

The background of the slide is a light, airy composition of numerous small, semi-transparent spheres and particles in various colors including white, light blue, yellow, orange, pink, and purple. These particles are scattered across the frame, with a denser cluster in the center where the text is located, creating a sense of depth and movement.

Angewandte Statistik

Julian Hinz — Universität Bielefeld

Session 2

Parametrische Regression

Wiederholung: Lineare Regression

Abgaben

1. Abgabe: 24. Mai (Aufgaben am 17. Mai)
2. Abgabe: 5. Juli (Aufgaben am 28. Juni)

Frist: [AoE](#) — “Anywhere on Earth”

an: abgabe-internationale-vwl@uni-bielefeld.de

Tutorien

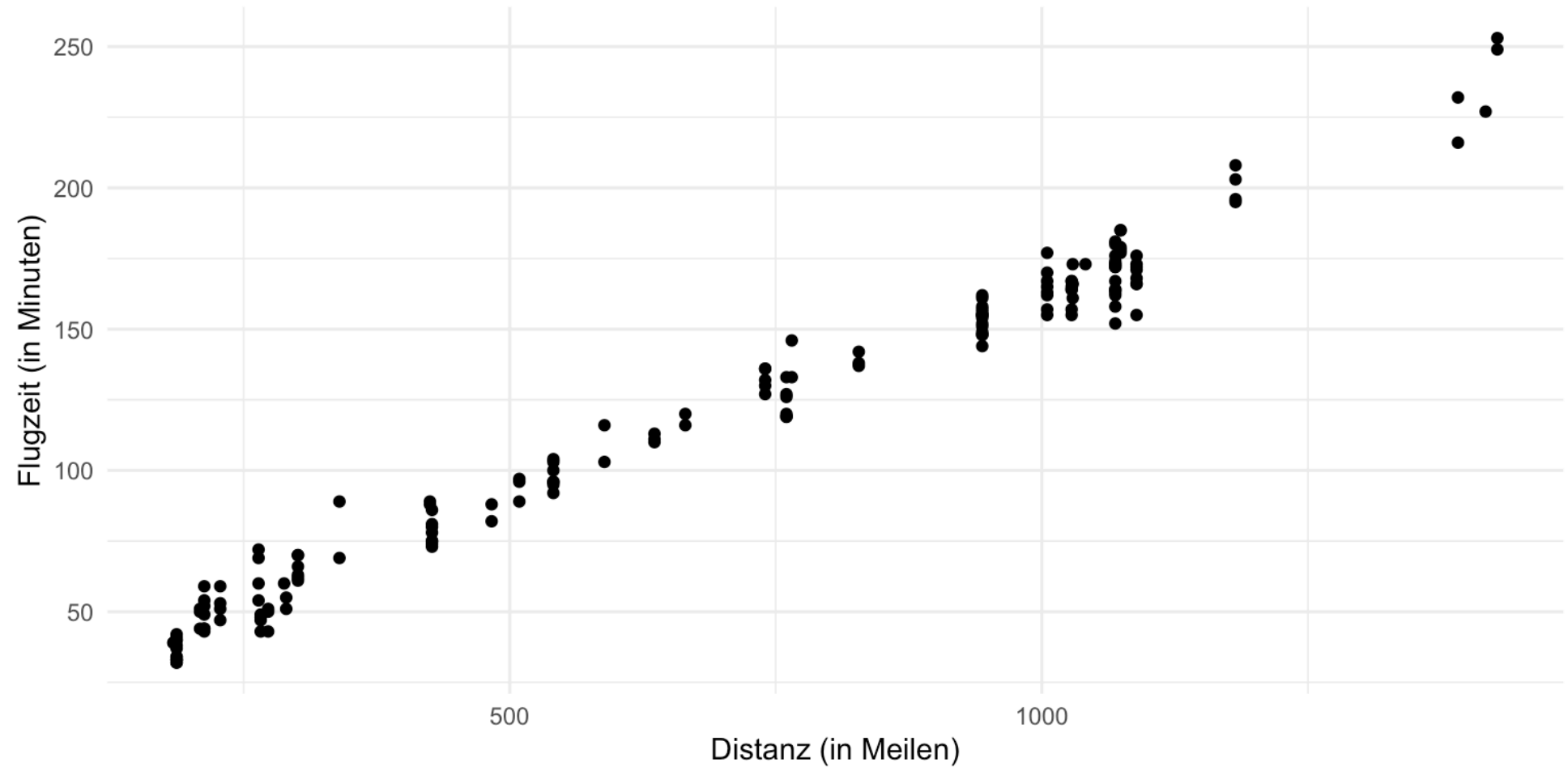
- 197 Anmeldungen für 120 Plätze
- kurzfristige Lösung: Live-Streaming per Zoom
- wir versuchen noch zusätzlichen Termin einzurichten

