

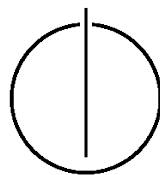


FAKULTÄT FÜR INFORMATIK  
DER TECHNISCHEN UNIVERSITÄT MÜNCHEN

Abschlussarbeit in Informatik

# **Effiziente statistische Methoden für Datenbanksysteme**

Thomas Heyenbrock







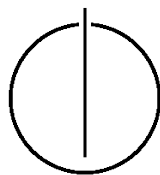
FAKULTÄT FÜR INFORMATIK  
DER TECHNISCHEN UNIVERSITÄT MÜNCHEN

Abschlussarbeit in Informatik

Effiziente statistische Methoden für Datenbanksysteme

Efficient statistical methods for database systems

Autor:	Thomas Heyenbrock
Aufgabensteller:	Prof. Alfons Kemper, Ph.D.
Betreuer:	Maximilian E. Schüle, M.Sc.
Datum:	15.02.2017





Ich versichere, dass ich diese Abschlussarbeit selbständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

München, den 12. Februar 2018

Thomas Heyenbrock



---

## Zusammenfassung

Das Ziel der vorliegenden Arbeit ist es, die Durchführung von statistischen Methoden in relationalen Datenbanksystemen zu demonstrieren. Das statistische Konzept, das dazu verwendet wird, ist die lineare und die logistische Regressionsanalyse. Dazu werden zuerst kurz die mathematischen Grundlagen erklärt. Danach wird die Umsetzung der Regressionsanalyse mit verschiedenen Programmiersprachen demonstriert, insbesondere in zwei relationalen Datenbanken. Diese Implementierungen werden daraufhin miteinander verglichen. Dabei stellt man fest, dass die lineare Regressionsanalyse bei kleinen Datenmengen in Datenbanken sehr gut funktioniert. Logistische Regression mit Gradientenverfahren kann dagegen nicht mehr mit anderen Implementierungen mithalten. Zum Abschluss wird ein mögliches Erweiterungspotenzial für relationale Datenbanksysteme erläutert, insbesondere eine Erweiterung um Matrix-Operationen.

## Abstract

The goal of this thesis is to demonstrate the implementation of statistical methods in relational database systems. The statistical concepts used are linear and logistic regression analysis. First, the mathematical basics are explained briefly. Then the implementation of the regression analysis will be demonstrated with different programming languages, especially in two relational databases. These implementations are then compared. It is found that the linear regression analysis works very well for small amounts of data in relational databases. On the other hand, logistic regression with gradient descent can no longer compete with other implementations. Finally, an expansion potential for relational database systems is explained, in particular an extension to matrix operations.





# Inhaltsverzeichnis

<b>Zusammenfassung</b>	<b>vii</b>
<b>1. Einführung und typische statistische Problemstellungen</b>	<b>1</b>
<b>2. Grundlagen statistischer Methoden</b>	<b>3</b>
2.1. Lineare Regression . . . . .	4
2.1.1. Einfache lineare Regression . . . . .	5
2.1.2. Multiple lineare Regression . . . . .	5
2.2. Logistische Regression . . . . .	6
2.2.1. Gradientenverfahren . . . . .	8
2.2.2. Gradient bei logistischer Regression . . . . .	8
<b>3. Anwendung statistischer Methoden</b>	<b>11</b>
3.1. Beispieldaten . . . . .	11
3.2. R-Projekt . . . . .	12
3.2.1. Grundprinzip . . . . .	12
3.2.2. Einfache lineare Regression . . . . .	12
3.2.3. Multiple lineare Regression . . . . .	13
3.2.4. Logistische Regression . . . . .	14
3.3. TensorFlow . . . . .	15
3.3.1. Grundprinzip . . . . .	15
3.3.2. Einfache lineare Regression . . . . .	16
3.3.3. Multiple lineare Regression . . . . .	17
3.3.4. Logistische Regression . . . . .	17
3.4. SQL . . . . .	18
3.4.1. Einfache lineare Regression . . . . .	18
3.4.2. Multiple lineare Regression . . . . .	19
3.4.3. Logistische Regression . . . . .	21
<b>4. Vergleich der verschiedenen Implementierungen</b>	<b>25</b>
4.1. Einfache lineare Regression . . . . .	25
4.2. Multiple lineare Regression . . . . .	26
4.3. Logistische Regression . . . . .	27
<b>5. Erweiterungspotenzial in Datenbanksystemen</b>	<b>29</b>
5.1. Einfache lineare Regression . . . . .	29
5.2. Multiple lineare Regression . . . . .	30
5.3. Logistische Regression . . . . .	31

<b>6. Fazit</b>	<b>33</b>
<b>Anhang</b>	<b>37</b>
<b>A. Python-Skript zum Generieren der Beispieldaten</b>	<b>37</b>
<b>B. R-Skripte</b>	<b>41</b>
B.1. Einfache lineare Regression . . . . .	41
B.2. Multiple lineare Regression . . . . .	43
B.3. Logistische Regression . . . . .	43
<b>C. TensorFlow-Skripte</b>	<b>47</b>
C.1. Einfache lineare Regression . . . . .	47
C.2. Multiple lineare Regression . . . . .	49
C.3. Logistische Regression . . . . .	52
<b>D. MySQL-Skripte</b>	<b>57</b>
D.1. Einfache lineare Regression . . . . .	57
D.2. Multiple lineare Regression . . . . .	58
D.3. Logistische Regression . . . . .	65
<b>E. PostgreSQL-Skripte</b>	<b>71</b>
E.1. Einfache lineare Regression . . . . .	71
E.2. Multiple lineare Regression . . . . .	72
E.3. Logistische Regression . . . . .	77
<b>F. Python-Skripte für das Benchmarking</b>	<b>83</b>
F.1. Berechnung der Benchmarks . . . . .	83
F.2. Auswertung der Benchmarks . . . . .	89
<b>Bibliografie</b>	<b>95</b>

# 1. Einführung und typische statistische Problemstellungen

Die Statistik ist ein Teilgebiet der Mathematik, in welchem Methoden zum Umgang und zur Verarbeitung von Daten behandelt werden. Dabei wird oft ein vorhandener Satz an Daten, auch Stichprobe genannt, betrachtet und analysiert, um daraus Vorhersagen für die Gesamtheit aller Daten zu treffen. Den Teilbereich der Statistik, welcher sich mit solchen Problemen befasst, nennt man induktive oder schließende Statistik.

Betrachtet man beispielsweise die Körpergröße und das Gewicht von 100 Testpersonen, dann kann man sich fragen, ob diese beiden Merkmale in Zusammenhang stehen. Insbesondere ist interessant, wie man einen möglichen Zusammenhang quantitativ darstellen kann und ob man neben der Körpergröße auch andere Faktoren für das Gewicht einer Person betrachten sollte. Das sind beispielhafte Typen von Fragen, welchen man in der Statistik oft begegnet.

Die Statistik zeigt Methoden und Vorgehensweisen auf, wie man solche Fragen angehen und beantworten kann. Ein oft verwendetes Verfahren ist die Regression bzw. die Regressionsanalyse. Hier sucht und analysiert man Beziehungen zwischen mehreren Variablen und versucht diese quantitativ zu beschreiben.

Greifen wir das obige Beispiel wieder auf: Bei der Regression sucht man nach einer Formel, welche für gegebene Körpergröße das Gewicht einer Person möglichst gut schätzt. Oft trifft man Annahmen über die Art der Beziehung zwischen den Dimensionen, um die Suche a priori einzugrenzen. Man beschränkt sich in vielen Fällen auf lineare Funktionen, da solche leicht zu behandeln sind. In der Praxis sind aber auch allgemeine Potenzfunktionen, exponentielle Funktionen oder logistische Funktionen bzw. Beziehungen oft anzutreffen.

Es gibt speziell für Statistik entwickelte Software, seien es einfache Programmiersprachen wie das R-Projekt oder komplexere Programme mit grafischem Interface wie SPSS von IBM. Die Daten, welche man als Basis für Analysen verwendet, müssen aber an einer anderen Stelle gespeichert und verwaltet werden. Oft liegen diese in einer Datenbank und müssen zuerst in das Analysetool importiert werden.

Konzeptuell ist eine Datenbank nicht für Regressionsanalyse geschaffen. Dennoch ist es in relationalen Datenbanksystemen mit SQL möglich, Regression direkt in der Datenbank durchzuführen. Damit übergeht man den eben genannten Schritt des Importierens. Man kann außerdem die Ergebnisse der Analyse direkt aus der Datenbank abfragen oder dort weiterverwenden.

Diese Arbeit wird zuerst das Konzept der Regression konkreter einführen und die mathematischen Grundlagen darlegen. Darauf aufbauend betrachten wir Implementierungen für Regressionsanalyse in verschiedenen Programmiersprachen. Hier soll insbesondere die Anwendung von Regressionsanalyse mit Hilfe von SQL demonstriert werden. Danach wollen wir die Sprachen bezüglich der Laufzeit miteinander vergleichen und noch kurz auf das Erweiterungspotenzial für relationale Datenbanksysteme eingehen.



## 2. Grundlagen statistischer Methoden

Bei der Regressionsanalyse geht es im Allgemeinen darum, das Verhalten einer Größe  $Y$  in Abhängigkeit einer oder mehrerer anderer Größen  $X_1, X_2, \dots, X_n$  zu modellieren. Die Größe  $Y$  wird abhängig genannt, die Größen  $X_j$  nennt man unabhängig. Für diese Arbeit wollen wir zunächst einige Annahmen treffen. Diese sollen immer gelten, falls nicht explizit etwas anderes festgelegt wird.

- Die genannten Größen sind Zufallsvariablen. Eine solche Zufallsvariable ist eine Funktion, deren Werte die Ergebnisse eines Zufallsvorgangs darstellt.
- Die Zufallsvariablen sind auf der Menge  $M = \{1, \dots, m\}$  definiert und bilden in die reellen Zahlen ab:

$$Y : M \rightarrow \mathbb{R}, \quad X_1 : M \rightarrow \mathbb{R}, \quad \dots, \quad X_n : M \rightarrow \mathbb{R}$$

Die Zufallsvariablen sind also metrisch skaliert. Die  $m$  Zahlen in der Menge  $M$  entsprechen den  $m$  Datenpunkten, die wir als Datenbasis für die Regressionsanalyse besitzen.

- Wir verwenden die folgenden Abkürzungen für die Werte der Zufallsvariablen:

$$\begin{aligned} y_i &:= Y(i) \quad \text{für alle } i \in M, \\ x_{i,j} &:= X_j(i) \quad \text{für alle } i \in M \text{ und } 1 \leq j \leq n \end{aligned}$$

Nun definieren wir ein Modell, mit dem der Zusammenhang zwischen der abhängigen und den unabhängigen Variablen dargestellt werden soll. Dazu verwenden wir eine Funktion  $f$ , welche für Werte von  $X_1$  bis  $X_n$  einen geschätzten Wert für  $Y$  liefert. Idealerweise existiert eine Funktion, die zum einen eine einfache Darstellung (etwa durch eine arithmetische Formel) besitzt und zum anderen alle unabhängigen Werte der Datenmenge exakt prognostiziert. Das bedeutet:

$$y_i = f(x_{i,1}, \dots, x_{i,n}) \quad \text{für alle } 1 \leq i \leq m$$

Im Allgemeinen ist es nicht möglich, eine solche Funktion zu finden. Man verwirft also die Anforderung der Exaktheit für alle Datenpunkte und versucht stattdessen eine einfache Funktion zu finden, die die Datenmenge möglichst gut approximiert. Wir definieren für jeden Datenpunkt den Fehler  $e_i$ , der sich durch die Ungenauigkeit der Modellfunktion  $f$  ergibt:

$$e_i := y_i - f(x_{i,1}, \dots, x_{i,n})$$

Je näher ein Fehlerterm bei null liegt, desto besser ist die Annäherung für den jeweiligen Datenpunkt. Um eine gute Approximation für die gesamte Datenmenge zu erhalten, sollten

die Fehlerterme global betrachtet möglichst klein bleiben. Das Ziel der Regressionsanalyse ist nun die Bestimmung der Funktion  $f$ . Dazu nimmt man an, dass  $f$  eine bestimmte Form hat. Diese Annahme kann sich je nach der Problemstellung unterscheiden. Oft arbeitet man mit linearen Funktionen  $f$ , aber auch quadratische, exponentielle oder logistische Funktionen sind nicht unüblich.

## 2.1. Lineare Regression

Bei der linearen Regression geht man von einem linearen Zusammenhang zwischen der abhängigen und den unabhängigen Variablen aus. Die Funktion  $f$  ist also von folgender Form:

$$f(x_1, \dots, x_n) = \alpha + \sum_{j=1}^n \beta_j \cdot x_j \quad \text{mit } \beta_j \in \mathbb{R}$$

Dabei sind  $\alpha$  und  $\beta_k$  für  $k = 1, \dots, n$  reelle Zahlen, die sogenannten Parameter der Funktion. Das Maß für die Qualität von  $f$  ist die Summe der quadrierten Fehlerterme. Diese Summe kann wiederum als Funktion in Abhängigkeit der Parameter definiert werden:

$$E(\alpha, \beta_1, \dots, \beta_n) := \sum_{i=1}^m e_i^2 = \sum_{i=1}^m (y_i - f(x_{i,1}, \dots, x_{i,n}))^2 = \sum_{i=1}^m \left( y_i - \alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right)^2$$

Bei der linearen Regression sucht man die Parameter  $\hat{\alpha}$  und  $\hat{\beta}_k$  für die  $E$  einen minimalen Wert annimmt:

$$E(\hat{\alpha}, \hat{\beta}_1, \dots, \hat{\beta}_n) = \min \{ E(\alpha, \beta_1, \dots, \beta_n) \mid \alpha \in \mathbb{R}, \beta_1 \in \mathbb{R}, \dots, \beta_n \in \mathbb{R} \}$$

Um dieses Minimierungsproblem zu lösen, berechnen wir die partiellen Ableitungen von  $E$  nach allen Parametern:

$$\begin{aligned} \frac{\partial E}{\partial \alpha} &= \sum_{i=1}^m 2 \cdot \left( y_i - \alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \cdot (-1) \\ &= -2 \cdot \sum_{i=1}^m \left( y_i - \alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \\ \frac{\partial E}{\partial \beta_k} &= \sum_{i=1}^m 2 \cdot \left( y_i - \alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \cdot (-x_{i,k}) \\ &= -2 \cdot \sum_{i=1}^m x_{i,k} \cdot \left( y_i - \alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \end{aligned}$$

Durch Nullsetzen der partiellen Ableitungen erhält man ein lineares Gleichungssystem mit  $(n+1)$  Gleichungen und ebenso vielen Unbekannten.

$$\frac{\partial E}{\partial \alpha} = 0, \quad \frac{\partial E}{\partial \beta_1} = 0, \quad \dots, \quad \frac{\partial E}{\partial \beta_n} = 0$$

Die Lösung dieses Gleichungssystems (falls diese existiert) ist das gesuchte Minimum. Damit findet man die gesuchten Parameter für die lineare Funktion  $f$ .

### 2.1.1. Einfache lineare Regression

Man spricht von einfacher linearer Regression, wenn man mit nur einer unabhängigen Variable arbeitet. Anschaulich möchte man hier die bestmögliche Schätzgerade durch eine gegebene Punktwolke legen.

Wir nennen die unabhängige Variable in diesem Kapitel statt  $X_1$  einfach nur  $X$ . Ebenso schreiben wir  $\beta := \beta_1$  und  $x_i := x_{i,1}$ . Dann können wir das lineare Gleichungssystem zum Auffinden des Minimums wie folgt aufschreiben:

$$\begin{aligned} 0 &= -2 \cdot \sum_{i=1}^m (y_i - \alpha - \beta \cdot x_i) \\ 0 &= -2 \cdot \sum_{i=1}^m x_i \cdot (y_i - \alpha - \beta \cdot x_i) \end{aligned}$$

Dieses Gleichungssystem kann explizit gelöst werden. Man erhält das folgende Ergebnis für die Lösung  $\hat{\alpha}$  und  $\hat{\beta}$ :

$$\begin{aligned} \hat{\beta} &= \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2} \\ \hat{\alpha} &= \bar{y} - \hat{\beta}\bar{x} \end{aligned}$$

Dabei bezeichnen  $\bar{x}$  und  $\bar{y}$  die Mittelwerte von  $X$  respektive  $Y$ , also:

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i, \quad \bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$$

Eine Herleitung dieser Lösung findet sich in Kapitel 3.6.2 in [2].

### 2.1.2. Multiple lineare Regression

Bei multipler linearer Regression existieren mindestens zwei unabhängige Variablen. Anstatt wieder explizite Formeln für jeden einzelnen Parameter anzugeben, berechnen wir alle gesuchten Parameter gleichzeitig mit Hilfe von Matrizenrechnung. Definieren wir dazu die folgenden Matrizen und Vektoren:

$$\begin{aligned} X &= \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m,1} & \dots & x_{m,n} \end{pmatrix} \in \mathbb{R}^{m \times (n+1)} \\ y &= \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^{m \times 1}, \quad b = \begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_n \end{pmatrix} \in \mathbb{R}^{(n+1) \times 1} \end{aligned}$$

Dabei ist  $b$  der Vektor mit den gesuchten Parametern für die Minimierung der kleinsten Quadrate bzw. der Funktion  $E$ . Falls die Matrix  $X^T X$  invertierbar ist, gilt die folgende Formel für die Berechnung der gesuchten Parameter:

$$b = (X^T X)^{-1} X^T y$$

Mehr dazu findet man auch in Kapitel 12.2.3 in [2].

### 2.2. Logistische Regression

Die logistische Regression findet Anwendung im Falle, dass die abhängige Variable eine binäre Variable ist, also eine Variable, die nur zwei Werte annehmen kann. Oft handelt es sich um eine Eigenschaft oder einen Gegenstand, den man entweder besitzt oder nicht, wie zum Beispiel das Geschlecht einer Person, ein Premium-Abonnement für einen Web-Service oder der Besitz eines Autos. Wir bezeichnen die beiden möglichen Werte einer solchen Variablen mit 0 und 1. Die Zuordnung vom Merkmal zur Zahl ist frei wählbar.

