: 120.3 on 1 and 98 DF, p-value: < 2.2e-16

Mit einem Bestimmtheitsmaß von $R^2 = 0.56$, haben wir mit der linearen Regression allein unter Verwendung der `Wohnflaeche` noch kein gutes Modell erzeugt.

Da zu Beginn der explorativen Datenanalyse festgestellt wurde, dass sich die Kaltmieten zwischen den beiden Orten Frankfurt am Main und Leipzig unter sonst gleichen Bedingungen bereits stark unterscheidet, soll diese zuerst in die Modellierung miteinbezogen werden.

```
km.lm2 <- lm(Kaltemiete ~ Wohnflaeche + Ort, data = mieten)
summary(km.lm2)
```

Call:

```
lm(formula = Kaltemiete ~ Wohnflaeche + Ort, data = mieten)
```


Residuals:

Min	1Q	Median	3Q	Max
-640.77	-278.36	38.32	189.69	1492.77

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	329.208	97.419	3.379	0.00105 **
Wohnflaeche	15.716	1.004	15.653	< 2e-16 ***
OrtLeipzig	-670.496	68.440	-9.797	3.68e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 342.1 on 97 degrees of freedom

Multiple R-squared: 0.7743, Adjusted R-squared: 0.7697

F-statistic: 166.4 on 2 and 97 DF, p-value: < 2.2e-16

```
km.lm2.coef <- coef(km.lm2)
print(km.lm2.coef)
```

```
(Intercept) Wohnflaeche OrtLeipzig
329.20836    15.71636   -670.49589
```

Mit dem Miteinbeziehen der Indikatorvariable x_2 beziehungsweise der kategorialen Variable des Ortes erhalten wir ein Bestimmtheitsmaß von $R^2 = 0.78$. Die Variation des Mietpreises kann also zu 78% durch den Mietpreis und dem Ort erklärt werden. An der Zusammenfassung der Ergebnisse lässt sich außerdem ablesen, dass in Leipzig die Kaltmiete im Mittelwert 655,13€ billiger als in Frankfurt ist. Das Modell lässt sich nun wie folgt darstellen:

$$\hat{y}_i = 329.21 + x_{1i} \cdot 15.72 - 670.5 \cdot \begin{cases} 1: & x_{2i} = \text{Leipzig} \\ 0: & x_{2i} = \text{Frankfurt} \end{cases}$$

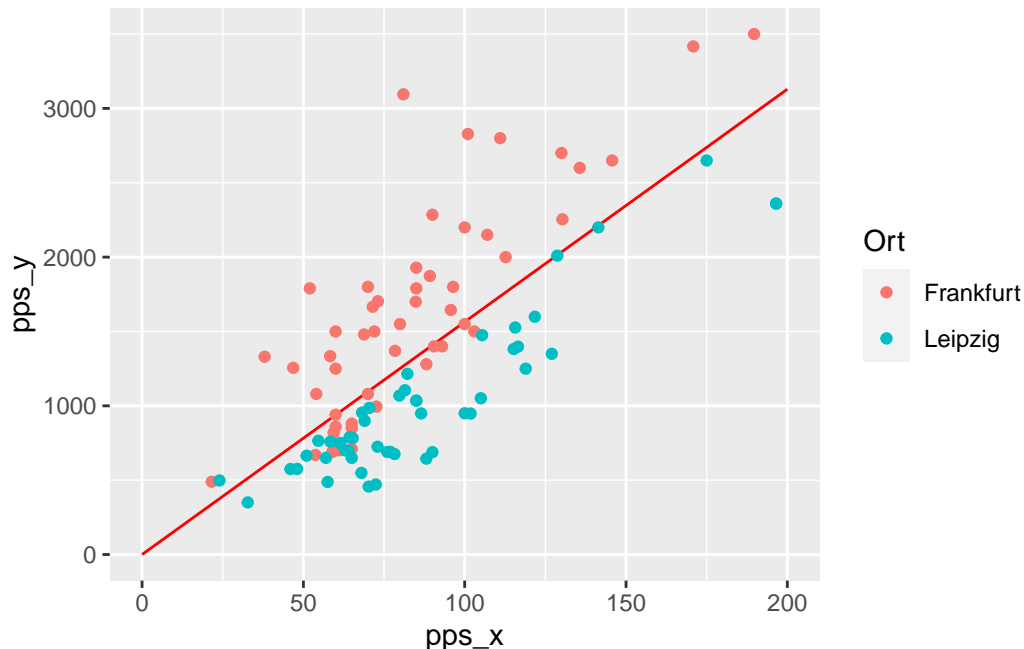
Wobei x_{1i} die Wohnflaeche und x_{2i} entsprechend den Wohnort darstellt.

Entsprechend dem berechneten P-Wert von $3.68 \cdot 10^{-16}$ kann das Ergebnis in Bezug auf den Wohnort unter Verwendung eines Signifikanzniveaus von 5% als statistisch signifikant bezeichnet werden. Die H_0 -Hypothese $\mu_{\text{Frankfurt}} = \mu_{\text{Leipzig}}$ kann somit verworfen werden.