Parametrische Regressionsanalyse

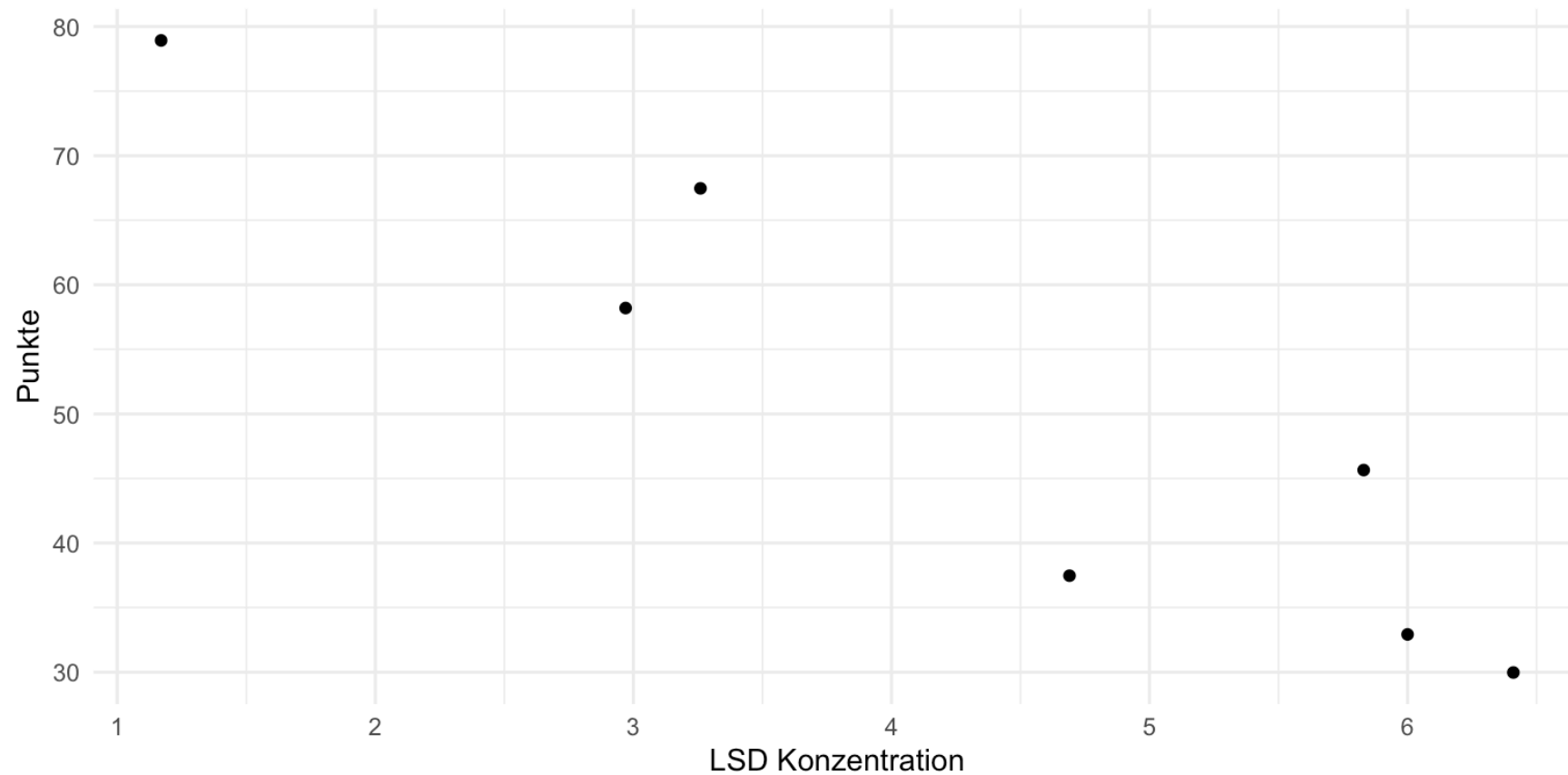
Parametrische Regressionsanalyse

(Wie) hängt bestimmte **Zielgröße** (Y) von erklärenden **Variablen** (x_1, \dots, x_p) ab?

Flugzeit vs. Flugstrecke



Punkte im Mathetest vs. LSD-Konzentration im Blut



Beispiele

- Wirtschaftliche Größe und Exporte
- Mietpreis hängt ab von Wohnfläche, Baujahr, Ausstattung, Lage, etc.
- WiWi-Note hängt ab von Abischnitt, Lernaufwand, etc.
- Gebrauchtwagenpreis hängt ab von Kilometerstand, Marke, Motorleistung, etc.
- ...

Analyse von Abhängigkeiten zwischen Variablen

Wir wollen aus Daten ein **Modell** lernen (d.h. schätzen), welches den Einfluss der erklärenden Variablen (x_1, \dots, x_p) auf die Zielgröße (Y) gut beschreibt:

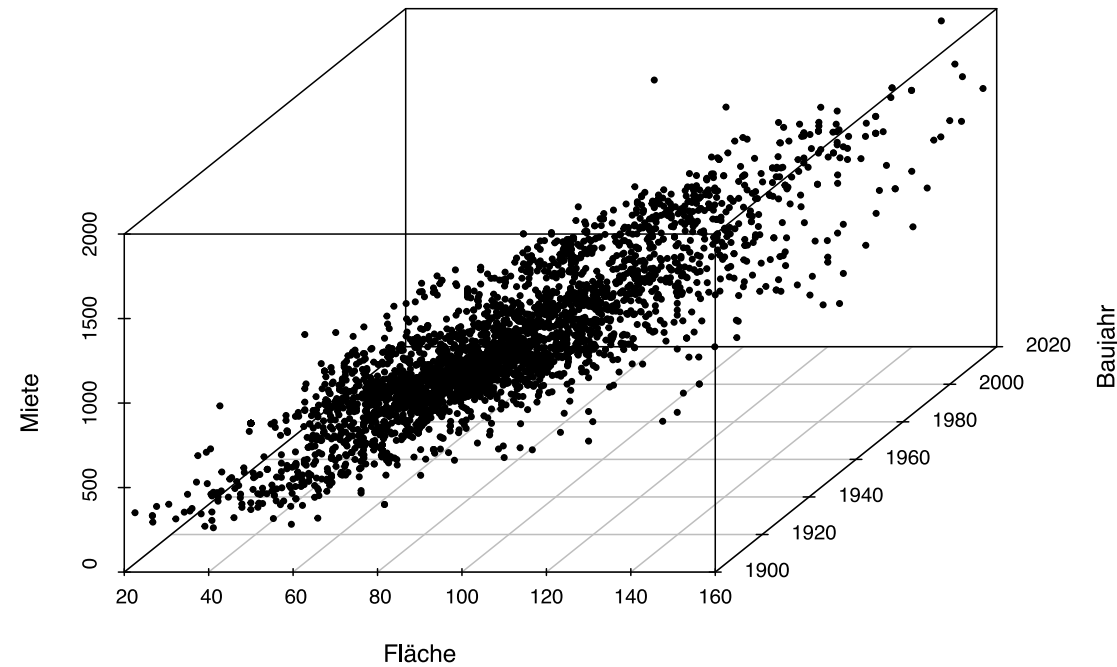
$$Y = f(x_1, \dots, x_p) + \text{zufälliger Fehler}$$

Wieso brauchen wir den Fehlerterm?

- Ein “perfekter” Zusammenhang liegt bei realen Prozessen so gut wie nie vor.
- Oftmals wird (Y) von weiteren, nicht beobachteten Variablen beeinflusst.

Im Idealfall erklärt (f) einen Großteil der Variabilität von (Y), so dass der Fehlerterm klein ist, was zu ziemlich genauen Vorhersagen führen kann.

Mietpreis (in Euro/Monat) vs. Wohnfläche (in qm) und Baujahr, für Bielefelder Mietwohnungen



Regression — allgemeine Formulierung und Ziele

Statistische Modellierung der Form

$$Y = f(x_1, \dots, x_p) + \epsilon$$

nennt man **Regression** (oder **Regressionsanalyse**)

Konkrete Ziele der Regression könnten u.a. sein

- Vorhersage von Werten (z.B. Flugzeiten)
- Zusammenhänge aufzeigen/testen ob ein Einfluss vorliegt (z.B. Einfluss von Drogen und kognitiver Leistung)
- Identifikation von Ausreißern (z.B. Mietpreise/Wucher)

Modellformulierung

Lineares Regressionsmodell

$$\underbrace{Y_i}_{\text{Zielgröße}} = \underbrace{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}_{\text{system. Effekt}} + \underbrace{\epsilon_i}_{\text{Fehler}}, \quad E(\epsilon_i) = 0, \quad i = 1, \dots, n.$$

Hier: *lineare* Form des systematischen Teils $f(x_{i1}, \dots, x_{ip})$

Warum ist diese Formulierung so wichtig?

- überraschend oft realistisch
- relativ robust gegenüber einzelnen Ausreißern
- sehr einfach umzusetzen und zu interpretieren
- viel flexibler als es auf den ersten Blick scheint

Einfaches lineare Modell ($p = 1$) in Matrixnotation

Die Modellspezifikation

$$Y_i = \beta_0 + \beta_1 x_{i1} + \epsilon_i, \quad i = 1, \dots, n,$$

umfasst genau genommen n Gleichungen, eine für jeden Datenpunkt:

$$Y_1 = \beta_0 + \beta_1 x_{11} + \epsilon_1$$

$$Y_2 = \beta_0 + \beta_1 x_{21} + \epsilon_2$$

$$\vdots$$

$$Y_n = \beta_0 + \beta_1 x_{n1} + \epsilon_n$$

Überführung in Matrixnotation

$$\begin{array}{l} Y_1 = \beta_0 + \beta_1 x_{11} + \epsilon_1 \\ \vdots \\ Y_n = \beta_0 + \beta_1 x_{n1} + \epsilon_n \end{array} \rightsquigarrow \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

wobei

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{n1} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Multiples lineares Modell ($p \geq 1$) in Matrixnotation