Lineare Regression eignet sich nicht zur Modellierung einer binären Variablen, da eine lineare Funktion entweder konstant oder unbeschränkt ist, in zweiten Fall also insbesondere Werte größer als 1 und kleiner als 0 annimmt. Um diesem Problem abzuweichen, wählen wir zusätzlich zu der linearen Funktion eine weitere Funktion, die beliebige Zahlen auf das Intervall  $[0, 1]$  abbildet. Im Falle der logistischen Regression verwendet man die gleichnamige logistische Funktion:

$$l : \mathbb{R} \rightarrow (0, 1), \quad x \mapsto \frac{1}{1 + e^{-x}}$$

Diese Funktion wendet man nun auf die Linearkombination aller unabhängigen Variablen mit Parametern  $\beta_1$  bis  $\beta_n$  und konstantem Term  $\alpha$  an. Zur Vereinfachung definieren wir für das restliche Kapitel die Variable  $c_i$  wie folgt:

$$c_i := \alpha + \sum_{j=1}^n \beta_j \cdot x_{i,j}$$

Das Ergebnis der Funktion  $l$  für den  $i$ -ten Datensatz bezeichnen wir mit  $\pi_i$ . Dieses ist wieder eine Funktion in Abhängigkeit der Parameter:

$$\pi_i = \pi_i(\alpha, \beta_1, \dots, \beta_n) := l\left(\alpha + \sum_{j=1}^n \beta_j \cdot x_{i,j}\right) = \frac{1}{1 + e^{-c_i}}$$

Wir stellen hierbei fest, dass folgende Identität für die Funktionen  $\pi_i$  gilt:

$$\begin{aligned} \pi_i(-\alpha, -\beta_1, \dots, -\beta_n) &= \frac{1}{1 + e^{c_i}} = \frac{1 + e^{c_i} - e^{c_i}}{1 + e^{c_i}} \\ &= 1 - \frac{e^{c_i}}{1 + e^{c_i}} = 1 - \frac{1}{e^{-c_i} + 1} \\ &= 1 - \pi_i(\alpha, \beta_1, \dots, \beta_n) \end{aligned}$$



Wir interpretieren  $\pi_i$  als die Wahrscheinlichkeit dafür, dass der Wert der abhängigen Variable eines Datensatzes mit Werten  $x_{i,1}, \dots, x_{i,n}$  der unabhängigen Variablen gleich 1 ist, also:

$$\pi_i = P(Y_i = 1 | X_1 = x_{i,1}, \dots, X_n = x_{i,n})$$

Man möchte die Parameter  $\alpha$  und  $\beta_k$  nun so schätzen, dass die Wahrscheinlichkeit für das Auftreten der vorhandenen Datenbasis maximiert wird. Diese Wahrscheinlichkeit ist gegeben durch:

$$\begin{aligned} L(\alpha, \beta_1, \dots, \beta_n) &= \prod_{i=1}^m P(Y_i = y_i | X_1 = x_{i,1}, \dots, X_n = x_{i,n}) \\ &= \prod_{i=1}^m y_i \cdot \pi_i(\alpha, \beta_1, \dots, \beta_n) + (1 - y_i) \cdot (1 - \pi_i(\alpha, \beta_1, \dots, \beta_n)) \\ &= \prod_{i=1}^m y_i \cdot \pi_i(\alpha, \beta_1, \dots, \beta_n) + (1 - y_i) \cdot \pi_i(-\alpha, -\beta_1, \dots, -\beta_n) \end{aligned}$$

Da alle  $y_i$  entweder gleich 0 oder gleich 1 sind, ist immer nur einer der beiden Summanden in obigem Produkt nicht null. Diese Fallunterscheidung kann man auch in das Vorzeichen der Parameter verschieben, da sich die beiden möglichen Faktoren nur darin unterscheiden. Es gilt also:

$$L(\alpha, \beta_1, \dots, \beta_n) = \prod_{i=1}^m \pi_i((2y_i - 1)\alpha, (2y_i - 1)\beta_1, \dots, (2y_i - 1)\beta_n)$$

Das Verfahren der Maximierung dieser Wahrscheinlichkeit bezeichnet man auch als Maximum-Likelihood-Methode. Die Funktion  $L$  nennt man Likelihoodfunktion. Oft maximiert man nicht  $L$  direkt, sondern eher den natürlichen Logarithmus von  $L$ :

$$\begin{aligned} L_{\log}(\alpha, \beta_1, \dots, \beta_n) &:= \ln(L(\alpha, \beta_1, \dots, \beta_n)) \\ &= \sum_{i=1}^m \ln \left( \pi_i((2y_i - 1)\alpha, (2y_i - 1)\beta_1, \dots, (2y_i - 1)\beta_n) \right) \end{aligned}$$

Der Sinn ist, dass man das Produkt damit in eine Summe einzelner Logarithmen umwandeln kann, welche wiederum einfacher abzuleiten ist. Die Maximierung von  $L$  ist äquivalent mit der von  $L_{\log}$ , da der Logarithmus eine stetig wachsende Funktion ist. Die Werte von  $L$  liegen stets zwischen 0 und 1, also ist  $L_{\log}$  wohldefiniert.

Um dieses Regressionsproblem zu lösen, muss man also die Parameter  $\hat{\alpha}$  und  $\hat{\beta}_k$  finden, für die gilt:

$$L(\hat{\alpha}, \hat{\beta}_1, \dots, \hat{\beta}_n) = \max \{ L(\alpha, \beta_1, \dots, \beta_n) \mid \alpha \in \mathbb{R}, \beta_1 \in \mathbb{R}, \dots, \beta_n \in \mathbb{R} \}$$

In diesem Fall kommt man nicht mehr an einer iterativen Lösung vorbei, da das lineare Gleichungssystem aus den partiellen Ableitungen nicht mehr exakt lösbar ist. Eine der einfachsten Methoden zur Lösung von Optimierungsproblemen ist das Gradientenverfahren, welches im kommenden Teilkapitel kurz eingeführt wird. Danach wird gezeigt, wie man das Gradientenverfahren für logistische Regression anwendet.

### 2.2.1. Gradientenverfahren

Das Gradientenverfahren ist ein iterativer Algorithmus zur Lösung von Optimierungsproblemen. Nachdem wir hier bei der logistischen Regression eine Funktion maximieren wollen, führen wir das Gradientenverfahren dementsprechend als Maximierungsalgorithmus ein. Man kann dasselbe Verfahren aber auch zur Lösung von Minimierungsproblemen einsetzen. Gegeben sei eine Funktion der folgenden Form, die maximiert werden soll:

$$f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}, \quad (\alpha, \beta_1, \dots, \beta_n) \mapsto f(\alpha, \beta_1, \dots, \beta_n)$$

Beim Gradientenverfahren beginnt man mit beliebigen Startwerten  $\alpha^0$  und  $\beta_k^0$  und einer Schrittweite  $s \in \mathbb{R}^+$ . Vom Startpunkt aus geht man nun etwas in die Richtung des steilsten Anstieges der Funktion und erhält dadurch neue Werte  $\alpha^1$  und  $\beta_k^1$ . Diese Richtung ist der Gradient der Funktion  $f$ .

Der Gradient ist ein Vektor, der sich aus den partiellen Ableitungen von  $f$  nach jeweils einer Variablen zusammensetzt und wird wie folgt notiert:

$$\nabla(f) = \begin{pmatrix} \partial f / \partial \alpha \\ \partial f / \partial \beta_1 \\ \vdots \\ \partial f / \partial \beta_n \end{pmatrix}$$

Der Gradient von  $f$  ist wieder eine Funktion, die Werte  $\alpha$  und  $\beta_k$  auf einen Vektor der Länge  $n + 1$  abbildet. Der iterative Schritt des Verfahrens definiert sich wie folgt:

$$\begin{pmatrix} \alpha^{i+1} \\ \beta_1^{i+1} \\ \vdots \\ \beta_n^{i+1} \end{pmatrix} = \begin{pmatrix} \alpha^i \\ \beta_1^i \\ \vdots \\ \beta_n^i \end{pmatrix} + s \cdot \nabla(f)(\alpha^i, \beta_1^i, \dots, \beta_n^i)$$

Will man ein Minimierungsproblem lösen, muss man nur das Vorzeichen des Gradienten vertauschen, also ein Minus statt einem Plus in der obigen Formel verwenden. Man geht also entgegengesetzt der Richtung des steilsten Anstiegs und damit in die Richtung des steilsten Abstiegs von  $f$ .

Danach muss noch getestet werden, dass  $L$  für die neuen Parameter auch wirklich einen größeren Wert annimmt also zuvor. Falls nicht, muss die Schrittweite  $s$  verkleinert werden, zum Beispiel um einen festen, zuvor definierten Faktor.

### 2.2.2. Gradient bei logistischer Regression

Um das Gradientenverfahren bei logistischer Regression einsetzen zu können, muss der Gradient für den Logarithmus der Likelihoodfunktion bekannt sein. In diesem Kapitel berechnen wir die partiellen Ableitungen nach allen Parametern.

Um  $L_{\log} = \ln(L)$  partiell ableiten zu können, berechnen wir zuerst die partiellen Ableitungen aller  $\pi_i$ . Dazu notieren wir die natürliche Exponentialfunktion als  $\exp$ . Für die partielle

Ableitung nach  $\alpha$  ergibt sich mit der Kettenregel folgende Funktion:

$$\begin{aligned}\frac{\partial \pi_i}{\partial \alpha} &= - \left( 1 + \exp \left( -\alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \right)^{-2} \cdot \exp \left( -\alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \cdot (-1) \\ &= \left( 1 + \exp \left( -\alpha - \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \right)^{-1} \cdot \left( 1 + \exp \left( \alpha + \sum_{j=1}^n \beta_j \cdot x_{i,j} \right) \right)^{-1} \\ &= \pi_i(\alpha, \beta_1, \dots, \beta_n) \cdot \pi_i(-\alpha, -\beta_1, \dots, -\beta_n)\end{aligned}$$

Die partiellen Ableitungen nach einem der  $\beta_k$  für  $k = 1, \dots, n$  kann fast analog gebildet werden. Bei der Anwendung der Kettenregel auf die innerste lineare Funktion bleibt jedoch noch der konstante Faktor  $x_{i,k}$  übrig.

$$\frac{\partial \pi_i}{\partial \beta_k} = x_{i,k} \cdot \pi_i(\alpha, \beta_1, \dots, \beta_n) \cdot \pi_i(-\alpha, -\beta_1, \dots, -\beta_n)$$

Bevor wir  $L_{log}$  ableiten, definieren wir Hilfsvariablen  $\tilde{\alpha} := (2y_i - 1)\alpha$  und  $\tilde{\beta}_k := (2y_i - 1)\beta_k$  für  $k = 1, \dots, n$ . Damit erhält man:

$$\begin{aligned}L_{log}(\alpha, \beta_1, \dots, \beta_n) &= \ln(L(\alpha, \beta_1, \dots, \beta_n)) \\ &= \sum_{i=1}^m \ln \left( \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n) \right)\end{aligned}$$

Leitet man nach  $\alpha$  ab, so erhält man:

$$\begin{aligned}\frac{\partial L_{log}}{\partial \alpha} &= \sum_{i=1}^m \frac{\partial}{\partial \alpha} \left( \ln \left( \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n) \right) \right) \\ &= \sum_{i=1}^m \left( \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n) \right)^{-1} \cdot \frac{\partial \pi_i}{\partial \tilde{\alpha}} \cdot \frac{\partial \tilde{\alpha}}{\partial \alpha} \\ &= \sum_{i=1}^m \left( \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n) \right)^{-1} \cdot \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n) \cdot (1 - \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n)) \cdot (2y_i - 1) \\ &= \sum_{i=1}^m (1 - \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n)) \cdot (2y_i - 1)\end{aligned}$$

Für die partielle Ableitung nach  $\beta_k$  erhält man analog:

$$\frac{\partial L_{log}}{\partial \beta_k} = \sum_{i=1}^m x_{i,k} \cdot (1 - \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n)) \cdot (2y_i - 1)$$

Wir betrachten nun die einzelnen Summanden der partiellen Ableitungen getrennt für die beiden möglichen Werten von  $y_i$ . Ist  $y_i = 0$  dann gilt:

$$\begin{aligned}(1 - \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n)) \cdot (2y_i - 1) &= (1 - \pi_i(-\alpha, -\beta_1, \dots, -\beta_n)) \cdot (-1) \\ &= -1 + (1 - \pi_i(\alpha, \beta_1, \dots, \beta_n)) \\ &= -\pi_i(\alpha, \beta_1, \dots, \beta_n)\end{aligned}$$

Für  $y_i = 1$  ergibt sich Folgendes:

$$\begin{aligned}(1 - \pi_i(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_n)) \cdot (2y_i - 1) &= (1 - \pi_i(\alpha, \beta_1, \dots, \beta_n)) \cdot (2 - 1) \\ &= 1 - \pi_i(\alpha, \beta_1, \dots, \beta_n)\end{aligned}$$

Damit können wir die partiellen Ableitungen weiter vereinfachen:

$$\begin{aligned}\frac{\partial L_{log}}{\partial \alpha} &= \sum_{i=1}^m y_i - \pi_i(\alpha, \beta_1, \dots, \beta_n) \\ \frac{\partial L_{log}}{\partial \beta_k} &= \sum_{i=1}^m x_{i,k} \cdot (y_i - \pi_i(\alpha, \beta_1, \dots, \beta_n))\end{aligned}$$

Diese Darstellung der partiellen Ableitungen verwenden wir später in SQL zur Berechnung des Gradienten. Der Term innerhalb der Summe wird für jedes Tupel der Relation der Datenpunkte berechnet. Danach wird das resultierende Attribut mit einer Gruppierung summiert.

## 3. Anwendung statistischer Methoden

In diesem Kapitel werden mehrere Programmiersprachen vorgestellt, die sich für Regressionsanalyse eignen. Dabei betrachten wir für Statistik entwickelte Programmiersprachen ebenso wie Softwarebibliotheken für maschinelles Lernen. Im letzten Teil dieses Kapitels wird demonstriert, wie man Regression mit vorhandener SQL-Syntax umsetzen und durchführen kann.

### 3.1. Beispieldaten

Wir arbeiten in dieser Arbeit mit Beispieldaten, welche mit einem Python-Skript erstellt wurden. Das Skript ist als Ganzes im Anhang A zu finden. Die Daten liegen in Form einer csv-Datei vor, welche in so gut wie jeder Sprache einfach eingelesen werden kann. Außerdem wird der Datensatz auch als INSERT-Abfrage in eine sql-Datei geschrieben, mit der die Daten in einem Datenbanksystem in eine Relation geschrieben werden können.

Wir betrachten hier fiktive Kunden eines Onlinehandels. Für jeden Kunden wissen wir sein Alter, die Anzahl seiner Käufe, die Summe des ausgegebenen Geldes und ob der Kunde eine Premium-Mitgliedschaft besitzt oder nicht. Der ausgegebene Betrag wird in Cent angegeben, um mit ganzen Zahlen rechnen zu können. Die Premium-Mitgliedschaft wird mit einer 1 symbolisiert, während eine 0 das Gegenteil bedeutet.

Insgesamt wurden für diese Arbeit 100000 solcher Datensätze erzeugt. Die später dargestellten Ergebnisse wurden mit den ersten 1000 Datenpunkten berechnet. In der folgenden Tabelle sind die ersten 10 dieser Datensätze dargestellt.

Tabelle 3.1.: Auszug aus den Beispieldaten

<b>age</b>	<b>purchases</b>	<b>money</b>	<b>premium</b>
30	1	4421	0
30	11	23346	1
33	1	4010	0
31	19	52517	1
29	3	8046	0
28	12	25295	0
41	16	38236	1
23	3	7098	1
25	1	2707	0
38	20	50976	1

Wir wählen außerdem drei Fragestellungen, welche wir mit den in Kapitel 2 vorgestellten Arten von Regressionen beantworten werden:

1. Als Erstes wollen wir wissen, ob das ausgegebene Geld mit der Anzahl der Käufe in linearem Zusammenhang steht. Diese Frage können wir mit einfacher linearer Regression beantworten. *money* ist hierbei die abhängige Variable und *purchases* ist die unabhängige Variable.
2. Die zweite Frage ist ähnlich der Ersten. Hier wollen wir wissen, ob neben der Anzahl der Käufe auch das Alter des Kunden einen Einfluss auf das ausgegebene Geld hat. Hier haben wir nun zwei unabhängige Variablen, nämlich *age* und *purchases*. Die abhängige Variable bleibt *money*. Diese Frage beantworten wir mit multipler linearer Regression.
3. Zuletzt interessiert uns, ob eine Premium-Mitgliedschaft mit der Summe des ausgegebenen Geldes zusammenhängt. *money* ist also nun die unabhängige Variable und *premium* ist die abhängige Variable. Nachdem *premium* eine binäre Variable ist, nutzen wir hier logistische Regression.

## 3.2. R-Projekt

Das R-Projekt oder einfach nur R ist eine Sprache für statistische Berechnungen und auch grafische Darstellung. Damit ist R wie geschaffen für Regressionsanalyse. Die für diese Arbeit erstellten R-Skripte sind im Anhang B.

### 3.2.1. Grundprinzip

In R sind Datenstrukturen wie Vektoren, Matrizen und Listen als Datentypen vorhanden. Darauf aufbauend existieren sogenannte Dataframes. Ein Dataframe ist eine Liste von Vektoren gleicher Länge und wird in R zur Repräsentation von Datentabellen verwendet. Die Vektoren der Liste entsprechen den Attributen einer Relation. Wir importieren zuerst immer die Beispieldaten aus der csv-Datei in einen solchen Dataframe.

In R lassen sich außerdem sogenannte Modelle definieren, welche als Eingabe nur die zugehörigen Daten und eine Formel benötigen. Eine Formel ist ein Zeichenfolge der Form "*y ~ modell*" und symbolisiert den funktionalen Zusammenhang zwischen der abhängigen und den unabhängigen Variablen.

Hat man ein Modell erstellt, so bietet R Funktionen, um die Parameter für das definierte Modell mittels Regressionsanalyse zu berechnen. Wir werden im Folgenden von den Funktionen *lm* (für "linear model") und *glm* (für "generalized linear model") Gebrauch machen. Die Implementierungen der beiden Funktionen wurden an die gleichnamigen Funktionen in der verwandten Programmiersprache S angelehnt. Mehr dazu findet man in [3].