Grafisch erhalten wir dadurch die zwei Geraden, abhängig des Ortes:

```
frankfurt.p1 = c(0, 200)
frankfurt.p2 = c(0, km.lm2.coef[0] * 200)
```

```
leipzig.p1 = c(0, 200)
pps_y = c(0, price_per_squaremeter * 200)
gf_line(pps_y ~ pps_x, color = "red") |>
  gf_point(Kaltmiete ~ Wohnflaeche, data = mieten, color = ~ Ort)
```



Andere Modellierung

Aus den zahlreichen Diagrammen des vorherigen Abschnitts der explorativen Datenanalyse konnten sich bereits diverse Zusammenhänge erkennen lassen. Dieser letzte Teil der Untersuchung der gegebenen Daten beschäftigt sich abschließend mit der Modellierung der Kaltmiete unter Verwendung der zur Verfügung stehenden Variablen wie der Wohnfläche, der Art der Heizung oder dem Vorhandensein eines Balkons. Ziel ist hierbei die Erstellung eines Modells, durch das die Variable `Kaltmiete` bestmöglich erklärt werden kann.

Modellierung über Durchschnittsmietpreis

Das erste Diagramm der explorativen Datenanalyse, in dem die `Kaltmiete` der Inserate zusammen mit deren `Wohnfläche` im Streudiagramm dargestellt wurden, lässt einen positiven Zusammenhang der `Kaltmiete` zur `Wohnfläche` vermuten. In einem ersten einfachen Modell, mit dem dieser Zusammenhang modelliert werden soll, kann aus den Daten beispielsweise der Durchschnitt des Quadratmeterpreises berechnet werden.

Damit ergibt sich als erstes einfaches Modell für die **Kaltniete** unter Verwendung der **Wohnfläche** als unabhängige Variable folgende Gleichung:

$$Kaltmiete = Wohnfläche \cdot 15.65\text{€}$$

Wir betrachten das Modell, indem die berechnete Gerade in das Streudiagramm der explorativen Datenanalyse eingezeichnet wird:

Dazu kann noch der Korrelationskoeffizient bestimmt werden.

Dabei wird festgestellt, dass die **Kaltniete** zu 74.23% mit der Wohnfläche korreliert. Das Modell kann demnach bereits für einen groben Richtwert verwendet werden. Ziel ist jedoch eine noch genauere Modellierung der Kaltmiete unter Berücksichtigung der weiteren Daten.

Modellierung über die lineare Regression

Um die weiteren Variablen miteinzubeziehen, verwenden wir die lineare Regression, die mit der `lm()`-Funktion auf die Daten angewendet werden kann. Wir starten zunächst erneut mit der **Kaltniete** und der **Wohnfläche**.

Mit einem Bestimmtheitsmaß von $R^2 = 0.56$, haben wir mit der linearen Regression allein unter Verwendung der **Wohnfläche** noch kein gutes Modell erzeugt.

Da zu Beginn der explorativen Datenanalyse festgestellt wurde, dass sich die Kaltmieten zwischen den beiden Orten Frankfurt am Main und Leipzig unter sonst gleichen Bedingungen bereits stark unterscheidet, soll diese zuerst in die Modellierung miteinbezogen werden.

Mit dem Miteinbeziehen der Indikatorvariable x_2 beziehungsweise der kategorialen Variable des Ortes erhalten wir ein Bestimmtheitsmaß von $R^2 = 0.78$. Die Variation des Mietpreises kann also zu 78% durch den Mietpreis und dem Ort erklärt werden. An der Zusammenfassung der Ergebnisse lässt sich außerdem ablesen, dass in Leipzig die Kaltmiete im Mittelwert 655,13€ billiger als in Frankfurt ist. Das Modell lässt sich nun wie folgt darstellen:

$$\hat{y}_i = 329.21 + x_{1_i} \cdot 15.72 - 670.5 \cdot \begin{cases} 1 : & x_{2_i} = \text{Leipzig} \\ 0 : & x_{2_i} = \text{Frankfurt} \end{cases}$$

Wobei x_{1_i} die **Wohnfläche** und x_{2_i} entsprechend den **Wohnort** darstellt.

Entsprechend dem berechneten P-Wert von $3.68 \cdot 10^{-16}$ kann das Ergebnis in Bezug auf den Wohnort unter Verwendung eines Signifikanzniveaus von 5% als statistisch signifikant bezeichnet werden. Die H_0 -Hypothese $\mu_{\text{Frankfurt}} = \mu_{\text{Leipzig}}$ kann somit verworfen werden.

Grafisch erhalten wir dadurch die zwei Geraden, abhängig des Ortes:

Zusammenfassung

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

Quellen und Hilfsmittel

Führen Sie hier die verwendeten Hilfsmittel sowie die verwendete Literatur auf.