Komplett analog für das multiple lineare Modell

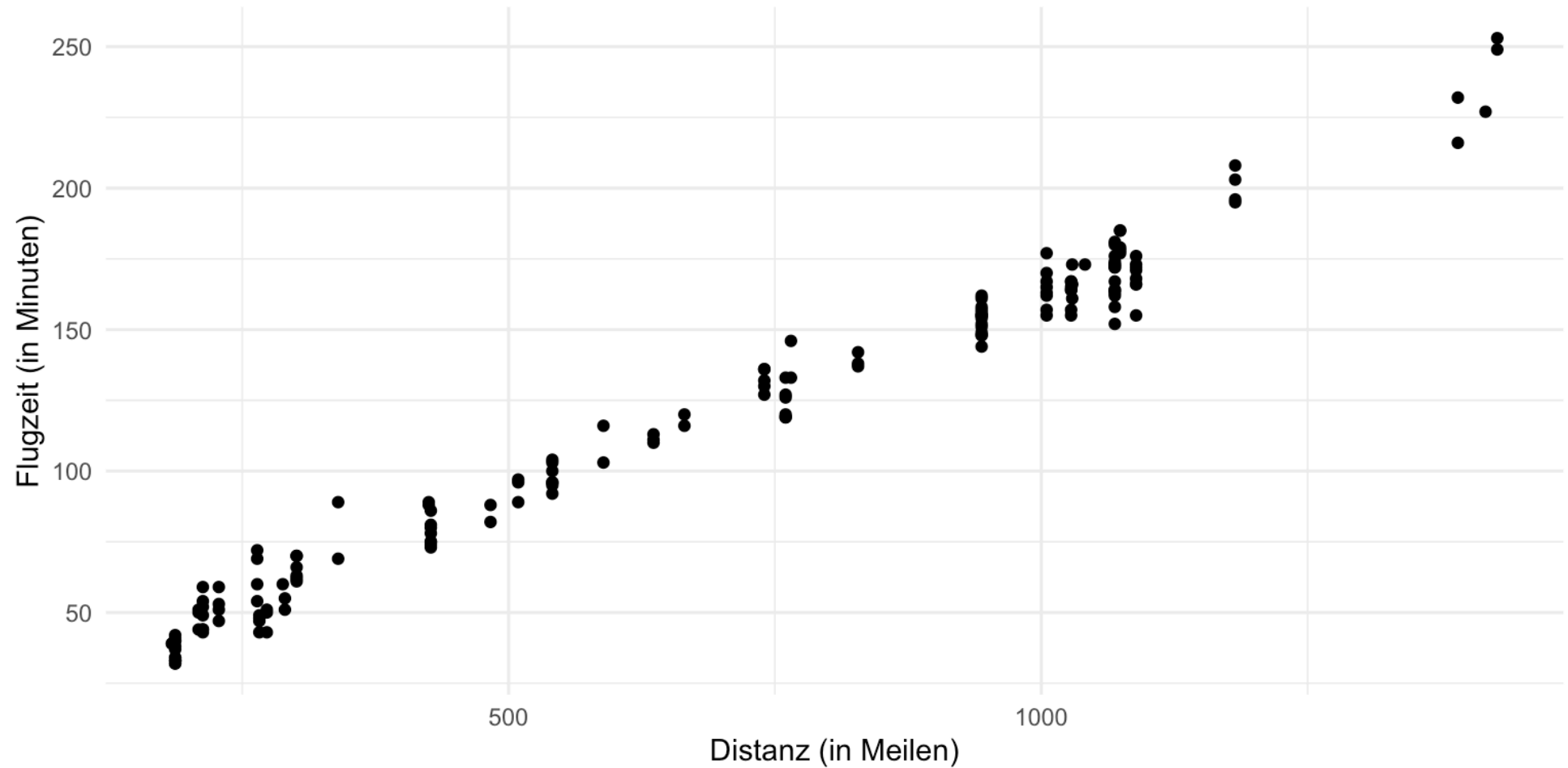
$$\begin{array}{l} Y_1 = \beta_0 + \beta_1 x_{11} + \dots + \beta_p x_{1p} + \epsilon_1 \\ \vdots \\ Y_n = \beta_0 + \beta_1 x_{n1} + \dots + \beta_p x_{np} + \epsilon_n \end{array} \rightsquigarrow \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