### 3.2.2. Einfache lineare Regression

Betrachten wir also die erste Frage aus dem vorherigen Teilkapitel. Die Formel lautet hier "*money ~ purchases*". Man liest also die Daten aus der csv-Datei, erstellt das Modell mit der genannten Formel und berechnet die Parameter mit der Funktion *lm*. Mit nur drei

Zeilen Code lässt sich also die komplette Regressionsanalyse implementieren. Die *print*-Funktion dient zur Ausgabe des Ergebnisses.

```

1 data <- read.csv2("sample.csv", sep = ",", header = TRUE)
2 modell <- as.formula("money ~ purchases")
3 slr <- lm(modell, data = data)
4 print(slr)

```

Das Ergebnis des obigen Codes ist Folgendes:

```

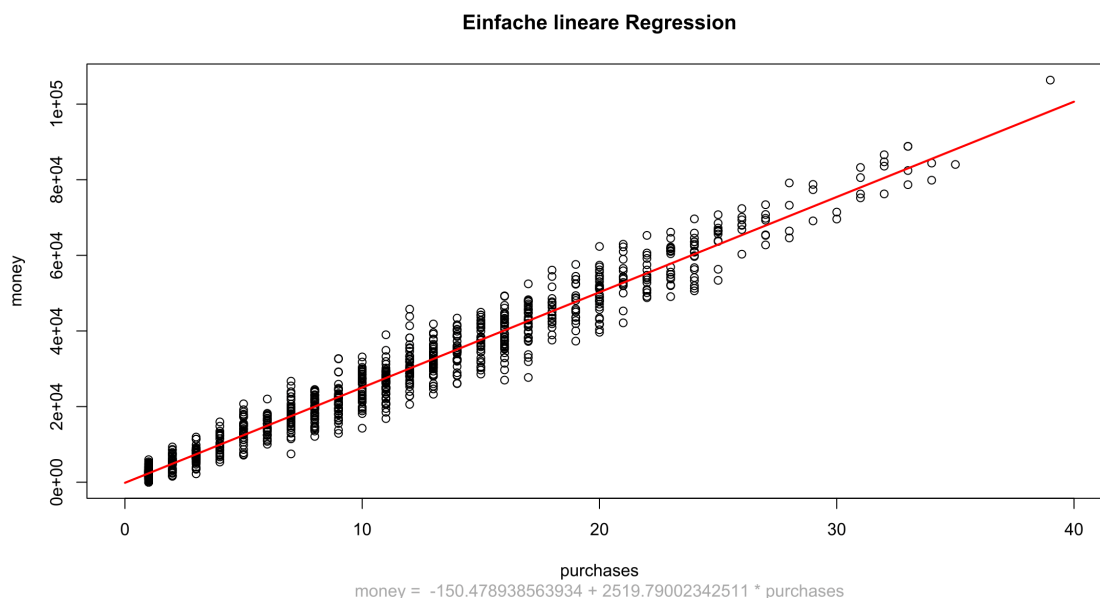
1 Coefficients:
2 (Intercept)    purchases
3      -150.5      2519.8

```

Der Wert unter "*Intercept*" entspricht dabei dem Parameter  $\alpha$  in unserer Notation, der Wert unter "*purchases*" entspricht  $\beta$ .

Der relativ kleine Wert für  $\alpha$  entspricht der Intuition, dass ein Kunde ohne Käufe auch kein Geld ausgegeben hat. Der relativ große Wert von ca. 2520 für  $\beta$  zeigt, dass die Anzahl der gekauften Artikel sehr einen großen Einfluss auf das ausgegebene Geld hat. Ein Kunde gibt pro gekauftem Artikel etwa 2520 Cent, also 25.20 Euro aus.

R verfügt auch über Möglichkeiten zur grafischen Darstellung. Lässt man die Datenpunkte und die lineare Ausgleichsfunktion mit den berechneten Parametern plotten, erhält man dieses Diagramm:



### 3.2.3. Multiple lineare Regression

Bei multipler linearer Regression unterscheidet sich der R-Code nur in der Wahl der Formel von dem Code aus dem vorherigen Teilkapitel. Hier wollen wir *money* durch eine

lineare Summe aus *purchases* und *age* modellieren, deshalb lautet die Formel hier "*money ~ purchases + age*". Man erhält das folgende Ergebnis:

1	Coefficients:		
2	(Intercept)	purchases	age
3	-766.7	2520.3	17.5

Der Wert für das  $\beta$  zu *purchases* ist fast exakt derselbe wie bei einfacher linearer Regression, was bei denselben Daten auch zu erwarten war. Der Wert für das  $\beta$  zu *age* ist dagegen nahe bei null. Das bedeutet, dass das Alter im Vergleich zu der Anzahl der Käufe keinen signifikanten Einfluss auf das ausgegebene Geld hat.

#### 3.2.4. Logistische Regression

Bei logistischer Regression nutzen wir nicht mehr wie bisher ein lineares Modell, sondern ein generalisiertes lineares Modell. Logistische Regression ist im Wesentlichen ein Spezialfall dieses Modelles. Dazu nutzen wir die *glm*-Funktion in R. Um logistische Regression damit betreiben zu können, wählt man den Parameter *family* dieser Funktion als *binomial*.

Man braucht wie bei linearer Regression eine Formel für das Modell. Diese bildet man analog wie bisher, indem mal die abhängige Variable mit den unabhängigen Variablen über eine Tilde verbindet. Im Fall der dritten Fragestellung wählt man die Formel als "*premium ~ money*".

Der gesamte R-Code für die logistische Regression ist wieder ähnlich zu dem Code aus 3.2.2 und lautet wie folgt:

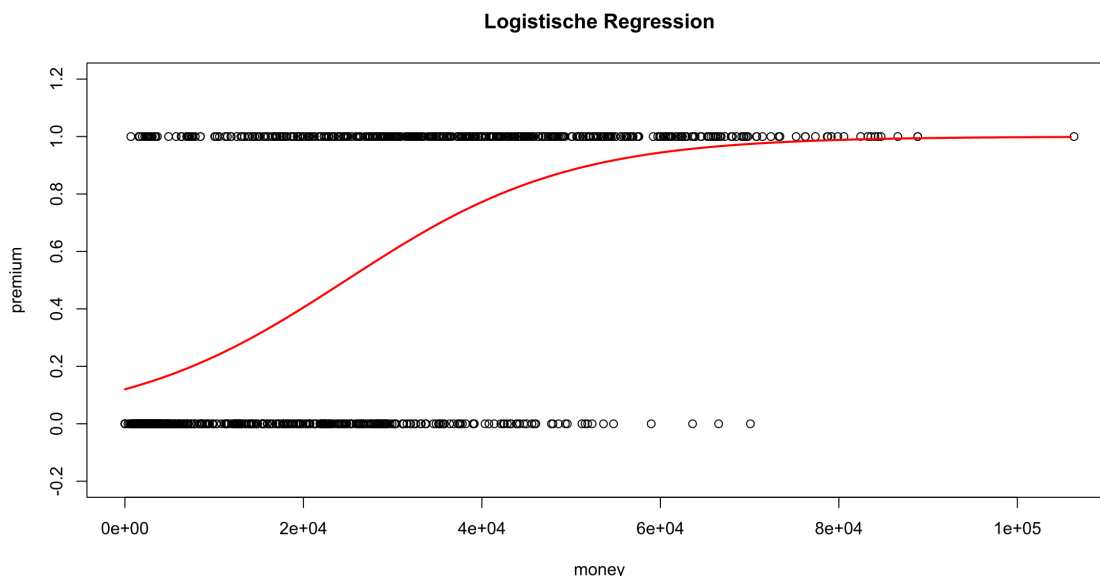
```
1 data <- read.csv2("sample.csv", sep = ",", header = TRUE)
2 modell <- as.formula("premium ~ money")
3 logit <- glm(modell, family = binomial, data = data)
4 print(logit)
```

Nach der Ausführung erhält man das folgende Ergebnis:

1	Coefficients:	
2	(Intercept)	money
3	-1.9910911	0.0000803
4		
5	Degrees of Freedom: 999 Total (i.e. Null); 998 Residual	
6	Null Deviance:	1386
7	Residual Deviance: 1006	AIC: 1010

Eine anschauliche Interpretation der zurückgegebenen Parameter ist nicht mehr so einfach. Wir lassen uns das Ergebnis daher wieder als Plot visualisieren:





Für Kunden, die weniger als 100 Euro ausgegeben haben ist die Wahrscheinlichkeit Premium-Mitglied zu sein mit etwa 25% relativ gering. Je höher der Wert für *money* aber wird, desto größer wird auch die genannte Wahrscheinlichkeit. So ist ein Kunde mit mehr als 800 Euro Ausgaben so gut wie immer ein Premium-Mitglied.

### 3.3. TensorFlow

TensorFlow ist eine Softwarebibliothek, die von Google für die Umsetzung von Algorithmen für maschinelles Lernen entwickelt wurde. Das umfasst insbesondere auch die Möglichkeit zur iterativen Optimierung von Kostenfunktionen, was wir nun zur Regressionsanalyse nutzen wollen.

TensorFlow bietet APIs für verschiedene Programmiersprachen an. Die Skripte, welche für diese Arbeit erstellt wurden, sind in Python geschrieben. Wir verwenden TensorFlows Implementierung eines Gradientenabstiegsverfahrens, für welches man die Anzahl der Schritte und die Abstiegsgeschwindigkeit selbst wählen muss. Die Genauigkeit des Ergebnisses hängt folglich von einer angemessenen Auswahl dieser Werte ab.

Die Python-Skripte umfassen nun zwischen 70 und 100 Codezeilen, daher findet man diese Skripte als Ganzes nur im Anhang C. Die wichtigsten Ausschnitte aus dem Code sollen in den folgenden Teilkapiteln aber einen Einblick in die Funktionsweise des Codes geben.

#### 3.3.1. Grundprinzip

TensorFlow arbeitet mit Tensoren als grundlegende Datenstruktur. Solche sind im Wesentlichen mehrdimensionale Matrizen mit festen Dimensionen. Tensoren können mit Hilfe verschiedenster Operatoren weiterverarbeitet werden.

Es gibt drei Möglichkeiten, Tensoren zu definieren, nämlich als Konstante, als Variable oder als Platzhalter. Während Konstanten ihren Wert nach der Initialisierung nicht mehr ändern können, sind die Werte der beiden Letztgenannten veränderbar. Der Unterschied

besteht darin, dass Variablen mit einem Startwert initiiert werden, Platzhalter besitzen dagegen anfangs keinen Wert. Wir werden Variablen nutzen, um die Parameter über die Iterationen zu speichern. Die Daten, mit denen wir das Modell trainieren werden, übergeben wir an Platzhalter.

Wie das Modell exakt definiert wird, zeigen die folgenden Teilkapitel. Wir stellen immer eine Kostenfunktion auf, die dann mit Hilfe eines Gradientenabstiegsverfahrens iterativ minimiert wird. Die Definition dieses sogenannten Trainingsschrittes sieht immer gleich aus:

```
1  train_step = tf.train
2                .GradientDescentOptimizer(learn_rate)
3                .minimize(cost)
```

Dabei ist *tf* die importierte TensorFlow-Bibliothek, *learn\_rate* ist die Geschwindigkeit bzw. Schrittweite des Verfahrens und *cost* ist die zuvor definierte Kostenfunktion.

Um wirklich Berechnungen durchführen zu können, muss in TensorFlow eine Session erzeugt werden. In dieser Session wird dann eine Schleife gestartet, die *train\_step* immer wieder mit der Datenbasis füttert, um den aktuellen Wert der Parameter zu verbessern. Je mehr Iterationen durchlaufen werden, desto exakter werden die Parameter approximiert.

#### 3.3.2. Einfache lineare Regression

Hier definieren wir unsere Platzhalter und Variablen wie folgt:

```
1  x = tf.placeholder(tf.float32, [None, 1])
2  y = tf.placeholder(tf.float32, [None, 1])
3  alpha = tf.Variable(tf.zeros([1]))
4  beta = tf.Variable(tf.zeros([1, 1]))
```

*x* und *y* sind die Platzhalter für die Datensätze für die später folgende Optimierung. *alpha* und *beta* sind die Variablen für die Parameter, welche mit Werten null initiiert werden. Daraus berechnen wir über den linearen funktionalen Zusammenhang die geschätzten *y* Werte:

```
1  y_calc = tf.matmul(x, beta) + alpha
```

Die Funktion *tf.matmul* führt Matrizenmultiplikation durch. *y\_calc* entspricht  $\pi_i$  aus dem Grundlagenkapitel 2.2. Nun definieren wir die Kostenfunktion als Mittelwert der Quadrate zwischen den wahren und geschätzten *y*-Werten:

```
1  cost = tf.reduce_mean(tf.square(y - y_calc))
```

Damit können wir die Session starten und unsere Parameter berechnen lassen. Mit einer Schrittweite von 0.0054 und 2000 Iterationen erhält man folgendes Ergebnis:

```

1  alpha:  -150.377670
2  beta:    2519.784424
3  cost:    15031988.000000

```

Die hier berechneten Werte für *alpha* und *beta* sind schon sehr nahe an den exakten Werten. Zusätzlich geben wir hier auch den aktuellen Wert der Kostenfunktion aus.

### 3.3.3. Multiple lineare Regression

In diesem Fall sind nun mehr unabhängige Variablen vorhanden, daher vergrößern wir die Dimension der Tensoren *x* und *beta* um eins.

```

1  x = tf.placeholder(tf.float32, [None, 2])
2  beta = tf.Variable(tf.zeros([2, 1]))

```

Die restlichen Variablen werden wie bei einfacher linearer Regression definiert. Das Gradientenverfahren ist bei zwei abhängigen Variablen komplexer, daher muss die Schrittweite verringert werden. Eine Folge davon ist, dass man mehr Schritte für dieselbe Präzision des Ergebnisses durchführen muss. Bei einer Schrittweite von 0.00071 und 50000 Schritten erhält man folgendes Ergebnis:

```

1  alpha:          -754.910095
2  beta_purchases: 2520.172363
3  beta_age:       17.224550
4  cost:           15004541.000000

```

### 3.3.4. Logistische Regression

Wir definieren unsere Tensoren exakt wie bei einfacher linearer Regression, da wie hier wieder mit je einer unabhängigen und einer abhängigen Variablen arbeiten. Die bisher verwendete Berechnung der *y*-Werte fügen wir nun zusätzlich in die logistische Funktion ein:

```

1  y_calc = 1 / (1 + tf.exp(- tf.matmul(x, beta) - alpha))

```

Bei der verwendeten Exponentialfunktion besteht die Gefahr, dass der Wert so nahe an null gerät, dass Python auf exakt null rundet und damit weiterrechnet. Um dieses Problem auszuschließen, werden die Werte der unabhängigen Variable zuvor linear in das Intervall  $[0, 1]$  transformiert. Die am Ende berechneten Parameter werden entsprechend linear zurücktransformiert, um den ursprünglichen Daten zu entsprechen.

Hier wollen wir die Likelihoodfunktion maximieren. TensorFlow bietet allerdings nur eine API für Minimierung, weshalb wir das Inverse der Likelihoodfunktion als Kostenfunktion verwenden. Zusätzlich wenden wir wieder den natürlichen Logarithmus auf die Funktion an. Die Kostenfunktion sieht also wie folgt aus:

```
1 cost = - tf.reduce_sum(  
2     tf.log(  
3         y * y_calc +  
4         (1 - y) * (1 - y_calc)  
5     )  
6 )
```

Wir wählen eine Schrittweite von 0.0001 und iterieren 1000 Schritte, um das folgende Ergebnis zu erhalten:

```
1 alpha:  -1.991089  
2 beta:   0.000080  
3 cost:   502.972290
```

## 3.4. SQL

Die "Structured Query Language" alias "SQL" ist eine Sprache zur Definition und Verarbeitung von Datenstrukturen in Datenbanksystemen und wird in nahezu allen Implementierungen relationaler Datenbanksysteme unterstützt. Oft liegen die Daten, welche man für Regressionsanalyse verwenden möchte in einem solchen Datenbanksystem.

SQL als Programmiersprache ist turingvollständig. Es ist also mit standardisierten SQL-Methoden möglich, Regression direkt in der Datenbank zu betreiben. Die konkrete Umsetzung hängt von der Art der Regression und dem spezifischen Datenbanksystem ab. In diesem Kapitel soll das nun in zwei Open-Source-Datenbanksystemen demonstriert werden, nämlich MySQL und PostgreSQL. Der vollständige SQL-Code befindet sich wegen der Länge wieder komplett im Anhang. Die MySQL-Skripte sind unter D zu finden, die Skripte für PostgreSQL liegen in E.

Im Gegensatz zu den beiden bisher vorgestellten Sprachen verfolgen wir hier kein einheitliches Grundprinzip, in dem sich alle Implementierungen ähnlich sind. Die einzige Gemeinsamkeit ist, dass wir in allen SQL-Skripten Prozeduren bzw. Funktionen definieren, welche bei Aufruf die Regressionsanalyse durchführen. Wir nehmen dazu an, dass die Daten für die Regression in einer Relation *sample* liegen.

Bei einfacher linearer Regression berechnen wir die Parameter exakt über die Formeln aus Kapitel 2.1.1. Bei multipler Regression verwenden wir die Matrixformel aus Kapitel 2.1.2. Hier müssen wir zusätzlich Algorithmen zur Transponierung, Multiplikation und Invertierung von Matrizen implementieren. Für logistische Regression steht uns keine explizite Formel zur Verfügung, weshalb wir ein Gradientenverfahren zur Maximierung der Likelihoodfunktion implementieren.

### 3.4.1. Einfache lineare Regression

Einfache lineare Regression kann ohne größeren Aufwand mit SQL umgesetzt werden. Wir berechnen zuerst die Mittelwerte über die Spalten *purchases* und *money*, dann die Summen

in Zähler und Nenner der Formel für  $\beta$ . Daraus können wir mit einfacher Arithmetik die beiden Parameter bestimmen.

Diese Berechnung kann man leicht in einer einzelnen Abfrage umsetzen. Das Skript für PostgreSQL tut das auch und definiert die genannten Berechnungsschritte als einzelne Views. In MySQL existiert die VIEW-Syntax nicht. Deshalb wird die Berechnung der Übersicht halber auf mehrere Abfragen aufgeteilt.

Führen wir die Prozeduren im jeweiligen Datenbanksystem aus, erhalten wir folgende Ergebnisse:

Tabelle 3.2.: Einfache lineare Regression in MySQL

variable	value
alpha	-150.478938563892081700000000000000
beta	2519.790023425114700000000000000000

Tabelle 3.3.: Einfache lineare Regression in PostgreSQL

variable	value
alpha	-150.478938563811511468617403503606400000000000000
beta	2519.7900234251071069143923667424

### 3.4.2. Multiple lineare Regression

Wir wollen zur Lösung dieses Regressionsproblems die Matrixformel aus Kapitel 2.1.2 anwenden. Dazu müssen Methoden für das Transponieren, Multiplizieren und Invertieren von Matrizen implementiert werden, da weder MySQL noch PostgreSQL über solche Funktionen verfügen. Wir müssen außerdem einen Weg finden, Matrizen im jeweiligen Datenbanksystem zu repräsentieren.

Die Matrizenrepräsentierung lösen wir unterschiedlich in den beiden Datenbanksystemen. In MySQL definieren wir uns temporäre Relationen, welche jeweils eine Matrix repräsentieren. Jede solche Relation besitzt dasselbe Schema und besteht aus den Attributen *row*, *column* und *value*. Die Ersten beiden enthalten die Indizes des Matrixelements, Letzteres enthält den Wert des jeweiligen Matrixelements.

Wir definieren insgesamt sieben solcher Relationen:

Tabelle 3.4.: Relationen für multiple lineare Regression in MySQL

Schema	entspricht folgender Matrix (vergleiche Berechnungsformel)
<i>matrix_X</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$X$
<i>matrix_y</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$y$
<i>matrix_transposed</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$X^T$
<i>matrix_product_1</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$X^T X$
<i>matrix_inverse</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$(X^T X)^{-1}$
<i>matrix_product_2</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$X^T y$
<i>matrix_results</i> ([ <i>row</i> , <i>column</i> , <i>value</i> ])	$(X^T X)^{-1} X^T y$

Die beiden erstgenannten Relationen werden mit den vorhandenen Werten aus der Relation *sample* befüllt. Die Relation *matrix\_transposed* wird mit einer einfachen Abfrage über der Relation *matrix\_X* berechnet, welche die Indizes für Zeile und Spalte vertauscht.

Die Relationen *matrix\_product\_1*, *matrix\_product\_2* und *matrix\_result* sind Ergebnisse von Matrizenmultiplikationen. Diese Produkte werden mit zwei Schleifen berechnet, die über die Zeilen und Spalten der Ergebnismatrix iterieren, also über die Tupel der Ergebnisrelation. Der jeweilige Wert zusammen mit dem Zeilen- und Spaltenindex der zu berechnenden Matrix wird über eine Abfrage eingefügt. Diese Abfrage nimmt das Kreuzprodukt beider zu multiplizierenden Relationen und betrachtet die entsprechende Zeile der ersten Relation und die entsprechende Spalte der zweiten Relation. Außerdem sollen das Attribut *row* der ersten Relation gleich dem Attribut *colum* der zweiten Relation sein. Der Wert für die zu berechnende Matrix ergibt sich als Summe über die Produkte der Werte der verbleibenden Tupel.

Für die Berechnung der inversen Matrix, also für die Relation *matrix\_inverse*, wird ein einfacher iterativer Algorithmus verwendet, welcher über alle Zeilen der zu invertierenden Matrix läuft. In jedem Schritt werden alle Elemente der Matrix so angepasst, sodass man am Ende das Inverse erhält. Details zu dem verwendeten Algorithmus findet man in [1].

Die komplette Berechnung wurde in eine Prozedur verpackt. Führt man diese aus, erhält man das folgende Ergebnis:

Tabelle 3.5.: Multiple lineare Regression in MySQL

variable	value
alpha	-766.736173784421688472047506524797
beta_purchases	2520.301741791306785710069856997055
beta_age	17.503722143360901662105040790569

Betrachten wir nun die Implementierung in PostgreSQL. Gegenüber MySQL hat man hier den Vorteil, dass mehrdimensionale Felder als Datentyp existieren. Wir verwenden also keine temporären Relationen für die Matrizen mehr, sondern speichern diese als zweidimensionales Array.

Neben der eigentlichen Prozedur zum Ausführen der Regression definieren wir drei weitere Funktionen:

Tabelle 3.6.: Funktionen für multiple lineare Regression in PostgreSQL

Name	Eingabe	Ausgabe	Beschreibung
<i>matrix_transpose</i>	<i>a numeric</i> (65, 30)[ ][ ]	<i>numeric</i> (65, 30)[ ][ ]	transponiert die Matrix <i>a</i>
<i>matrix_multiplication</i>	<i>a numeric</i> , (65, 30)[ ][ ], <i>b numeric</i> (65, 30)[ ][ ]	<i>numeric</i> (65, 30)[ ][ ]	multipliziert die Matrizen <i>a</i> und <i>b</i>
<i>matrix_inversion</i>	<i>a numeric</i> (65, 30)[ ][ ]	<i>numeric</i> (65, 30)[ ][ ]	invertiert die Matrix <i>a</i>

Mit diesen drei Funktionen lässt sich die Formel aus Kapitel 2.1.2 einfach umsetzen. Mit zwei auf der *sample*-Relation erzeugten Matrizen  $x$  und  $y$  nutzen wir die genannten Funktionen, um die Matrix mit den gesuchten Parametern zu berechnen.

Führt man die multiple lineare Regressionsanalyse in PostgreSQL durch, erhält man das folgende Ergebnis:

Tabelle 3.7.: Multiple lineare Regression in PostgreSQL

variable	value
alpha	-766.736173784421688472047506524797
beta_purchases	2520.301741791306785710069856997055
beta_age	17.503722143360901662105040790569

### 3.4.3. Logistische Regression

Zur Lösung des logistischen Regressionsproblems möchten wir in SQL ein Gradientenverfahren implementieren. Wir verwenden denselben Algorithmus in MySQL und PostgreSQL. Die Skripte unterscheiden sich lediglich in der datenbankspezifischen Syntax.

Wir verwenden für das Verfahren mehrere Relationen, in denen gewisse Informationen gespeichert und verarbeitet werden. Dabei weisen wir jedem Datenpunkt eine *id* zu, um verschiedene Relationen sinnvoll verknüpfen zu können.

Tabelle 3.8.: Tabellen für logistische Regression

Schema	Beschreibung
<i>datapoints</i> ([ <i>id</i> , <i>variable</i> , <i>value</i> ])	Ein Tupel enthält den Wert einer bestimmten unabhängigen Variable für einen bestimmten Datenpunkt.
<i>binary_values</i> ([ <i>id</i> , <i>value</i> ])	Ein Tupel enthält den Wert für die abhängige Variable eines bestimmten Datenpunktes.
<i>parameters</i> ([ <i>variable</i> , <i>old</i> , <i>new</i> ])	Ein Tupel enthält den alten und den neuen Wert des Parameters einer bestimmten Variable.
<i>logits</i> ([ <i>id</i> , <i>old</i> , <i>new</i> ])	Ein Tupel enthält den Wert der logistischen Funktion (also $\pi_i$ ) berechnet mit den alten oder neuen Parameterwerten für einen bestimmten Datenpunkt.
<i>gradient</i> ([ <i>variable</i> , <i>value</i> ])	Ein Tupel enthält den Wert der partiellen Ableitung nach einer bestimmten Variable.

Wir definieren außerdem einige Hilfsfunktionen. Diese Funktionen besitzen keinen Rückgabewert. Stattdessen bearbeiten sie die oben genannten Tabellen.

Tabelle 3.9.: Funktionen für logistische Regression

Name	Eingabe	Beschreibung
<i>calculate_logits</i>		berechnet den Wert der logistischen Funktion bzw. $\pi_i$ mit den alten und neuen Parameterwerten für jeden Datenpunkt
<i>calculate_gradient</i>		berechnet den Wert der partiellen Ableitungen nach allen Variablen
<i>calculate_new_parameters</i>	<i>step</i> <i>numeric</i> (65, 30)	berechnet die neuen Parameterwerte aus dem Gradienten und der übergebenen Schrittweite

Wir führen wie schon bei TensorFlow zuerst eine Lineartransformation für die Werte aus *money* durch und bilden diese linear auf das Intervall  $[0, 1]$  ab. Das dient erneut dazu, dass die Werte der Exponentialfunktion nicht zu nahe an null geraten. Die transformierten Werte fügen wir in die Relation *datapoints* ein. Danach werden die anderen Relationen erzeugt und initiale Werte eingefügt.

Es folgt eine *while*-Schleife, die solange läuft, bis entweder die vorgegebene Anzahl an Schritten erreicht wurde, oder die Schrittweite zu klein für die gewählte Präzision der Komazahlen wird. Wir berechnen zuerst den Gradienten, dann die neuen Parameter. Dann ist ein Aufruf von *calculate\_logits* nötig, um die logistische Funktion für die neuen Parameter zu berechnen.

Wir überprüfen, ob die neuen Parameter wirklich besser sind als die alten. Falls nicht, wird die Schrittweite halbiert, danach werden die Parameter und die Werte der logistischen Funktion erneut berechnet. Das wird solange wiederholt, bis die neuen Parameter besser sind oder die Schrittweite unter die Präzisionsgrenze fällt.

Haben wir die neuen Parameter erfolgreich berechnet, werden die Werte der *old*-Attribute in den Relationen *parameters* und *logits* mit den neuen Werten überschrieben und der Iterationsschritt ist beendet. Nachdem die Schleife beendet wurde, werden die Parameter wieder entsprechend linear transformiert, um den tatsächlichen Werten für *money* zu entsprechen.

Wir entscheiden uns für 1000 Iterationen und wählen eine Schrittweite von 0.008. Führt man die Prozeduren im jeweiligen Datenbanksystem aus, erhält man das folgende Ergebnis:

Tabelle 3.10.: Logistische Regression in MySQL

variable	value
alpha	-1.991090115311453143480846099933
beta_money	0.000080298536993602280846099933



Tabelle 3.11.: Logistische Regression in PostgreSQL

<b>variable</b>	<b>value</b>
beta_money	0.000080298565326972737807102042
alpha	-1.991090799730035106155294219916



## 4. Vergleich der verschiedenen Implementierungen

Nun wollen wir die implementierten Algorithmen im Bezug auf ihre Laufzeit miteinander vergleichen. Die Berechnung der Laufzeiten erfolgt erneut mit einem Python-Skript. Dieses Skript führt die jeweilige Regression über einen Kommandozeilenbefehl aus und misst die Zeit dieser Operation. Die Anzahl der Wiederholungen und die Anzahl der verwendeten Datenpunkte können als Parameter übergeben werden. Das Skript findet man im Anhang unter F.1.

Die Laufzeiten der Berechnungen werden als csv-Datei gespeichert. Pro Berechnung wird eine Zeile in die Ergebnisdatei geschrieben. Diese Zeile enthält die Programmiersprache, die Art der Regression, die Anzahl der verwendeten Datenpunkte und die Dauer der Berechnung.

Für jede Kombination der ersten drei Werte wurden 100 Berechnungen durchgeführt, falls die Dauer einer Berechnung höchstens 100 Sekunden betrug. Für Berechnungen, die länger als 100 Sekunden dauerten, wurde die Anzahl auf 50 Berechnungen reduziert. Skripte mit einer Laufzeit von länger als 1000 Sekunden wurden nur zehn mal ausgeführt. Für das Benchmarking wurde ein MacBook Pro (Mitte 2012, Betriebssystem macOS High Sierra 10.13.2) mit einem 2,9 GHz Intel Core i7 Prozessor und 8 GB 1600 MHz DDR3 Arbeitsspeicher verwendet.

Auch für das Auswerten der berechneten Benchmarks wurde ein Python-Skript erstellt, welches im Anhang F.2 abgedruckt ist. Dieses Skript liest die csv-Datei mit Benchmarks ein und gibt eine Tabelle mit den durchschnittlichen Laufzeiten pro Art der Regression, Programmiersprache und Anzahl verwendeter Datenpunkte aus. Außerdem wird pro Typ der Regression ein Plot zum Vergleich der Laufzeiten erzeugt.

### 4.1. Einfache lineare Regression

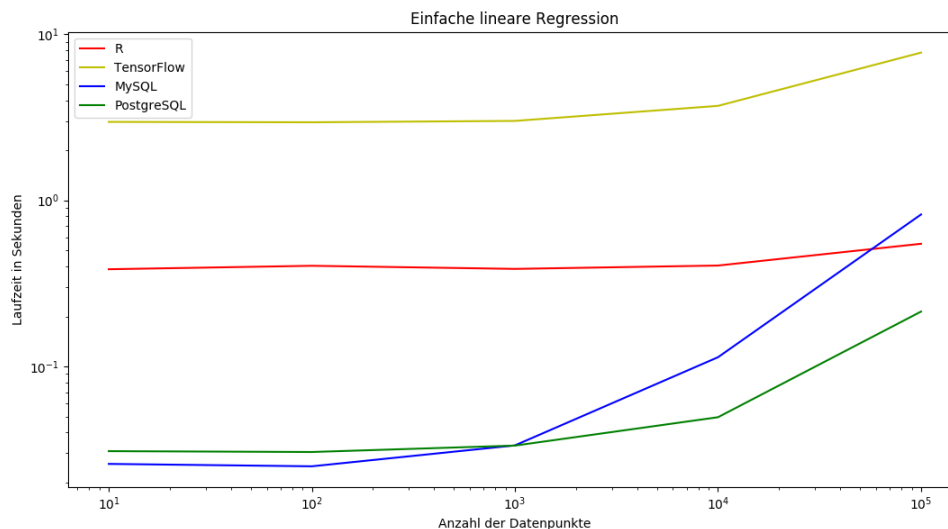
Folgende Ergebnisse erhalten wir bei einfacher linearer Regression:

Tabelle 4.1.: Laufzeiten in Sekunden für einfache lineare Regression

	<b>10</b>	<b>100</b>	<b>1000</b>	<b>10000</b>	<b>100000</b>
<b>r</b>	0.38453477	0.40336191	0.38638819	0.40475887	0.54730985
<b>tensorflow</b>	2.96190362	2.94665535	3.00290149	3.69807061	7.73616447
<b>mysql</b>	0.02591379	0.02506811	0.03348756	0.11362195	0.82149320
<b>postgresql</b>	0.03092405	0.03056195	0.03340294	0.04948901	0.21394655

Der zugehörige Graph sieht folgendermaßen aus:

#### 4. Vergleich der verschiedenen Implementierungen



TensorFlow und R haben eine relativ konstante Laufzeit, auch bei größeren Datenmengen. Dabei ist TensorFlow mit iterativer Berechnung wie erwartet mit Abstand am langsamsten. Die SQL-Implementierungen sind bei geringer Anzahl an Datenpunkten sogar die schnellsten Skripte. Für größer werdende Datenmengen erkennt man aber einen rapiden Anstieg in der Laufzeit.

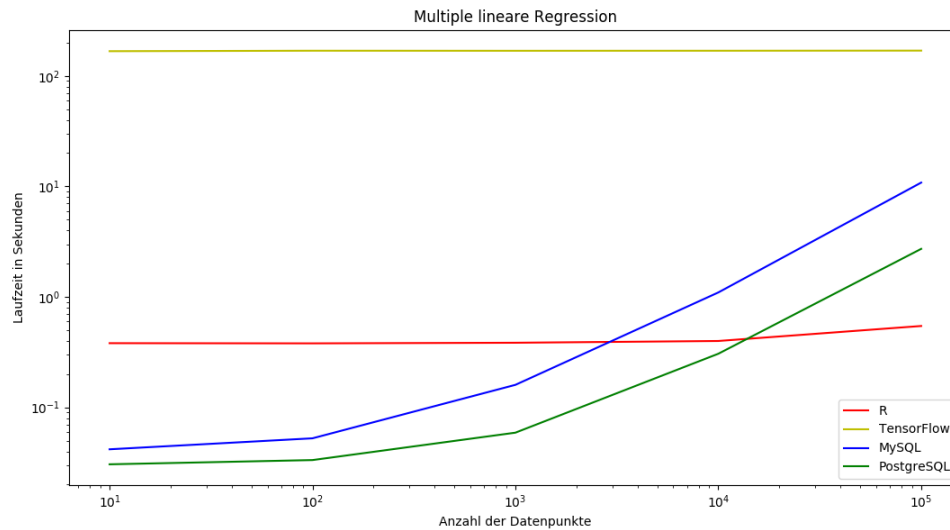
#### 4.2. Multiple lineare Regression

Die Skripte für multiple lineare Regression besitzen folgende Laufzeiten:

Tabelle 4.2.: Laufzeiten in Sekunden für multiple lineare Regression

	10	100	1000	10000	100000
<b>r</b>	0.38108478	0.37967851	0.38452695	0.39888490	0.54517789
<b>tensorflow</b>	167.220486	168.885465	168.668689	168.821214	169.254182
<b>mysql</b>	0.04174929	0.05248516	0.15990342	1.09634777	10.8214030
<b>postgresql</b>	0.03053954	0.03336087	0.05904620	0.30572490	2.71992954

Das geplottete Ergebnis sieht wie folgt aus:



Hier zeigt sich ein ähnliches Bild wie schon bei einfacher linearer Regression. Wieder liefern R und TensorFlow relativ konstante Laufzeiten, wobei die Laufzeit des TensorFlow-Skriptes wegen der 50000 durchgeführten Iterationen dieses Mal extrem langsam ist. Wieder sind die SQL-Skripte bei kleinen Datenmengen am schnellsten. Bei größeren Datenmengen werden sie allerdings von R geschlagen. Interessant ist außerdem, dass die PostgreSQL-Implementierung noch schneller als die Variante in MySQL. Die Arrays in PostgreSQL arbeiten also effizienter als die temporären Relationen in MySQL.

### 4.3. Logistische Regression

Betrachten wir zuletzt noch die Laufzeiten für logistische Regression:

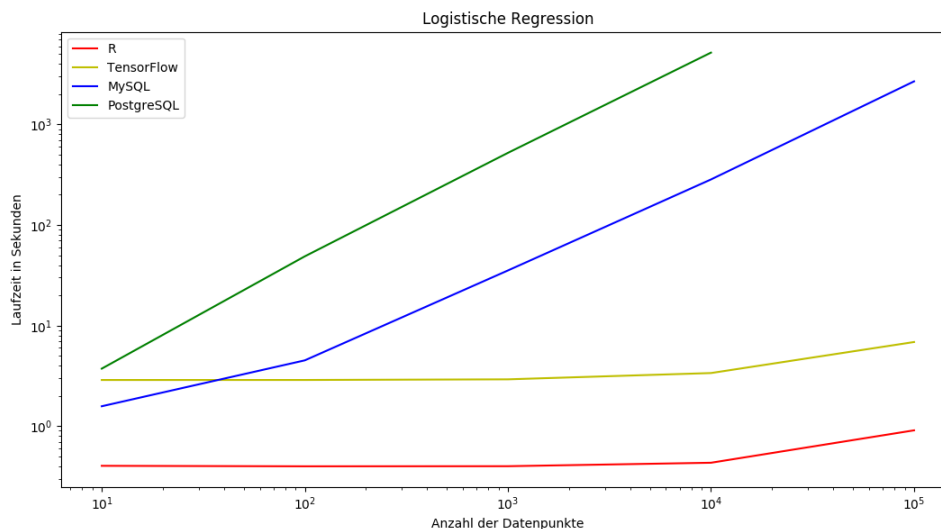
Tabelle 4.3.: Laufzeiten in Sekunden für logistische Regression

	10	100	1000	10000	100000
<b>r</b>	0.40403676	0.39953323	0.40047518	0.43356463	0.91014730
<b>tensorflow</b>	2.88191271	2.88278436	2.92470042	3.37794654	6.87024382
<b>mysql</b>	1.57908838	4.51631200	35.3472805	283.849979	2680.43056
<b>postgresql</b>	3.73521209	48.8452754	521.625671	5175.87997	

Die Visualisierung der Tabelle sieht so aus:

#### 4. Vergleich der verschiedenen Implementierungen

---



Wieder erkennt man eine Ähnlichkeit zu den vorherigen Diagrammen. Die SQL-Implementierungen sind nun aber von Anfang an deutlich langsamer als die Skripte in R und TensorFlow. Die Laufzeit steigt außerdem sehr schnell weiter an. So wurden für 100000 Datenpunkte in PostgreSQL gar keine Benchmarks mehr berechnet, da die erwartete Laufzeit etwa 50000 Sekunden, also knapp 14 Stunden beträgt. Klarer Gewinner ist hier R, wo auch 100000 Datenpunkte in weniger als einer Sekunde verarbeitet werden können.