wobei

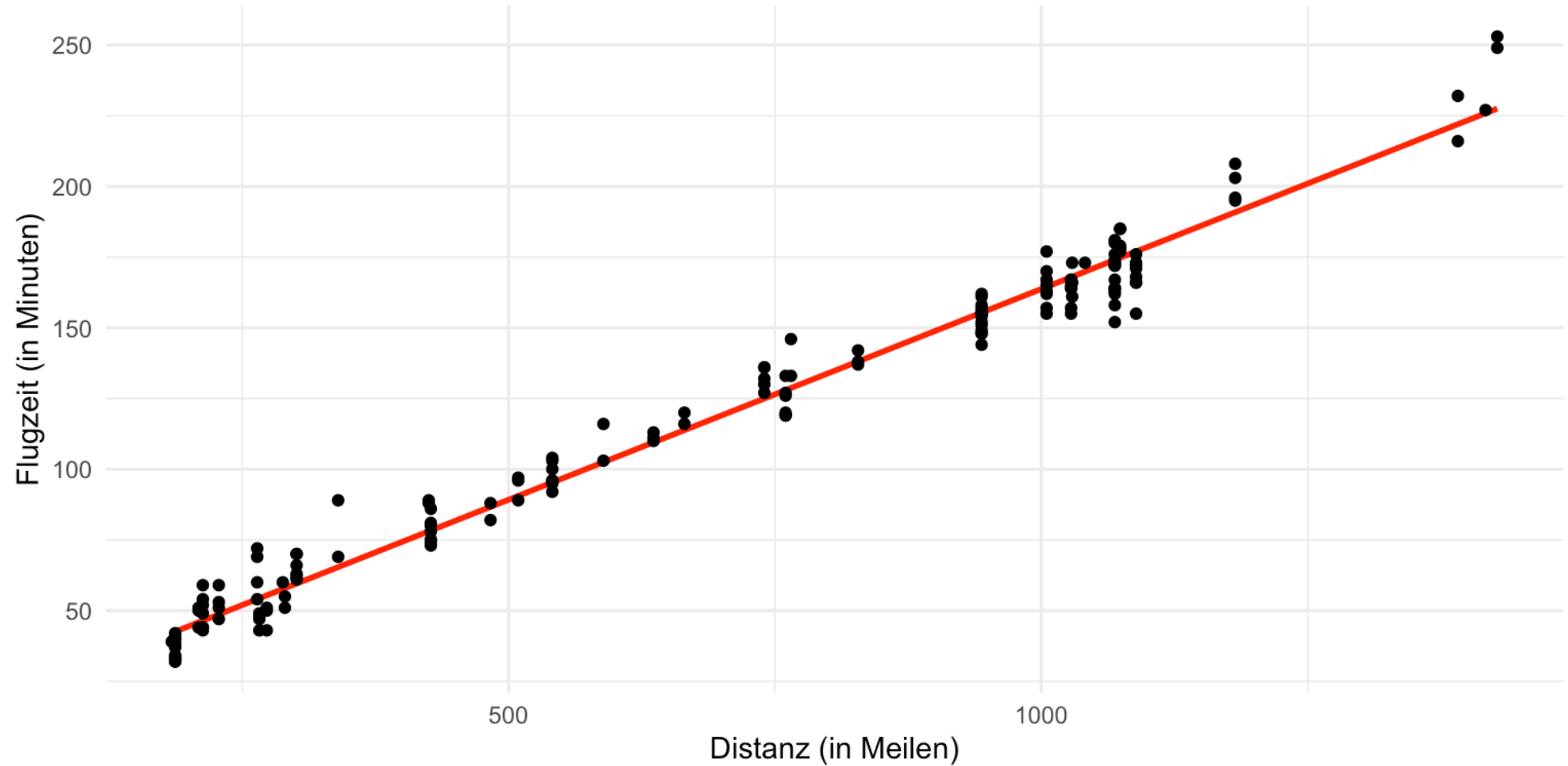
$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & & \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Parameterschätzung