## 5. Erweiterungspotenzial in Datenbanksystemen

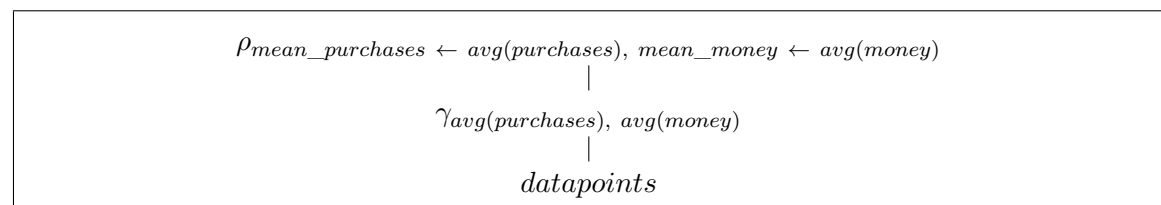
In Kapitel 3.4 wurde gezeigt, wie Regressionsanalyse in SQL durchgeführt werden kann. Die dazu implementierten Funktionen können als Erweiterungspotenzial für relationale Datenbanksysteme gesehen werden.

In diesen Funktionen sind die zu verwendende Relation und deren Attribute aktuell noch fest implementiert. Um diese Funktionen in der Praxis nutzbar zu machen, müsste man weitere Parameter einfügen, mit denen man Relation und Attribute zum Zeitpunkt der Ausführung bestimmen kann.

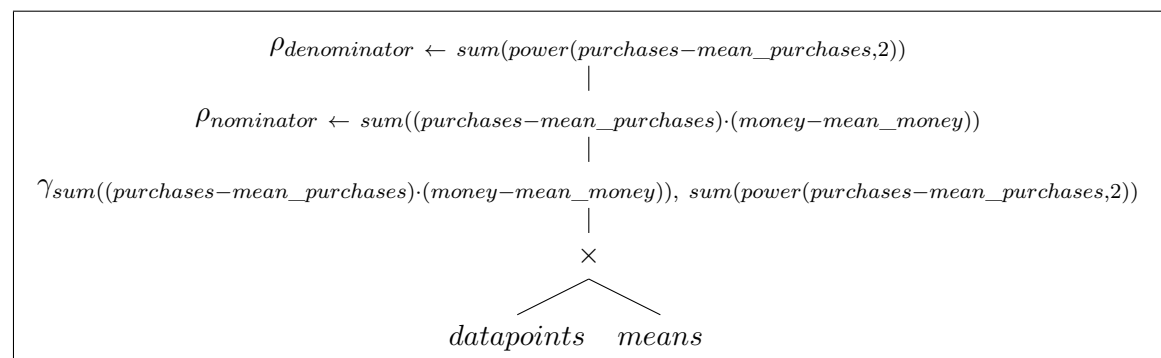
Wir wollen nun die Abfragen der in Kapitel 3.4 implementierten Funktionen mit Hilfe von Operatorbäumen darstellen und beschreiben.

### 5.1. Einfache lineare Regression

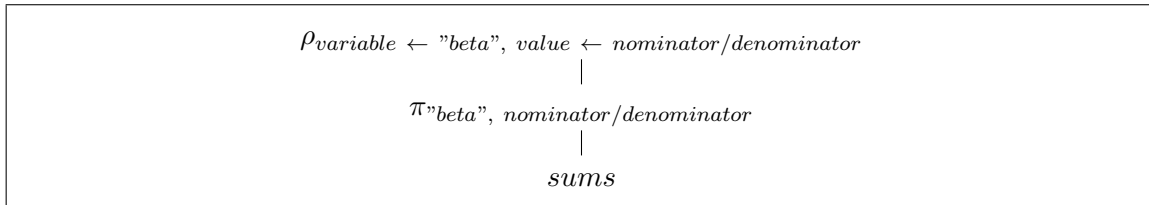
Die einfache lineare Regression besteht aus mehreren kleinen Teilabfragen. Zuerst berechnet man die Mittelwerte der beiden Attribute für die Regression. Das Ergebnis ist die Relation *means*:



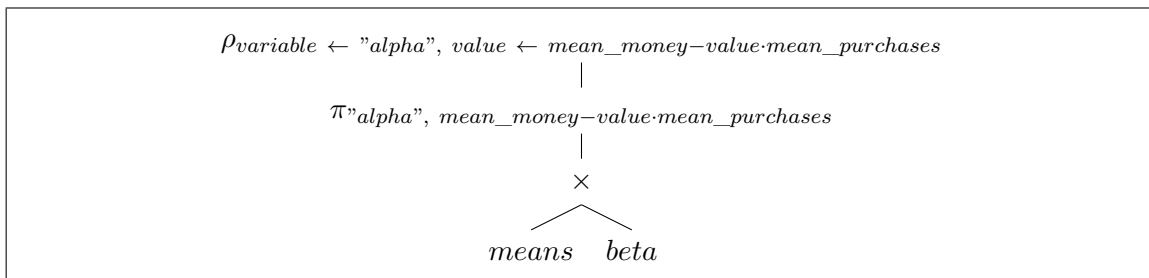
Mit diesen Mittelwerten berechnet man die Summen im Nenner und Zähler der Lösungsformel für  $\beta$  aus 2.1.1. Die Relation *sums* ist das Ergebnis dieser Abfrage:



Damit berechnet man einen Wert für *beta*:



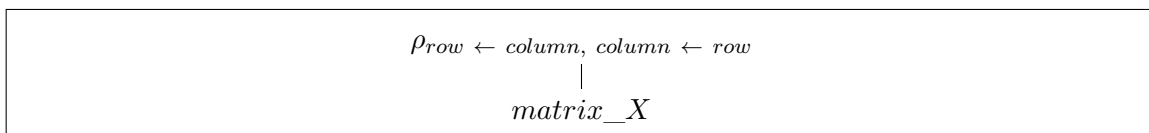
Mit Hilfe der Relation *beta* kann man nun auch *alpha* berechnen:



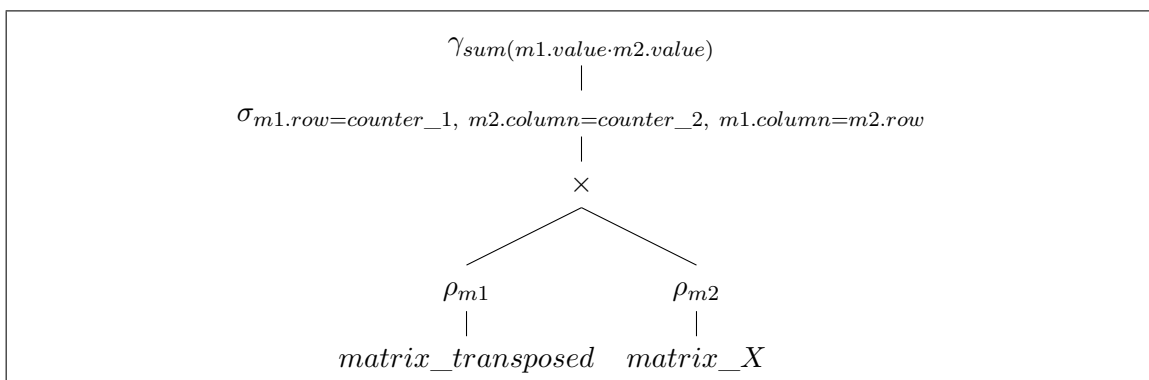
## 5.2. Multiple lineare Regression

Die hier gezeichneten Bäume entsprechen der Implementierung in MySQL. Dort verwenden wir Relationen und verarbeiten diese mit den Abfragen der hier dargestellten Bäume. In PostgreSQL verwenden wir dagegen Arrays und Schleifen zur Berechnung.

Wir beginnen mit der Berechnung der transponierten Matrix von *X*. Die zugehörige Relation wird *matrix\_transposed* genannt:



Die Matrixprodukte werden auch in MySQL mit Schleifen berechnet. Dabei wird ein Element der zu berechnenden Matrix mit folgender Abfrage bestimmt. Die Iteratoren *counter\_1* und *counter\_2* sind durch die Schleifen gegeben.



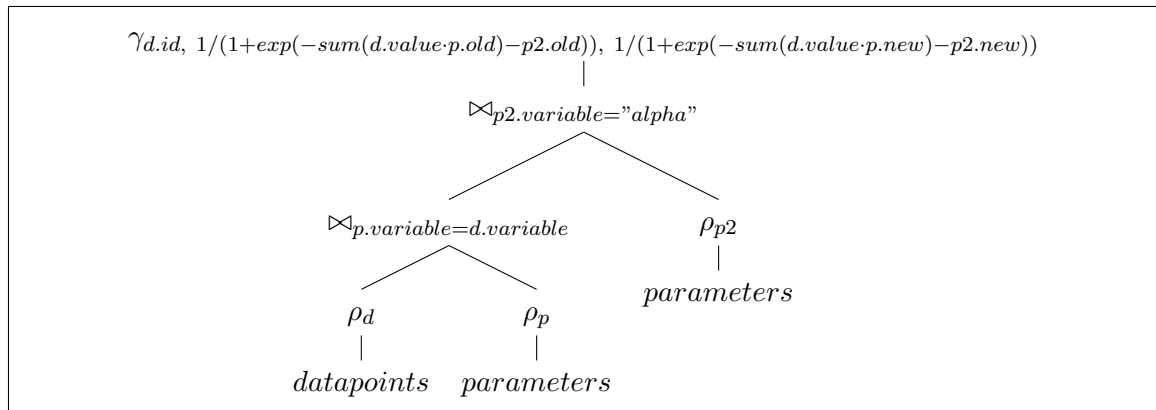


Die beiden Matrixprodukte zur Berechnung von *matrix\_product\_2* und *matrix\_result* besitzen denselben Operatorbaum, nur dass die beiden Relationen *matrix\_transposed* und *matrix\_y* bzw. *matrix\_inverse* und *matrix\_product\_2* verwendet werden.

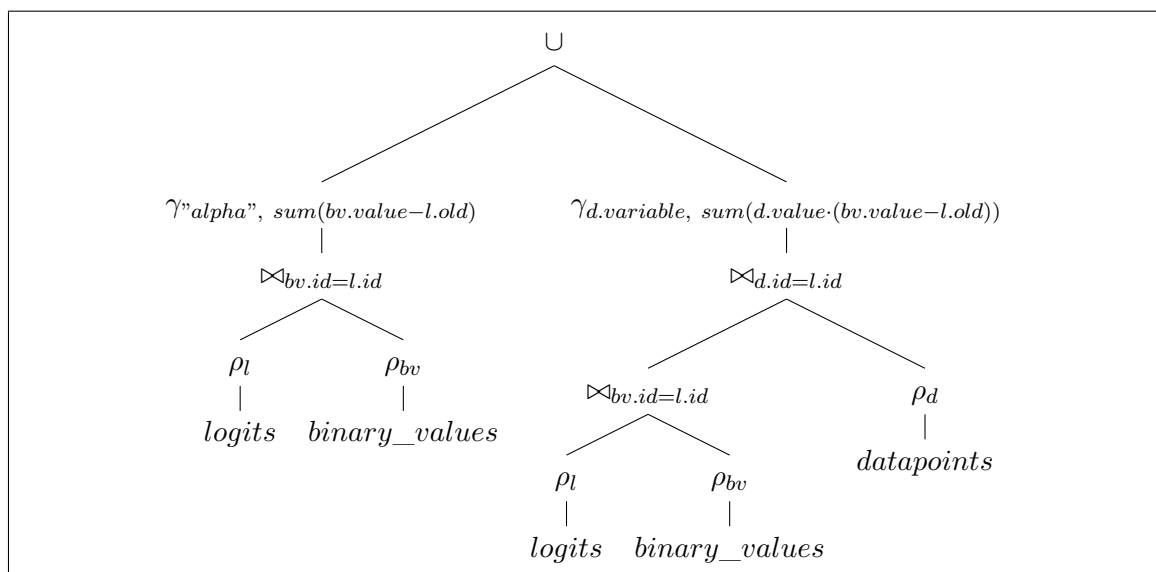
Zur Berechnung der inversen Matrix für die Relation *matrix\_inverse* wird in PostgreSQL und MySQL ein iterativer Algorithmus verwendet, dessen Äquivalent wir hier auf Grund der Komplexität nicht als Operatorbaum darstellen wollen.

### 5.3. Logistische Regression

Die Prozedur für die logistische Regression ist in verschiedene Teilprozeduren aufgeteilt, die jeweils eine Abfrage durchführen. In *calculate\_logit* wird die Relation *logits* mit folgender Abfrage befüllt:

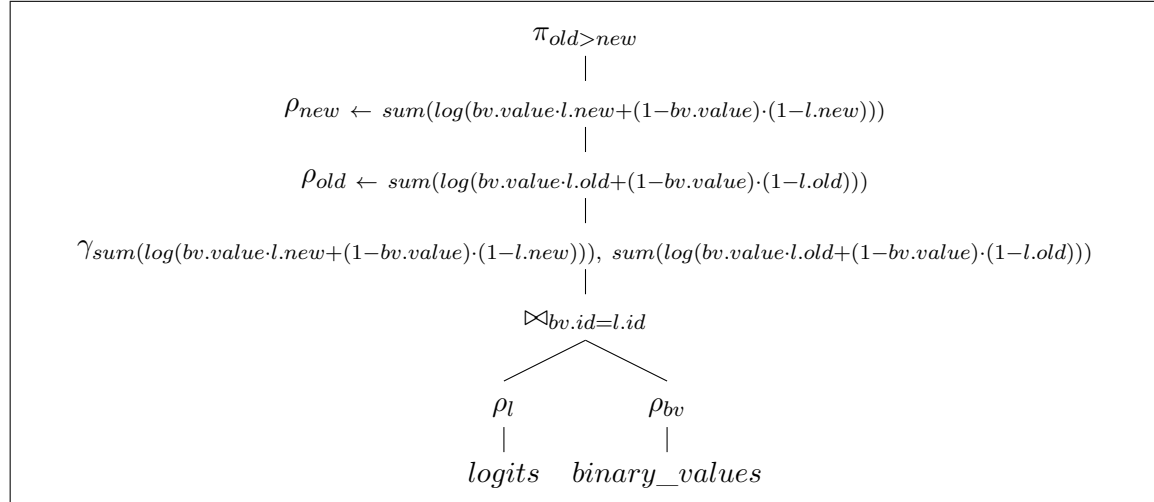


Die Prozedur *calculate\_gradient*, die den Gradienten in die Relation *gradient* schreibt, führt diese Abfrage aus:



Die Prozedur *calculate\_new\_parameters* führt eine Update-Abfrage durch. Hierfür zeichnen wir keinen Operatorbaum.

In der Hauptprozedur *logistic\_regression* wird in jeder Iteration überprüft, ob die mit der aktuellen Schrittweite berechneten neuen Parameter ein besseres Ergebnis liefern als die alten Parameter. Das geschieht mit folgender Abfrage:



Ansonsten werden in *logistic\_regression* nur die anderen Prozeduren aufgerufen und einfache Update-Abfragen gestellt.

## 6. Fazit

Wir haben in dieser Arbeit mit der Motivation für Regressionsanalyse begonnen. Danach haben wir die mathematischen Grundlagen der Regressionsanalyse erläutert. Darauf aufbauend haben wir die Umsetzung in R, TensorFlow und SQL demonstriert.

Hierbei haben zwei verschiedene Arten zur Berechnung der Parameter verwendet. Bei expliziter Berechnung wurden verschiedene Berechnungsformeln ausgewertet. Bei iterativer Berechnung wurden Optimierungsverfahren verwendet, um die Werte der gesuchten Parameter Schritt für Schritt besser zu approximieren.

In R haben wir ausschließlich explizite Berechnungen durchgeführt. Dazu haben wir die vorhandenen Funktionen *lm* für lineare Regression und *glm* für logistische Regression verwendet. Das Kürzel *lm* steht für "linear model", *glm* steht für "generalized linear model".

In TensorFlow kamen dagegen ausschließlich iterative Berechnungen zum Einsatz. Wir haben ein von TensorFlow implementiertes Gradientenverfahren genutzt, um eine von uns definierte Kostenfunktion zu minimieren. Die Kostenfunktionen waren dabei die Summe der kleinsten Quadrate bei linearer Regression und die inverse Likelihoodfunktion bei logistischer Regression.

In SQL haben wir beide Berechnungsarten umgesetzt. Bei linearer Regression haben wir die expliziten Formeln aus den Kapiteln 2.1.1 und 2.1.2 verwendet. Bei logistischer Regression wurde wiederum ein Gradientenverfahren angewandt. Dieses Mal wurde das Verfahren eigens implementiert.

Danach haben wir die Implementierungen bezüglich der Laufzeit miteinander verglichen. Dabei haben wir festgestellt, dass die SQL-Skripte bei kleinen Datenmengen und bei Verwendung von expliziten Berechnungsformeln durchaus mit R mithalten können. Besonders bei logistischer Regression erkannte man aber, dass die Implementierung in R von der Laufzeit her deutlich überlegen ist. TensorFlow war durch die strikte Verwendung iterativer Berechnungsarten immer langsamer als R.

Eine Implementierung der hier gezeigten Funktionen für Regression direkt in Datenbanksystemen und die Verwendung von effizienteren Algorithmen könnte die Performanz noch steigern und würde Regressionsanalyse in SQL auch praktisch nutzbar machen. Diese Arbeit ist ein erster "Proof of Concept" zu diesem Thema.



# Anhang



## A. Python-Skript zum Generieren der Beispieldaten

```
1 import sys
2 import random
3 import math
4
5 # Erzeuge Funktion, die die Daten in eine csv-Datei schreibt.
6 def outputCsv(data):
7     # Erzeuge ein Array mit den zu schreibenden Zeilen und übergebe die
8     ↪ Spaltennamen.
9     output = ["%s,%s,%s,%s" % (
10         "age",
11         "purchases",
12         "money",
13         "premium"
14     )]
15
16     # Iteriere über alle Datenpunkte.
17     for datapoint in data:
18         # Hänge eine Zeile an das output-Array an.
19         output.append("%s,%s,%s,%a" % (
20             datapoint["age"],
21             datapoint["purchases"],
22             datapoint["money"],
23             datapoint["premium"]
24         ))
25
26     # Öffne eine csv-Datei.
27     f = open("sample.csv", "w")
28
29     # Schreibe das output-Array als mit Zeilenumbrüchen gejointen String in
30     ↪ die Datei.
31     f.write("\n".join(output))
32
33     return
34
35 # Erzeuge Funktion, die die Daten in eine sql-Datei schreibt.
36 def outputSql(data):
37     # Erzeuge ein Array mit den zu schreibenden Zeilen.
```

```

36  # Füge SQL-Abfragen ein, die eine eventuell bestehende Tabelle löscht
    ↪ und neu erstellt.
37  # Beginne mit der INSERT-Abfrage.
38  output = [
39      "DROP TABLE IF EXISTS sample;",
40      "",
41      "CREATE TABLE sample (",
42      "    age INTEGER,",
43      "    purchases INTEGER,",
44      "    money INTEGER,",
45      "    premium INTEGER",
46      ");",
47      "",
48      "INSERT INTO sample (%s,%s,%s,%s) VALUES" % (
49          "age",
50          "purchases",
51          "money",
52          "premium"
53      )
54  ]
55
56  # Iteriere über alle Datenpunkte.
57  for datapoint in data:
58      # Füge eine Zeile an die INSERT-Abfrage an.
59      output.append("(%s,%s,%s,%s)," % (
60          datapoint["age"],
61          datapoint["purchases"],
62          datapoint["money"],
63          datapoint["premium"]
64      ))
65
66  # Ersetze das letzte Komma durch ein Semikolon, um die INSERT-Abfrage zu
    ↪ beenden.
67  output[-1] = output[-1][: -1] + ";"
68
69  # Öffne eine sql-Datei.
70  f = open("sample.sql", "w")
71
72  # Schreibe das output-Array als mit Zeilenumbrüchen gejointen String in
    ↪ die Datei.
73  f.write("\n".join(output))
74
75  return
76
77  def main(argv):

```



---

```

78  # Beende die Ausführung, wenn die Anzahl der zu erzeugenden Datenpunkte
    ↪ nicht übergeben wurde.
79  if len(argv) < 2:
80      print("Please provide number of datapoints that shall be generated.")
81      return
82
83  # Erzeuge ein leeres Array für die Daten.
84  data = []
85
86  # Iteriere über die Anzahl der zu erzeugenden Datenpunkte.
87  for i in range(0, int(argv[1])):
88      # Bestimme pseudo-zufällige Werte für den Datenpunkt.
89      age = int(max(random.normalvariate(25, 10) + 10, 18))
90      purchases = int(max(random.normalvariate(10, 10), 1))
91      money = int(max(purchases * 25 + random.normalvariate(0,
    ↪ (math.log(purchases) + 1) * 12), 0.01) * 100)
92      if random.uniform(0, 1) > math.exp(0.2 * purchases - 2) / (1 +
    ↪ math.exp(0.2 * purchases - 2)):
93          premium = 0
94      else:
95          premium = 1
96
97      # Hänge den Datenpunkt an das Array an.
98      data.append({
99          "age": age,
100         "purchases": purchases,
101         "money": money,
102         "premium": premium
103     })
104
105     # Schreibe die generierten Daten in eine .csv und eine .sql Datei.
106     outputCsv(data)
107     outputSql(data)
108     return
109
110     # Führe die main-Funktion aus.
111     if __name__ == "__main__":
112         main(sys.argv)

```

---



## B. R-Skripte

### B.1. Einfache lineare Regression

```
1  # Lese die übergebenen Argumente ein.
2  args = commandArgs(trailingOnly = TRUE)
3
4  # Setzte default-Werte für die Anzahl der Datenpunkte und ob geplottet
   ↪ werden soll.
5  n <- 1000
6  plot <- TRUE
7
8  # Ändere die default-Werte, falls entsprechende Argumente übergeben
   ↪ wurden.
9  if (length(args) == 1) {
10     if (substr(args[1], 1, 1) == "-") {
11         plot <- FALSE
12     } else {
13         n = strtoi(args[1])
14     }
15 }
16 if (length(args) == 2) {
17     if (substr(args[1], 1, 1) == "-") {
18         n = strtoi(args[2])
19     } else {
20         n = strtoi(args[1])
21     }
22     plot <- FALSE
23 }
24
25 # Speichere die aktuelle Zeit zur Zeitmessung.
26 start_time <- Sys.time()
27
28 # Lies die Daten aus der csv-Datei ein.
29 data <- head(read.csv2("../data/sample.csv", sep = ",", header = TRUE), n)
30
31 # Definiere das Modell.
32 modell <- as.formula("money ~ purchases")
33
34 # Führe die Regression durch.
35 slr <- lm(modell, data = data)
```

```
36
37 # Speichere die aktuelle Zeit zur Zeitmessung.
38 end_time <- Sys.time()
39
40 # Drucke die Ergebnisse der Regressionsanalyse und die Laufzeit.
41 print(slr)
42 print(end_time - start_time)
43
44 # Erstelle einen Plot.
45 if (plot) {
46   # Bestimme die Grenzen für die unabhängige Variable.
47   xmin <- min(data$purchases)
48   xmax <- max(data$purchases)
49
50   # Bestimme die Grenzen für die abhängige Variable.
51   ymin <- min(data$money)
52   ymax <- max(data$money)
53
54   # Bestimme die Koeffizienten aus der Regressionsanalyse.
55   b0 <- coef(slr["coefficients"])[1]
56   b1 <- coef(slr["coefficients"])[2]
57
58   # Erzeuge Vektoren zum Plot der linearen Funktion.
59   xplot <- c(xmin - 1, xmax + 1)
60   yplot <- c(b0 + (xmin - 1) * b1, b0 + (xmax + 1) * b1)
61
62   # Erstelle den Plot.
63   plot(
64     c(xmin - 1, xmax + 1),
65     c(ymin - 1, ymax + 1),
66     type = "n",
67     xlab = "purchases",
68     ylab = "money",
69     main = "Einfache lineare Regression",
70     sub = paste("money = ", b0, "+", b1, "* purchases"),
71     col.sub = "darkgray"
72   )
73   # Füge die Datenpunkte ein.
74   lines(
75     data$purchases,
76     data$money,
77     type="p"
78   )
79   # Füge die Ausgleichsgerade ein.
80   lines(
81     xplot,
```

```

82     yplot,
83     col = "red",
84     lwd = 2
85   )
86 }

```

## B.2. Multiple lineare Regression

```

1  # Lese die übergebenen Argumente ein.
2  args = commandArgs(trailingOnly = TRUE)
3
4  # Setzte default-Wert für die Anzahl der Datenpunkte.
5  n = 1000
6
7  # Ändere den default-Wert, falls ein entsprechendes Argumente übergeben
   ↳ wurde.
8  if (length(args) > 0) {
9    n = strtoi(args[1])
10 }
11
12 # Speichere die aktuelle Zeit zur Zeitmessung.
13 start_time <- Sys.time()
14
15 # Lies die Daten aus der csv-Datei ein.
16 data <- head(read.csv2("./data/sample.csv", sep = ",", header = TRUE), n)
17
18 # Definiere das Modell.
19 modell <- as.formula("money ~ purchases + age")
20
21 # Führe die Regression durch.
22 mlr <- lm(modell, data = data)
23
24 # Speichere die aktuelle Zeit zur Zeitmessung.
25 end_time <- Sys.time()
26
27 # Drucke die Ergebnisse der Regressionsanalyse und die Laufzeit.
28 print(mlr)
29 print(end_time - start_time)

```

## B.3. Logistische Regression

```

1  # Lese die übergebenen Argumente ein.
2  args = commandArgs(trailingOnly = TRUE)
3

```

```
4  # Setzte default-Werte für die Anzahl der Datenpunkte und ob geplottet
   ↪ werden soll.
5  n <- 1000
6  plot <- TRUE
7
8  # Ändere die default-Werte, falls entsprechende Argumente übergeben
   ↪ wurden.
9  if (length(args) == 1) {
10     if (substr(args[1], 1, 1) == "-") {
11         plot <- FALSE
12     } else {
13         n = strttoi(args[1])
14     }
15 }
16 if (length(args) == 2) {
17     if (substr(args[1], 1, 1) == "-") {
18         n = strttoi(args[2])
19     } else {
20         n = strttoi(args[1])
21     }
22     plot <- FALSE
23 }
24
25 # Speichere die aktuelle Zeit zur Zeitmessung.
26 start_time <- Sys.time()
27
28 # Lies die Daten aus der csv-Datei ein.
29 data <- head(read.csv2("./data/sample.csv", sep = ",", header = TRUE), n)
30
31 # Definiere das Modell.
32 modell <- as.formula("premium ~ money")
33
34 # Führe die Regression durch.
35 logit <- glm(modell, family = binomial, data = data)
36
37 # Speichere die aktuelle Zeit zur Zeitmessung.
38 end_time <- Sys.time()
39
40 # Drucke die Ergebnisse der Regressionsanalyse und die Laufzeit.
41 print(logit)
42 print(end_time - start_time)
43
44 # Erstelle einen Plot.
45 if (plot) {
46     # Bestimme die Grenzen für die unabhängige Variable.
47     xmin <- min(data$money)
```

```
48  xmax <- max(data$money)
49
50  # Bestimme eine Funktion, die den Wert der logistischen Funktion
   ↪ berechnet.
51  logitFunction <- function(x){
52    # Verwende die Parameter aus der Regressionsanalyse.
53    b0 <- coef(logit["coefficients"])[1]
54    b1 <- coef(logit["coefficients"])[2]
55    c <- b0 + x * b1
56    return(exp(c) / (1 + exp(c)))
57  }
58
59  # Erzeuge Vektoren zum Plot der logistischen Funktion.
60  xplot <- seq(xmin - 1, xmax + 1, 1000)
61  yplot <- logitFunction(xplot)
62
63  # Erstelle den Plot.
64  plot(
65    c(xmin - 1, xmax + 1),
66    c(-0.2, 1.2),
67    type = "n",
68    xlab = "money",
69    ylab = "premium",
70    main = "Logistische Regression"
71  )
72  # Füge die Datenpunkte ein.
73  lines(
74    data$money,
75    data$premium,
76    type="p"
77  )
78  # Füge die logistische Funktion ein.
79  lines(
80    xplot,
81    yplot,
82    col = "red",
83    lwd = 2
84  )
85 }
```





## C. TensorFlow-Skripte

### C.1. Einfache lineare Regression

```
1  import numpy as np
2  import tensorflow as tf
3  import os.path as p
4  import csv
5  import sys
6  import matplotlib.pyplot as plt
7  from time import time
8
9  # Erstelle eine Funktion zum einlesen der Daten aus der csv-Daten.
10 def get_data(n_samples):
11     # Bestimme den Pfad der csv-Datei und öffne die Datei.
12     filename = p.abspath(p.join(p.dirname(p.realpath(__file__)), "..",
13     ↪ "data", "sample.csv"))
14     csvfile = open(filename, newline="")
15     csvreader = csv.reader(csvfile, delimiter=",", quotechar="|")
16
17     # Erstelle leere Arrays für die Daten.
18     x = []
19     x_plot = []
20     y = []
21
22     # Iteriere über die Zeilen der csv-Datei.
23     for row in csvreader:
24         # Überspringe die erste Zeile.
25         if not row[0] == "age":
26             # Füge die Daten der Zeile zum jeweiligen Array hinzu.
27             x.append([int(row[1])])
28             x_plot.append(int(row[1]))
29             y.append(int(row[2]))
30
31     # Gib die Arrays bis zu der gewünschten Menge an Datenpunkten zurück.
32     return (np.array(x[:n_samples]), x_plot[:n_samples],
33     ↪ np.transpose([y[:n_samples]]))
34
35 def main(argv):
36     # Bestimme default-Werte für die Anzahl der Datenpunkte und ob geplottet
37     ↪ werden soll.
```

```
35 datapoint_size = 1000
36 plot = True
37
38 # Ändere die default-Werte, wenn entsprechende Argumente übergeben
   ↪ wurden.
39 if len(argv) == 2:
40     if argv[1] == "-":
41         plot = False
42     else:
43         datapoint_size = int(argv[1])
44 elif len(argv) == 3:
45     plot = False
46     if argv[1] == "-":
47         datapoint_size = int(argv[2])
48     else:
49         datapoint_size = int(argv[1])
50
51 # Bestimme die aktuelle Zeit zur Zeitmessung.
52 start_time = time()
53
54 # Bestimme die Anzahl der Iterationen und die Schrittweite (abhängig von
   ↪ der Anzahl der Datenpunkte.)
55 steps = 2000
56 if datapoint_size <= 10:
57     learn_rate = 0.0076
58 elif datapoint_size <= 100:
59     learn_rate = 0.0064
60 elif datapoint_size <= 1000:
61     learn_rate = 0.0056
62 elif datapoint_size <= 10000:
63     learn_rate = 0.0054
64 elif datapoint_size <= 100000:
65     learn_rate = 0.0054
66
67 # Deklariere die Platzhalter und Variablen.
68 x = tf.placeholder(tf.float32, [None, 1])
69 y = tf.placeholder(tf.float32, [None, 1])
70 alpha = tf.Variable(tf.zeros([1]))
71 beta = tf.Variable(tf.zeros([1, 1]))
72 y_calc = tf.matmul(x, beta) + alpha
73
74 # Definiere die Kostenfunktion und die Minimierungsoperation.
75 cost = tf.reduce_mean(tf.square(y - y_calc))
76 train_step =
   ↪ tf.train.GradientDescentOptimizer(learn_rate).minimize(cost)
77
```

```

78  # Importiere die Daten.
79  (all_xs, plot_xs, all_ys) = get_data(datapoint_size)
80
81  # Starte eine Session in TensorFlow.
82  sess = tf.Session()
83  init = tf.global_variables_initializer()
84  sess.run(init)
85
86  # Iteriere und trainiere.
87  for i in range(steps):
88      feed = { x: all_xs, y: all_ys }
89      sess.run(train_step, feed_dict=feed)
90
91  # Bestimme die aktuellen Parameterwerte nach der Anzahl der
92  ↪ Iterationen.
93  (curr_alpha, curr_beta, curr_cost) = sess.run([alpha, beta, cost],
94  ↪ feed_dict=feed)
95
96  # Bestimme die aktuelle Zeit zur Zeitmessung.
97  end_time = time()
98
99  # Drucke die Ergebnisse.
100 print("alpha:  %f" % curr_alpha)
101 print("beta:   %f" % curr_beta)
102 print("cost:   %f" % curr_cost)
103 print("")
104 print("time elapsed:  %f sec" % (end_time - start_time))
105
106 # Erstelle einen Plot (falls gewünscht).
107 if plot:
108     plt.plot(plot_xs, all_ys, "ro", label="Original data")
109     plt.plot(plot_xs, curr_beta * all_xs + curr_alpha, label="Fitted
110     ↪ line")
111     plt.legend()
112     plt.show()
113
114 # Führe die main-Funktion aus.
115 if __name__ == "__main__":
116     main(sys.argv)

```

## C.2. Multiple lineare Regression

```

1  import numpy as np
2  import tensorflow as tf
3  import os.path as p

```

```
4 import csv
5 import sys
6 from time import time
7
8 # Erstelle eine Funktion zum einlesen der Daten aus der csv-Daten.
9 def get_data(n_samples):
10     # Bestimme den Pfad der csv-Datei und öffne die Datei.
11     filename = p.abspath(p.join(p.dirname(p.realpath(__file__)), "..",
12     ↪ "data", "sample.csv"))
13     csvfile = open(filename, newline="")
14     csvreader = csv.reader(csvfile, delimiter=",", quotechar="|")
15
16     # Erstelle leere Arrays für die Daten.
17     x = []
18     y = []
19
20     # Iteriere über die Zeilen der csv-Datei.
21     for row in csvreader:
22         # Überspringe die erste Zeile.
23         if not row[0] == "age":
24             # Füge die Daten der Zeile zum jeweiligen Array hinzu.
25             x.append([int(row[1]), int(row[0])])
26             y.append(int(row[2]))
27
28     # Gib die Arrays bis zu der gewünschten Menge an Datenpunkten zurück.
29     return (np.array(x[:n_samples]), np.transpose([y[:n_samples]]))
30
31 def main(argv):
32     # Bestimme default-Wert für die Anzahl der Datenpunkte.
33     datapoint_size = 1000
34
35     # Ändere den default-Wert, wenn ein entsprechendes Argument übergeben
36     ↪ wurde.
37     if len(argv) == 2:
38         datapoint_size = int(argv[1])
39
40     # Bestimme die aktuelle Zeit zur Zeitmessung.
41     start_time = time()
42
43     # Bestimme die Anzahl der Iterationen und die Schrittweite (abhängig von
44     ↪ der Anzahl der Datenpunkte.)
45     steps = 50000
46     if datapoint_size <= 10:
47         learn_rate = 0.00093
48     elif datapoint_size <= 100:
49         learn_rate = 0.00078
```

```

47     elif datapoint_size <= 1000:
48         learn_rate = 0.0007
49     elif datapoint_size <= 10000:
50         learn_rate = 0.00071
51     elif datapoint_size <= 100000:
52         learn_rate = 0.00071
53
54     # Deklariere die Platzhalter und Variablen.
55     x = tf.placeholder(tf.float32, [None, 2])
56     y = tf.placeholder(tf.float32, [None, 1])
57     alpha = tf.Variable(tf.zeros([1]))
58     beta = tf.Variable(tf.zeros([2, 1]))
59     y_calc = tf.matmul(x, beta) + alpha
60
61     # Definiere die Kostenfunktion und die Minimierungsoperation.
62     cost = tf.reduce_mean(tf.square(y - y_calc))
63     train_step =
64         ↪ tf.train.GradientDescentOptimizer(learn_rate).minimize(cost)
65
66     # Importiere die Daten.
67     (all_xs, all_ys) = get_data(datapoint_size)
68
69     # Starte eine Session in TensorFlow.
70     sess = tf.Session()
71     init = tf.global_variables_initializer()
72     sess.run(init)
73
74     # Iteriere und trainiere.
75     for i in range(steps):
76         feed = { x: all_xs, y: all_ys }
77         sess.run(train_step, feed_dict=feed)
78
79     # Bestimme die aktuellen Parameterwerte nach der Anzahl der
80     ↪ Iterationen.
81     (curr_alpha, curr_beta, curr_cost) = sess.run([alpha, beta, cost],
82         ↪ feed_dict=feed)
83
84     # Bestimme die aktuelle Zeit zur Zeitmessung.
85     end_time = time()
86
87     # Drucke die Ergebnisse.
88     print("alpha:           %f" % curr_alpha)
89     print("beta_purchases:  %f" % curr_beta[0])
90     print("beta_age:         %f" % curr_beta[1])
91     print("cost:             %f" % curr_cost)
92     print("")

```

```
90     print("time elapsed: %f sec" % (end_time - start_time))
91
92     # Führe die main-Funktion aus.
93     if __name__ == "__main__":
94         main(sys.argv)
```