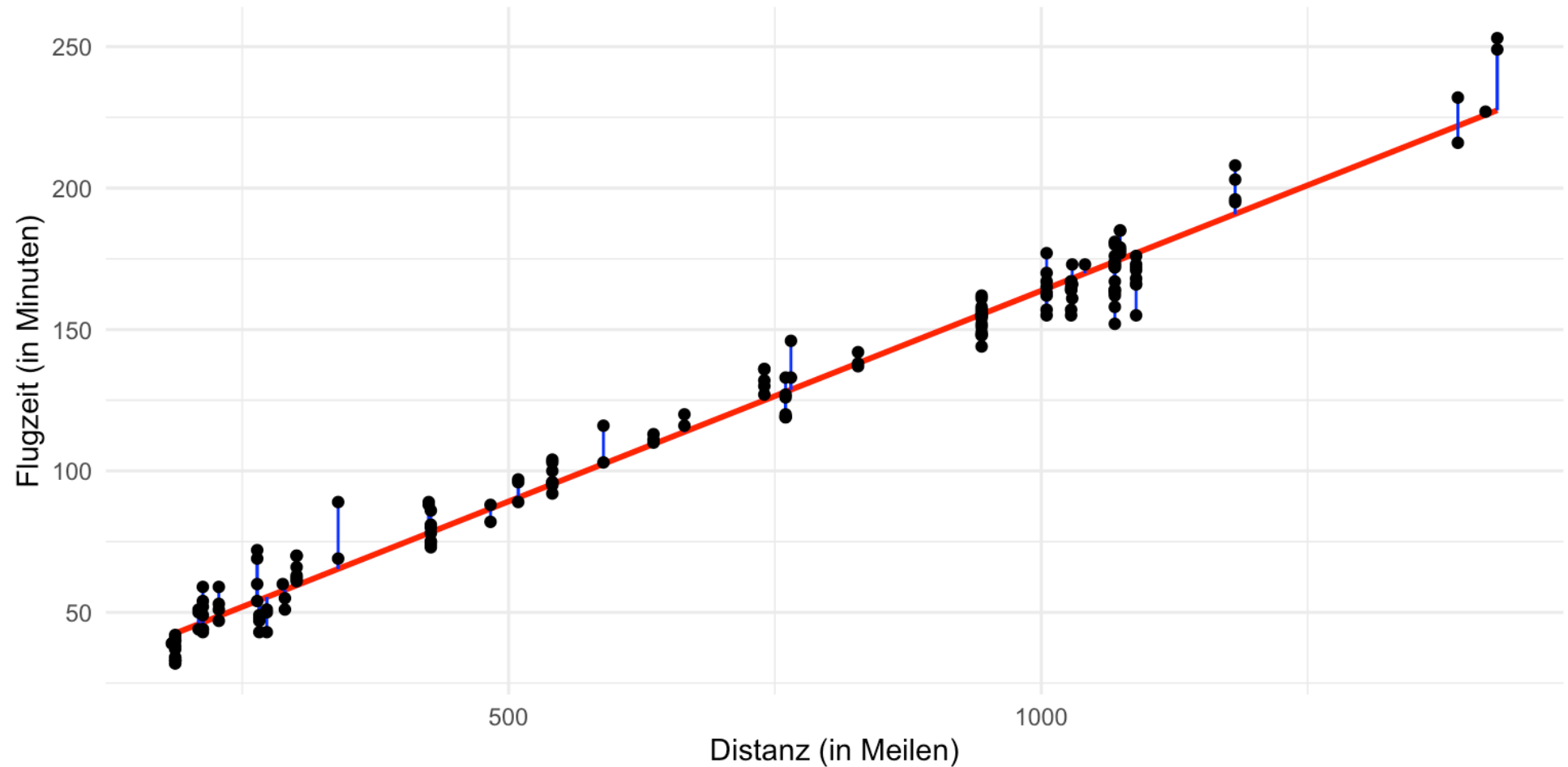
Schätzung der Regressionskoeffizienten



Schätzung der Regressionskoeffizienten



Schätzung der Regressionskoeffizienten



Schätzung der Regressionskoeffizienten

Als “Gesamtabstand” zwischen Modell und Datenpunkten definieren wir die Summe der quadrierten Fehler:

$$S(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n \left(y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \right)^2$$

Der **Kleinste-Quadrate-Schätzer** (KQS),

$$\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p) = \underset{\beta_0, \beta_1, \dots, \beta_p}{\operatorname{argmin}} S(\beta_0, \beta_1, \dots, \beta_p),$$

minimiert diesen Abstand.

KQS für das multiple lineare Regressionsmodell

Für das multiple lineare Regressionsmodell

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad E(\epsilon_i) = 0, \quad i = 1, \dots, n,$$

in Matrixnotation

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

ist der **Kleinste-Quadrate-Schätzer** gegeben durch

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

KQS in R für LSD-Beispiel

Mit Hilfe von `lm`-Funktion:

```
1 library(data.table)
2
3 daten = fread("../data/data2/LSD.txt")
4
5 regression <- lm(Punkte ~ LSD,
6                  data = daten)
7 print(regression)
```

Call:

```
lm(formula = Punkte ~ LSD, data = daten)
```

Coefficients:

(Intercept)	LSD
89.124	-9.009

KQS in R für LSD-Beispiel

“zu Fuß” mit Matrix Algebra:

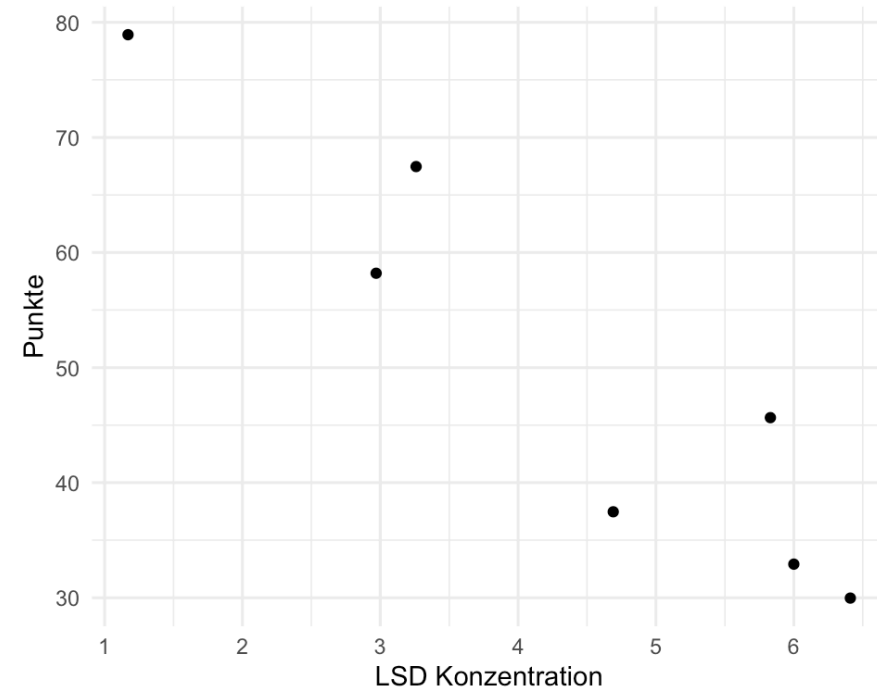
```
1 daten = fread("../data/data2/LSD.txt")
2
3 X = cbind(rep(1, 7), daten$LSD)
4 y = daten$Punkte
5
6 beta = solve(t(X) %*% X) %*% t(X) %*% y
7 print(beta)
```

```
      [,1]
[1,] 89.123874
[2,] -9.009466
```

KQS in R für LSD-Beispiel

“zu Fuß” mit Matrix Algebra:

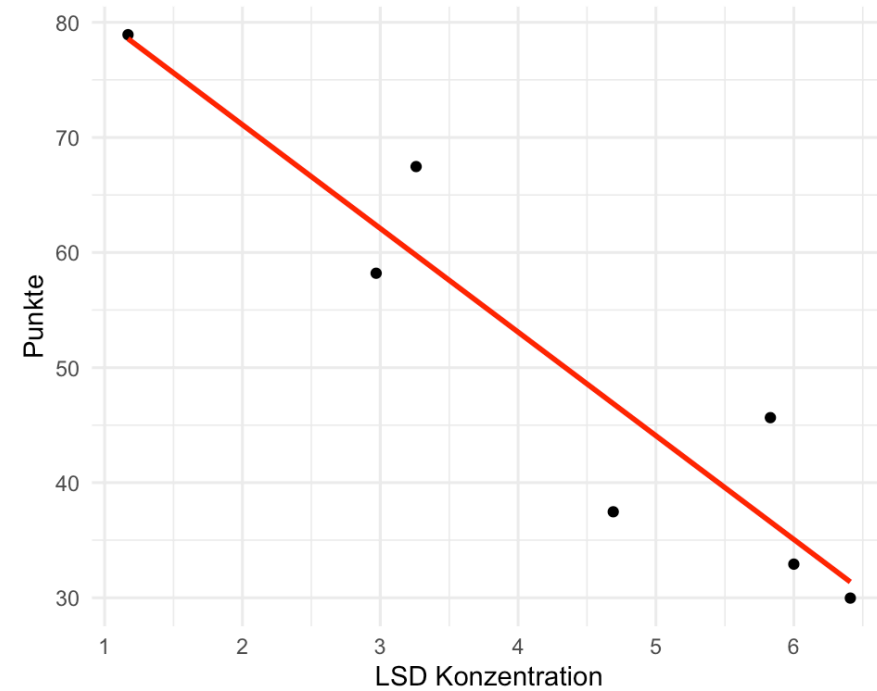
```
1 p_load(ggplot2)
2 p_load(data.table)
3
4 # daten laden
5 data = fread("../data/data2/LSD.txt")
6
7 # plotten
8 ggplot(data, aes(x = LSD, y = Punkte)) +
9   theme_minimal() +
10   geom_point() +
11   scale_x_continuous("LSD Konzentration") +
12   scale_y_continuous("Punkte")
```



KQS in R für LSD-Beispiel

“zu Fuß” mit Matrix Algebra:

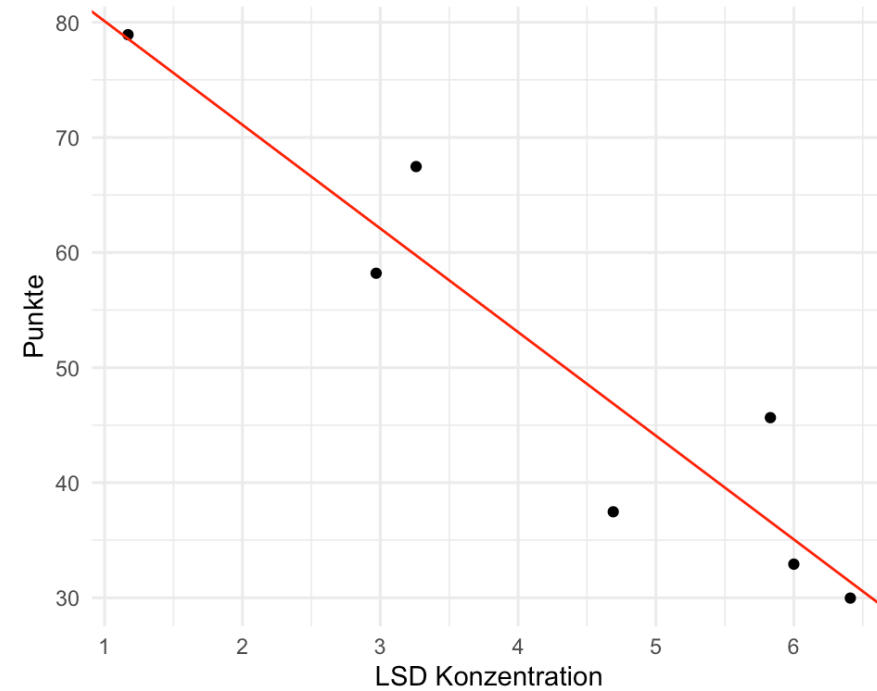
```
1 p_load(ggplot2)
2 p_load(data.table)
3
4 # daten laden
5 data = fread("../data/data2/LSD.txt")
6
7 # plotten
8 ggplot(data, aes(x = LSD, y = Punkte)) +
9   theme_minimal() +
10   geom_point() +
11   geom_smooth(method = "lm", color = "red",
12   scale_x_continuous("LSD Konzentration") +
13   scale_y_continuous("Punkte")
```



KQS in R für LSD-Beispiel

“zu Fuß” mit Matrix Algebra:

```
1 p_load(ggplot2)
2 p_load(data.table)
3
4 # daten laden
5 data = fread("../data/data2/LSD.txt")
6
7 # regression manuell rechnen
8 X = cbind(rep(1, 7), daten$LSD)
9 y = daten$Punkte
10 beta = solve(t(X) %*% X) %*% t(X) %*% y
11
12 # plotten
13 ggplot(data, aes(x = LSD, y = Punkte)) +
14   theme_minimal() +
15   geom_point() +
16   geom_abline(intercept = beta[1], slope = beta[2]) +
17   scale_x_continuous("LSD Konzentration") +
18   scale_y_continuous("Punkte")
```



KQS in R für Mietspiegel-Beispiel

Mit Hilfe von - Funktion:

```
1 library(data.table)
2
3 # daten laden
4 daten = fread("../data/data2/rents.csv")
5
6 # lm funktion
7 regression = lm(rent ~ area + year,
8                 data = daten)
9 summary(regression)
```

Call:

```
lm(formula = rent ~ area + year, data = daten)
```

Residuals:

Min	1Q	Median	3Q	Max
-386.54	-73.32	-3.15	69.16	1062.81

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.147e+03	1.459e+02	-14.71	<2e-16

area	7.906e+00	9.265e-02	85.33	<2e-16

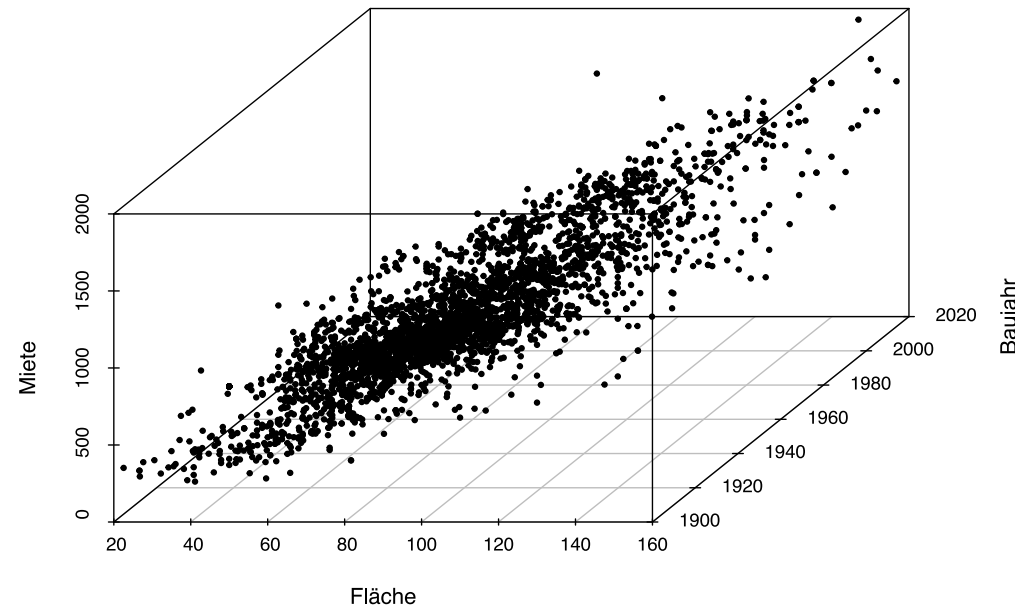
year	1.086e+00	7.452e-02	14.58	<2e-16

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

```
1 library(data.table)
2
3 # daten laden
4 daten = fread("../data/data2/rents.csv")
5
6 # zu fuß
7 X = cbind(rep(1,3854), daten$area, daten$year)
8 y = daten$rent
9 beta = solve(t(X) %*% X) %*% t(X) %*% y
10
11 print(beta)
```

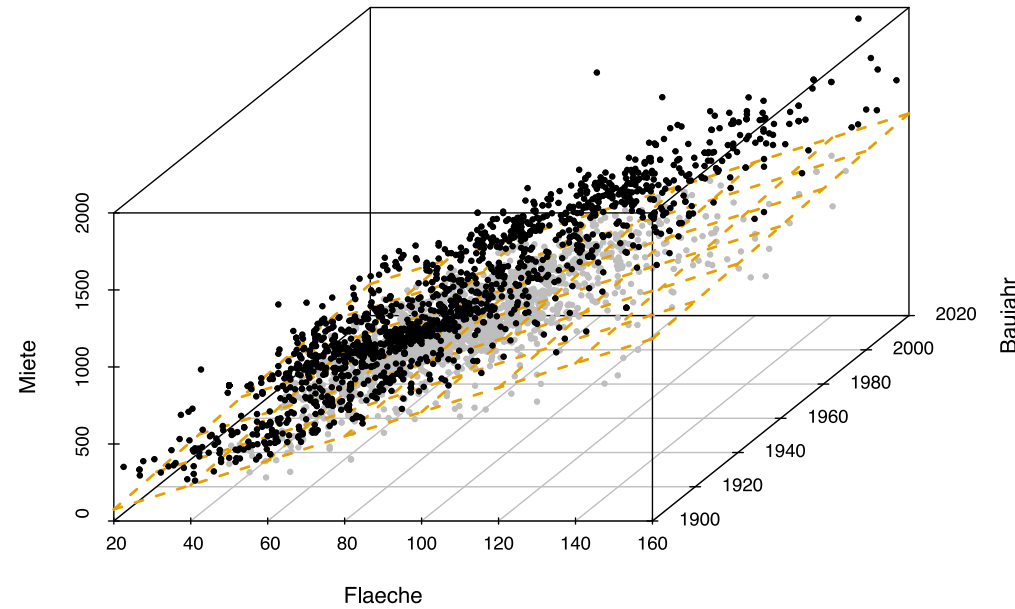
	[,1]
[1,]	-2146.724893
[2,]	7.906339
[3,]	1.086172

Angepasstes Modell im Mietspiegel-Beispiel



$$\text{Miete}_i = -2147 + 7.91 \cdot \text{Fläche}_i + 1.09 \cdot \text{Baujahr}_i + \epsilon_i$$

Angepasstes Modell im Mietspiegel-Beispiel



$$\text{Miete}_i = -2147 + 7.91 \cdot \text{Fläche}_i + 1.09 \cdot \text{Baujahr}_i + \epsilon_i$$