### C.3. Logistische Regression

```
1  import numpy as np
2  import tensorflow as tf
3  import os.path as p
4  import csv
5  import sys
6  import matplotlib.pyplot as plt
7  from time import time
8
9  # Erstelle eine Funktion zum einlesen der Daten aus der csv-Daten.
10 def get_data(n_samples):
11     # Bestimme den Pfad der csv-Datei und öffne die Datei.
12     filename = p.abspath(p.join(p.dirname(p.realpath(__file__)), "..",
13     ↪ "data", "sample.csv"))
14     csvfile = open(filename, newline="")
15     csvreader = csv.reader(csvfile, delimiter=",", quotechar="|")
16
17     # Erstelle leere Arrays für die Daten.
18     x = []
19     x_plot = []
20     y = []
21
22     # Iteriere über die Zeilen der csv-Datei.
23     for row in csvreader:
24         # Überspringe die erste Zeile.
25         if not row[0] == "age":
26             # Füge die Daten der Zeile zum jeweiligen Array hinzu.
27             x.append([int(row[2])])
28             x_plot.append(int(row[2]))
29             y.append(int(row[3]))
30
31     # Gib die Arrays bis zu der gewünschten Menge an Datenpunkten zurück.
32     return (np.array(x[:n_samples]), x_plot[:n_samples],
33     ↪ np.transpose([y[:n_samples]]))
34
35 def main(argv):
36     # Bestimme default-Werte für die Anzahl der Datenpunkte und ob geplottet
37     ↪ werden soll.
```

```

35 datapoint_size = 1000
36 plot = True
37
38 # Ändere die default-Werte, wenn entsprechende Argumente übergeben
   ↪ wurden.
39 if len(argv) == 2:
40     if argv[1] == "-":
41         plot = False
42     else:
43         datapoint_size = int(argv[1])
44 elif len(argv) == 3:
45     plot = False
46     if argv[1] == "-":
47         datapoint_size = int(argv[2])
48     else:
49         datapoint_size = int(argv[1])
50
51 # Bestimme die aktuelle Zeit zur Zeitmessung.
52 start_time = time()
53
54 # Bestimme die Anzahl der Iterationen und die Schrittweite (abhängig von
   ↪ der Anzahl der Datenpunkte.)
55 steps = 1000
56 if datapoint_size == 10:
57     learn_rate = 1
58 elif datapoint_size == 100:
59     learn_rate = 0.1
60 elif datapoint_size == 1000:
61     learn_rate = 0.01
62 elif datapoint_size == 10000:
63     learn_rate = 0.001
64 elif datapoint_size == 100000:
65     learn_rate = 0.0001
66
67 # Deklariere die Platzhalter und Variablen.
68 x = tf.placeholder(tf.float32, [None, 1])
69 y = tf.placeholder(tf.float32, [None, 1])
70 alpha = tf.Variable(tf.zeros([1]))
71 beta = tf.Variable(tf.zeros([1, 1]))
72 y_calc = 1 / (1 + tf.exp(- tf.matmul(x, beta) - alpha))
73
74 # Definiere die Kostenfunktion und die Minimierungsoperation.
75 cost = - tf.reduce_sum(
76     tf.log(
77         y * y_calc +
78         (1 - y) * (1 - y_calc)

```

```
79     )
80 )
81 train_step =
82     ↪ tf.train.GradientDescentOptimizer(learn_rate).minimize(cost)
83
84 # Importiere die Daten.
85 (all_xs, plot_xs, all_ys) = get_data(datapoint_size)
86
87 # Transformiere die unabhängige Variable linear ins Interval zwischen 0
88 ↪ und 1.
89 min_x = min(all_xs)
90 max_x = max(all_xs)
91 all_xs = (all_xs - min_x) / (max_x - min_x)
92
93 # Starte eine Session in TensorFlow.
94 sess = tf.Session()
95 init = tf.global_variables_initializer()
96 sess.run(init)
97
98 # Iteriere und trainiere.
99 for i in range(steps):
100     feed = { x: all_xs, y: all_ys }
101     sess.run(train_step, feed_dict=feed)
102
103 # Bestimme die aktuellen Parameterwerte nach der Anzahl der
104 ↪ Iterationen.
105 (curr_alpha, curr_beta, curr_cost) = sess.run([alpha, beta, cost],
106     ↪ feed_dict=feed)
107
108 # Transformiere die Parameter linear, um den originalen Daten zu
109 ↪ entsprechen.
110 curr_beta = curr_beta / (max_x - min_x)
111 curr_alpha = curr_alpha - curr_beta * min_x
112
113 # Bestimme die aktuelle Zeit zur Zeitmessung.
114 end_time = time()
115
116 # Drucke die Ergebnisse.
117 print("alpha:  %f" % curr_alpha)
118 print("beta:   %f" % curr_beta)
119 print("cost:   %f" % curr_cost)
120 print("")
121 print("time elapsed: %f sec" % (end_time - start_time))
122
123 # Erstelle einen Plot (falls gewünscht).
124 if plot:
```



```
120     all_xs = all_xs * (max_x - min_x) + min_x
121     plot_ys = 1 / (1 + np.exp(- curr_beta * all_xs - curr_alpha))
122     plot_order = np.argsort(plot_xs)
123     plt.plot(plot_xs, all_ys, "ro", label="Original data")
124     plt.plot(np.array(plot_xs)[plot_order], np.array(plot_ys)[plot_order],
125              ↪ label="Fitted line")
126     plt.legend()
127     plt.show()
128
129 # Führe die main-Funktion aus.
130 if __name__ == "__main__":
131     main(sys.argv)
```



## D. MySQL-Skripte

### D.1. Einfache lineare Regression

```
1  -- Lösche die bestehende Prozedur, falls vorhanden.
2  DROP PROCEDURE IF EXISTS simple_linear_regression;
3
4  DELIMITER ;;
5
6  -- Erstelle die Prozedur für einfache lineare Regression.
7  CREATE PROCEDURE `simple_linear_regression`(IN number_datapoints INT(11))
8  BEGIN
9
10     -- Deklariere die verwendeten Variablen.
11     DECLARE purchases_mean DECIMAL(65, 30);
12     DECLARE money_mean DECIMAL(65, 30);
13     DECLARE alpha DECIMAL(65, 30);
14     DECLARE beta DECIMAL(65, 30);
15
16     -- Erstelle eine temporäre Relation für die zu verwendenden Datenpunkte.
17     DROP TEMPORARY TABLE IF EXISTS datapoints;
18     CREATE TEMPORARY TABLE datapoints (
19         purchases INT(11),
20         money INT(11)
21     );
22
23     -- Füge die gewünschte Anzahl der Datenpunkte in die temporäre Relation
24     ↪ ein.
25     INSERT INTO datapoints
26     SELECT purchases, money
27     FROM sample
28     LIMIT number_datapoints;
29
30     -- Berechne die Mittelwerte der abhängigen und unabhängigen Variable.
31     SET purchases_mean = (
32         SELECT AVG(purchases)
33         FROM datapoints
34     );
35     SET money_mean = (
36         SELECT AVG(money)
37         FROM datapoints
```

```
37 );
38
39 -- Berechne beta.
40 SET beta = (
41     SELECT SUM((purchases - purchases_mean) * (money - money_mean))
42     FROM datapoints
43 );
44 SET beta = beta / (
45     SELECT SUM(POWER(purchases - purchases_mean, 2))
46     FROM datapoints
47 );
48
49 -- Berechne alpha.
50 SET alpha = money_mean - (beta * purchases_mean);
51
52 -- Gib eine Relation mit Parameternamen und zugehörigem Wert zurück.
53 SELECT 'alpha' AS `variable`, alpha AS `value`
54 UNION
55 SELECT 'beta' AS `variable`, beta AS `value`;
56
57 -- Lösche die temporäre Relation mit den Datenpunkten wieder.
58 DROP TEMPORARY TABLE IF EXISTS datapoints;
59
60 END;;
61
62 DELIMITER ;
```

## D.2. Multiple lineare Regression

```
1 -- Lösche die bestehende Prozedur, falls vorhanden.
2 DROP PROCEDURE IF EXISTS multiple_linear_regression;
3
4 DELIMITER ;;
5
6 -- Erstelle die Prozedur für multiple lineare Regression.
7 CREATE PROCEDURE multiple_linear_regression(IN number_datapoints INT(11))
8 BEGIN
9
10 -- Deklarieren die verwendeten Variablen.
11 DECLARE m INT(11);
12 DECLARE n INT(11);
13 DECLARE counter_1 INT(11);
14 DECLARE counter_2 INT(11);
15 DECLARE counter_3 INT(11);
16 DECLARE pivot DECIMAL(65, 30);
```

```

17
18 -- Bestimme die Dimensionsn für die Matrix X.
19 SET m = number_datapoints;
20 SET n = 3;
21
22 -- Lösche vorhandene temporäre Relationen.
23 DROP TEMPORARY TABLE IF EXISTS matrix_X;
24 DROP TEMPORARY TABLE IF EXISTS matrix_transposed;
25 DROP TEMPORARY TABLE IF EXISTS matrix_product_1;
26 DROP TEMPORARY TABLE IF EXISTS matrix_inverse;
27 DROP TEMPORARY TABLE IF EXISTS matrix_product_2;
28 DROP TEMPORARY TABLE IF EXISTS matrix_y;
29 DROP TEMPORARY TABLE IF EXISTS matrix_result;
30
31 -- Erstelle temporäre Relationen für die zu berechnenden Matrizen.
32 CREATE TEMPORARY TABLE matrix_X (
33     `row` INT(11),
34     `column` INT(11),
35     `value` DECIMAL(65, 30)
36 );
37 CREATE TEMPORARY TABLE matrix_transposed (
38     `row` INT(11),
39     `column` INT(11),
40     `value` DECIMAL(65, 30)
41 );
42 CREATE TEMPORARY TABLE matrix_product_1 (
43     `row` INT(11),
44     `column` INT(11),
45     `value` DECIMAL(65, 30)
46 );
47 CREATE TEMPORARY TABLE matrix_inverse (
48     `row` INT(11),
49     `column` INT(11),
50     `value` DECIMAL(65, 30)
51 );
52 CREATE TEMPORARY TABLE matrix_product_2 (
53     `row` INT(11),
54     `column` INT(11),
55     `value` DECIMAL(65, 30)
56 );
57 CREATE TEMPORARY TABLE matrix_y (
58     `row` INT(11),
59     `column` INT(11),
60     `value` DECIMAL(65, 30)
61 );
62 CREATE TEMPORARY TABLE matrix_result (

```

```
63     `row` INT(11),
64     `column` INT(11),
65     `value` DECIMAL(65, 30)
66 );
67
68 -- Füge Werte der unabhängigen Variablen in die Relation matrix_X ein.
69 SET @id = 0;
70
71 INSERT INTO matrix_X
72 SELECT
73     @id := (@id + 1) AS `row`,
74     1 AS `column`,
75     1 AS `value`
76 FROM sample
77 LIMIT number_datapoints;
78
79 SET @id = 0;
80
81 INSERT INTO matrix_X
82 SELECT
83     @id := (@id + 1) AS `row`,
84     2 AS `column`,
85     purchases AS `value`
86 FROM sample
87 LIMIT number_datapoints;
88
89 SET @id = 0;
90
91 INSERT INTO matrix_X
92 SELECT
93     @id := (@id + 1) AS `row`,
94     3 AS `column`,
95     age AS `value`
96 FROM sample
97 LIMIT number_datapoints;
98
99 -- Füge Werte der abhängigen Variable in die Relation matrix_y ein.
100 SET @id = 0;
101
102 INSERT INTO matrix_y
103 SELECT
104     @id := (@id + 1) AS `row`,
105     1 AS `column`,
106     money AS `value`
107 FROM sample
108 LIMIT number_datapoints;
```

```

109
110 -- Berechne matrix_transposed.
111 INSERT INTO matrix_transposed
112 SELECT
113     `column` AS `row`,
114     `row` AS `column`,
115     `value` AS `value`
116 FROM matrix_X;
117
118 -- Berechne matrix_product_1. Iteriere dazu über alle Zeilen und Spalten
119   ↳ der Ergebnismatrix.
120 SET counter_1 = 1;
121 WHILE counter_1 <= n DO
122
123     SET counter_2 = 1;
124
125     WHILE counter_2 <= n DO
126
127         -- Berechne den Wert des aktuellen Matrixelements.
128         INSERT INTO matrix_product_1 VALUES (
129             counter_1,
130             counter_2,
131             (
132                 SELECT SUM(matrix_X.`value` * matrix_transposed.`value`)
133                 FROM matrix_X, matrix_transposed
134                 WHERE matrix_X.`column` = counter_2
135                     AND matrix_transposed.`row` = counter_1
136                     AND matrix_transposed.`column` = matrix_X.`row`
137             )
138         );
139
140         SET counter_2 = counter_2 + 1;
141
142     END WHILE;
143
144     SET counter_1 = counter_1 + 1;
145
146 END WHILE;
147
148 -- Berechne matrix_inverse. Verwende dazu den in der Arbeit referenzierten
149   ↳ Algorithmus.
150 INSERT INTO matrix_inverse
151 SELECT *
152 FROM matrix_product_1;

```

```
153 SET counter_1 = 0;
154
155 WHILE counter_1 < n DO
156
157     SET counter_1 = counter_1 + 1;
158
159     DROP TEMPORARY TABLE IF EXISTS pivot_row;
160     CREATE TEMPORARY TABLE pivot_row (
161         `column` INT(11),
162         `value` DECIMAL(65, 30)
163     );
164
165     INSERT INTO pivot_row
166     SELECT `column`, `value`
167     FROM matrix_inverse
168     WHERE `row` = counter_1;
169
170     SET pivot = (
171         SELECT `value`
172         FROM matrix_inverse
173         WHERE `row` = counter_1 AND `column` = counter_1
174     );
175
176     UPDATE matrix_inverse
177     SET `value` = `value` / pivot
178     WHERE `row` = counter_1 AND `column` <> counter_1;
179
180     UPDATE matrix_inverse
181     SET `value` = - `value` / pivot
182     WHERE `row` <> counter_1 AND `column` = counter_1;
183
184     SET counter_2 = 1;
185
186     WHILE counter_2 <= n DO
187
188         IF counter_2 <> counter_1 THEN
189
190             SET counter_3 = 1;
191
192             WHILE counter_3 <= n DO
193
194                 IF counter_3 <> counter_1 THEN
195
196                     SET pivot = (
197                         SELECT `value`
198                         FROM pivot_row
```



```

199         WHERE `column` = counter_3
200     ) * (
201         SELECT `value`
202         FROM matrix_inverse
203         WHERE `row` = counter_2 AND `column` = counter_1
204     );
205
206     UPDATE matrix_inverse
207     SET `value` = `value` + pivot
208     WHERE `row` = counter_2 AND `column` = counter_3;
209
210     END IF;
211
212     SET counter_3 = counter_3 + 1;
213
214     END WHILE;
215
216     END IF;
217
218     SET counter_2 = counter_2 + 1;
219
220     END WHILE;
221
222     UPDATE matrix_inverse
223     SET `value` = 1 / `value`
224     WHERE `row` = counter_1 AND `column` = counter_1;
225
226     END WHILE;
227
228     -- Berechne matrix_product_2. Iteriere dazu über alle Zeilen der
229     ↪ Ergebnismatrix.
230     SET counter_1 = 1;
231
232     WHILE counter_1 <= n DO
233
234         -- Berechne den Wert des aktuellen Matrixelements.
235         INSERT INTO matrix_product_2 VALUES (
236             counter_1,
237             1,
238             (
239                 SELECT SUM(matrix_y.`value` * matrix_transposed.`value`)
240                 FROM matrix_y, matrix_transposed
241                 WHERE matrix_transposed.`row` = counter_1
242                     AND matrix_transposed.`column` = matrix_y.`row`
243             )
244         );

```

```
244
245     SET counter_1 = counter_1 + 1;
246
247 END WHILE;
248
249 -- Berechne matrix_result. Iteriere dazu über alle Zeilen der
    ↳ Ergebnismatrix.
250 SET counter_1 = 1;
251
252 WHILE counter_1 <= n DO
253
254     -- Berechne den Wert des aktuellen Matricelements.
255     INSERT INTO matrix_result VALUES (
256         counter_1,
257         1,
258         (
259             SELECT SUM(matrix_product_2.`value` * matrix_inverse.`value`)
260             FROM matrix_product_2, matrix_inverse
261             WHERE matrix_inverse.`row` = counter_1
262                   AND matrix_inverse.`column` = matrix_product_2.`row`
263         )
264     );
265
266     SET counter_1 = counter_1 + 1;
267
268 END WHILE;
269
270 -- Gib eine Relation mit Parameternamen und zugehörigen Werten zurück.
271 SELECT
272     CASE row
273         WHEN 1 THEN 'alpha'
274         WHEN 2 THEN 'beta_purchases'
275         WHEN 3 THEN 'beta_age'
276     END AS `variable`,
277     value
278 FROM matrix_result;
279
280 -- Lösche die temporären Relationen wieder.
281 DROP TEMPORARY TABLE IF EXISTS matrix_X;
282 DROP TEMPORARY TABLE IF EXISTS matrix_transposed;
283 DROP TEMPORARY TABLE IF EXISTS matrix_product_1;
284 DROP TEMPORARY TABLE IF EXISTS matrix_inverse;
285 DROP TEMPORARY TABLE IF EXISTS matrix_product_2;
286 DROP TEMPORARY TABLE IF EXISTS matrix_y;
287 DROP TEMPORARY TABLE IF EXISTS matrix_result;
288
```

```

289 END;;
290
291 DELIMITER ;

```

## D.3. Logistische Regression

```

1  -- Lösche die bestehenden Prozeduren, falls vorhanden.
2  DROP PROCEDURE IF EXISTS calculate_gradient;
3  DROP PROCEDURE IF EXISTS calculate_new_parameters;
4  DROP PROCEDURE IF EXISTS calculate_logit;
5  DROP PROCEDURE IF EXISTS logistic_regression;
6
7  DELIMITER ;;
8  -- Erstelle eine Prozedur zur Berechnung der Werte der logistischen
   ↳ Funktion.
9  CREATE PROCEDURE `calculate_logit`()
10 BEGIN
11
12  -- Deklariere die benötigten Variablen.
13  DECLARE alpha_old DECIMAL(65, 30);
14  DECLARE alpha_new DECIMAL(65, 30);
15
16  -- Bestimme den alten und neuen Wert von alpha.
17  SET alpha_old = (
18      SELECT old
19      FROM parameters
20      WHERE variable = 'alpha'
21  );
22  SET alpha_new = (
23      SELECT new
24      FROM parameters
25      WHERE variable = 'alpha'
26  );
27
28  -- Berechne die Werte der logistischen Funktion für alle Datenpunkte.
29  DELETE FROM logits;
30  INSERT INTO logits
31      SELECT
32          d.id,
33          1 / (1 + exp(- SUM(d.value * p.old) - alpha_old)) AS `old`,
34          1 / (1 + exp(- SUM(d.value * p.new) - alpha_new)) AS `new`
35      FROM datapoints d
36      JOIN parameters p ON p.variable = d.variable
37      GROUP BY d.id;
38

```

```
39  END;;
40
41  -- Erstelle eine Prozedur zur Berechnung des Gradienten.
42  CREATE PROCEDURE `calculate_gradient`()
43  BEGIN
44
45  DELETE FROM gradient;
46
47  -- Berechne die partielle Ableitung nach alpha.
48  INSERT INTO gradient
49  SELECT 'alpha' AS `variable`, SUM(bv.value - l.old) AS `value`
50  FROM logits l
51  JOIN binary_values bv ON bv.id = l.id;
52
53  -- Berechne die partielle Ableitung nach allen beta-Parametern.
54  INSERT INTO gradient
55  SELECT d.variable, SUM(d.value * (bv.value - l.old)) AS `value`
56  FROM logits l
57  JOIN binary_values bv ON bv.id = l.id
58  JOIN datapoints d ON d.id = l.id
59  GROUP BY d.variable;
60
61  END;;
62
63  -- Erzeuge eine Prozedur zur Berechnung der neuen Parameter abhängig von
64  -- ↳ der aktuellen Schrittweite.
65  CREATE PROCEDURE `calculate_new_parameters`(IN step DECIMAL(65, 30))
66  BEGIN
67
68  UPDATE parameters
69  JOIN gradient ON gradient.variable = parameters.variable
70  SET parameters.new = parameters.old + step * gradient.value;
71
72  END;;
73
74  -- Erstelle die Prozedur für logistische Regression.
75  CREATE PROCEDURE `logistic_regression`(IN number_datapoints INT(11), IN
76  -- ↳ rounds INT(11), step DECIMAL(65, 30))
77  BEGIN
78
79  -- Deklariere die verwendeten Variablen.
80  DECLARE min INT(11);
81  DECLARE max INT(11);
82  DECLARE transform DECIMAL(65, 30);
83  DECLARE better INT(1);
84  DECLARE counter INT(11);
```

```

83
84  -- Erstelle eine temporäre Relation für die Werte der unabhängigen
    ↳ Variablen.
85  DROP TEMPORARY TABLE IF EXISTS datapoints;
86  CREATE TEMPORARY TABLE datapoints (
87      id INT(11),
88      variable VARCHAR(32),
89      value DECIMAL(65, 30),
90      PRIMARY KEY (id, variable)
91  );
92
93  -- Berechne Minimum und Maximum der unabhängigen Variable.
94  SET min = (SELECT MIN(money) FROM sample);
95  SET max = (SELECT MAX(money) FROM sample);
96
97  -- Füge die linear transformierten Werte der unabhängigen Variablen in die
    ↳ Relation datapoints ein.
98  SET @counter = 0;
99  INSERT INTO datapoints
100  SELECT
101      @counter := @counter + 1 AS `id`,
102      'beta_money' AS `variable`,
103      (money - min) / (max - min) AS `value`
104  FROM sample
105  LIMIT number_datapoints;
106
107  -- Erstelle eine temporäre Relation für die (binären) Werte der abhängigen
    ↳ Variablen.
108  DROP TEMPORARY TABLE IF EXISTS binary_values;
109  CREATE TEMPORARY TABLE binary_values (
110      id INT(11),
111      value INT(1),
112      PRIMARY KEY (id)
113  );
114
115  -- Füge die Werte der abhängigen Variable ein.
116  SET @counter = 0;
117  INSERT INTO binary_values
118  SELECT
119      @counter := @counter + 1 AS `id`,
120      premium AS `value`
121  FROM sample
122  LIMIT number_datapoints;
123
124  -- Erstelle eine temporäre Relation für die alten und neuen
    ↳ Parameterwerte.

```

```
125 DROP TEMPORARY TABLE IF EXISTS parameters;
126 CREATE TEMPORARY TABLE parameters (
127     variable VARCHAR(32),
128     old DECIMAL(65, 30),
129     new DECIMAL(65, 30),
130     PRIMARY KEY (variable)
131 );
132
133 -- Füge die Initialwerte der Parameter ein.
134 INSERT INTO parameters VALUES
135     ('alpha', 0, 0),
136     ('beta_money', 0, 0);
137
138 -- Erstelle eine temporäre Relation für die Werte der logistischen
139 ↪ Funktion für alle Datenpunkte.
140 DROP TEMPORARY TABLE IF EXISTS logits;
141 CREATE TEMPORARY TABLE logits (
142     id INT(11),
143     old DECIMAL(65, 30),
144     new DECIMAL(65, 30),
145     PRIMARY KEY (id)
146 );
147
148 -- Befülle die Relation für die Werte der logistischen Funktion.
149 CALL calculate_logit();
150
151 -- Erstelle eine temporäre Relation für den Gradienten.
152 DROP TEMPORARY TABLE IF EXISTS gradient;
153 CREATE TEMPORARY TABLE gradient (
154     variable VARCHAR(32),
155     value DECIMAL(65, 30),
156     PRIMARY KEY (variable)
157 );
158
159 -- Iteriere über die Anzahl der gewünschten Iterationen.
160 SET counter = 0;
161 WHILE counter < rounds AND step > 0.000000000000000000000000000001 DO
162
163     -- Berechne den Gradienten und die neuen Parameter mit der aktuellen
164     ↪ Schrittweite.
165     CALL calculate_gradient();
166     CALL calculate_new_parameters(step);
167     CALL calculate_logit();
168
169     -- Verringere die Schrittweite solange, bis die neuen Parameter ein
170     ↪ besseres Ergebnis liefern als die alten.
```



```
212 DROP TEMPORARY TABLE IF EXISTS logits;
213 DROP TEMPORARY TABLE IF EXISTS gradient;
214
215 END;;
216 DELIMITER ;
```



## E. PostgreSQL-Skripte

### E.1. Einfache lineare Regression

```
1  -- Erstelle (oder ersetze falls vorhanden) die Prozedur für einfache
   ↳ lineare Regression.
2  CREATE OR REPLACE FUNCTION simple_linear_regression(number_datapoints
   ↳ INTEGER)
3  RETURNS TABLE (
4      variable VARCHAR(50),
5      value NUMERIC(65, 30)
6  ) AS $$
7  BEGIN
8
9      -- Erstelle eine Relation für die zu verwendenden Datenpunkte.
10     DROP TABLE IF EXISTS datapoints;
11     CREATE TEMPORARY TABLE datapoints (
12         purchases INTEGER,
13         money INTEGER
14     );
15
16     -- Füge die gewünschte Anzahl der Datenpunkte in die temporäre Relation
   ↳ ein.
17     INSERT INTO datapoints
18     SELECT purchases, money
19     FROM sample
20     LIMIT number_datapoints;
21
22     RETURN QUERY
23     WITH
24         -- Berechne die Mittelwerte der abhängigen und unabhängigen Variable.
25         means AS (
26             SELECT
27                 AVG(purchases) AS mean_purchases,
28                 AVG(money) AS mean_money
29             FROM datapoints
30         ),
31         -- Berechne die Summen im Nenner und Zähler der Formel für beta.
32         sums AS (
33             SELECT
```

```
34      SUM((purchases - mean_purchases) * (money - mean_money)) AS
      ↪ nominator,
35      SUM(POWER(purchases - mean_purchases, 2)) AS denominator
36  FROM datapoints, means
37 ),
38  -- Berechne beta.
39  beta AS (
40      SELECT
41          'beta'::VARCHAR(50) AS variable,
42          nominator / denominator AS value
43      FROM sums
44  ),
45  -- Berechne alpha.
46  alpha AS (
47      SELECT
48          'alpha'::VARCHAR(50) AS variable,
49          mean_money - beta.value * mean_purchases AS value
50      FROM means, beta
51  )
52  -- Gib eine Relation mit Parametername und zugehörigem Wert zurück.
53  SELECT *
54  FROM alpha
55  UNION
56  SELECT *
57  FROM beta;
58
59  -- Lösche die temporäre Relation mit den Datenpunkten wieder.
60  DROP TABLE IF EXISTS datapoints;
61
62  END;
63  $$ LANGUAGE plpgsql;
```

## E.2. Multiple lineare Regression

```
1  -- Erstelle Funktion für die Berechnung der transponierten Matrix.
2  CREATE OR REPLACE FUNCTION matrix_transpose(a NUMERIC(65, 30)[][])
3  RETURNS NUMERIC(65, 30)[][] AS $$
4  DECLARE
5      rows_a INTEGER := array_length(a, 1);
6      columns_a INTEGER := array_length(a, 2);
7      i INTEGER;
8      j INTEGER;
9      c NUMERIC(65, 30)[][];
10     new_row NUMERIC(65, 30)[];
```

```

11 BEGIN
12
13 -- Iteriere über alle Zeilen und Spalten der ursprünglichen Matrix.
14 i := 1;
15 WHILE i <= columns_a LOOP
16
17     j := 1;
18     WHILE j <= rows_a LOOP
19         -- Erzeuge ein Array mit der neuen Zeile der transponierten Matrix aus
20         --   ↳ der Spalte der ursprünglichen Matrix.
21         new_row[j] := a[j][i];
22         j := j + 1;
23     END LOOP;
24
25     -- Füge die Zeile in die Ergebnismatrix ein.
26     c := array_cat(c, array[new_row]);
27     i := i + 1;
28 END LOOP;
29 RETURN c;
30
31 END;
32 $$ LANGUAGE plpgsql;
33
34 -- Erstelle Funktion für die Berechnung des Produktes zweier Matrizen.
35 CREATE OR REPLACE FUNCTION matrix_multiplication(a NUMERIC(65, 30)[][], b
36   ↳ NUMERIC(65, 30)[][])
37 RETURNS NUMERIC(65, 30)[][] AS $$
38 DECLARE
39     rows_a INTEGER := array_length(a, 1);
40     columns_a INTEGER := array_length(a, 2);
41     columns_b INTEGER := array_length(b, 2);
42     new_row NUMERIC(65, 30)[];
43     c NUMERIC(65, 30)[][];
44     counter_1 INTEGER;
45     counter_2 INTEGER;
46     counter_3 INTEGER;
47 BEGIN
48     -- Iteriere über die Zeilen und Spalten der Ergebnismatrix.
49     counter_1 := 1;
50     WHILE counter_1 <= rows_a LOOP
51
52         counter_2 := 1;
53         WHILE counter_2 <= columns_b LOOP

```

```
54
55     -- Initiiere den Wert des aktuellen Elementes der Ergebnismatrix mit
56     ↪ 0.
57     new_row[counter_2] := 0;
58
59     -- Iteriere über die Summanden zur Berechnung des aktuellen
60     ↪ Matrixelements.
61     counter_3 := 1;
62     WHILE counter_3 <= columns_a LOOP
63         -- Addiere den aktuellen Summanden zum Wert des aktuellen
64         ↪ Elementes.
65         new_row[counter_2] := new_row[counter_2] + a[counter_1][counter_3] *
66         ↪ b[counter_3][counter_2];
67         counter_3 := counter_3 + 1;
68     END LOOP;
69
70     counter_2 := counter_2 + 1;
71     END LOOP;
72
73     c := array_cat(c, array[new_row]);
74     counter_1 := counter_1 + 1;
75     END LOOP;
76
77     RETURN c;
78
79 END;
80 $$ LANGUAGE plpgsql;
81
82 -- Erstelle Funktion für die Berechnung der inversen Matrix.
83 CREATE OR REPLACE FUNCTION matrix_inversion(a NUMERIC(65, 30)[][])
84 RETURNS NUMERIC(65, 30)[] AS $$
85 DECLARE
86     n INTEGER := array_length(a, 1);
87     p INTEGER := 0;
88     i INTEGER;
89     j INTEGER;
90     c NUMERIC(65, 30)[][] := a;
91     o NUMERIC(65, 30)[][];
92 BEGIN
93
94     -- Verwende den in der Arbeit referenzierten Algorithmus.
95     WHILE p < n LOOP
96
97         p := p + 1;
```

```

95   o := c;
96
97   j := 1;
98   WHILE j <= n LOOP
99       IF j <> p THEN
100         c[p][j] := c[p][j] / c[p][p];
101       END IF;
102       j := j + 1;
103   END LOOP;
104
105   i := 1;
106   WHILE i <= n LOOP
107       IF i <> p THEN
108         c[i][p] := - c[i][p] / c[p][p];
109       END IF;
110       i := i + 1;
111   END LOOP;
112
113   i := 1;
114   WHILE i <= n LOOP
115       IF i <> p THEN
116         j := 1;
117         WHILE j <= n LOOP
118             IF j <> p THEN
119                 c[i][j] := c[i][j] + o[p][j] * c[i][p];
120             END IF;
121             j := j + 1;
122         END LOOP;
123       END IF;
124       i := i + 1;
125   END LOOP;
126
127   c[p][p] := 1 / c[p][p];
128
129 END LOOP;
130
131 RETURN c;
132
133 END;
134 $$ LANGUAGE plpgsql;
135
136 -- Erstelle die Funktion für multiple lineare Regression.
137 CREATE OR REPLACE FUNCTION multiple_linear_regression(number_datapoints
138 ↪ INTEGER)
138 RETURNS TABLE (

```

```
139     variable VARCHAR(50),
140     value NUMERIC(65, 30)
141 ) AS $$
142 DECLARE
143     -- Erzeuge die Matrix X mit der gewünschten Anzahl an Datenpunkten.
144     x INTEGER[] [] := (
145         SELECT ARRAY(
146             SELECT ARRAY[1, purchases, age]
147             FROM sample
148             LIMIT number_datapoints
149         )
150     );
151     -- Erzeuge die Matrix y mit der gewünschten Anzahl an Datenpunkten.
152     y INTEGER[] := (
153         SELECT ARRAY(
154             SELECT ARRAY[money]
155             FROM sample
156             LIMIT number_datapoints
157         )
158     );
159     b NUMERIC(65, 30)[] [];
160 BEGIN
161
162     -- Berechne die Lösungsformel unter Verwendung der zuvor definierten
163     ↪ Funktionen.
164     b := matrix_multiplication(
165         matrix_inversion(
166             matrix_multiplication(
167                 matrix_transpose(x),
168                 x
169             ),
170         matrix_multiplication(
171             matrix_transpose(x),
172             y
173         )
174     );
175
176     -- Gib eine Relation mit Parameternamen und zugehörigen Werten zurück.
177     RETURN QUERY
178     SELECT 'alpha'::VARCHAR(50) AS variable, b[1][1] AS value
179     UNION
180     SELECT 'beta_purchases'::VARCHAR(50) AS variable, b[2][1] AS value
181     UNION
182     SELECT 'beta_age'::VARCHAR(50) AS variable, b[3][1] AS value;
```

```

183
184 END;
185 $$ LANGUAGE plpgsql;

```

### E.3. Logistische Regression

```

1  -- Erstelle (oder ersetze falls vorhanden) eine Prozedur zur Berechnung
2  ↳ der Werte der logistischen Funktion.
3  CREATE OR REPLACE FUNCTION calculate_logit()
4  RETURNS void AS $$
5  BEGIN
6
7
8  WITH
9      -- Bestimme den alten und neuen Wert von alpha.
10     alpha_old AS (
11         SELECT old
12         FROM parameters
13         WHERE variable = 'alpha'
14     ),
15     alpha_new AS (
16         SELECT new
17         FROM parameters
18         WHERE variable = 'alpha'
19     )
20     -- Berechne die Werte der logistischen Funktion für alle Datenpunkte.
21     INSERT INTO logits
22     SELECT
23         d.id,
24         1 / (1 + EXP(- SUM(d.value * p.old) - (SELECT old FROM alpha_old))) AS
25         ↳ old,
26         1 / (1 + EXP(- SUM(d.value * p.new) - (SELECT new FROM alpha_new))) AS
27         ↳ new
28     FROM datapoints d
29     JOIN parameters p ON p.variable = d.variable
30     GROUP BY d.id;
31
32 RETURN;
33
34 END;
35 $$ LANGUAGE plpgsql;

```

```
35  -- Erstelle (oder ersetze falls vorhanden) eine Prozedur zur Berechnung
    ↳ des Gradienten.
36  CREATE OR REPLACE FUNCTION calculate_gradient()
37  RETURNS void AS $$
38  BEGIN
39
40  DELETE FROM gradient;
41
42  -- Berechne die partielle Ableitung nach alpha.
43  INSERT INTO gradient
44  SELECT 'alpha' AS variable, SUM(bv.value - l.old) AS value
45  FROM logits l
46  JOIN binary_values bv ON bv.id = l.id;
47
48  -- Berechne die partielle Ableitung nach allen beta-Parametern.
49  INSERT INTO gradient
50  SELECT d.variable, SUM(d.value * (bv.value - l.old)) AS value
51  FROM logits l
52  JOIN binary_values bv ON bv.id = l.id
53  JOIN datapoints d ON d.id = l.id
54  GROUP BY d.variable;
55
56  RETURN;
57
58  END;
59  $$ LANGUAGE plpgsql;
60
61  -- Erzeuge (oder ersetze falls vorhanden) eine Prozedur zur Berechnung der
    ↳ neuen Parameter abhängig von der aktuellen Schrittweite.
62  CREATE OR REPLACE FUNCTION calculate_new_parameters(step NUMERIC(65, 30))
63  RETURNS void AS $$
64  BEGIN
65
66  UPDATE parameters
67  SET new = old + step * gradient.value
68  FROM gradient
69  WHERE gradient.variable = parameters.variable;
70
71  RETURN;
72
73  END;
74  $$ LANGUAGE plpgsql;
75
76  -- Erstelle die Prozedur für logistische Regression.
```



```

77 CREATE OR REPLACE FUNCTION logistic_regression(number_datapoints INTEGER,
↪ rounds INTEGER, step NUMERIC(65, 30))
78 RETURNS TABLE (
79     variable VARCHAR(50),
80     value NUMERIC(65, 30)
81 ) AS $$
82 DECLARE
83     counter INTEGER;
84 BEGIN
85
86     -- Erstelle eine Relation für die Werte der unabhängigen Variablen.
87     DROP TABLE IF EXISTS datapoints;
88     CREATE TEMPORARY TABLE datapoints (
89         id INTEGER,
90         variable VARCHAR(50),
91         value NUMERIC(65, 30)
92     );
93
94     -- Füge die linear transformierten Werte der unabhängigen Variablen in die
↪ Relation datapoints ein.
95     INSERT INTO datapoints
96     SELECT
97         row_number() OVER () AS id,
98         'beta_money' AS variable,
99         (money - (
100             SELECT MIN(money) FROM sample
101         ))::NUMERIC(65, 30) / ((
102             SELECT MAX(money) FROM sample
103         ) - (
104             SELECT MIN(money) FROM sample
105         ))::NUMERIC(65, 30) AS value
106     FROM sample
107     LIMIT number_datapoints;
108
109     -- Erstelle eine Relation für die (binären) Werte der abhängigen
↪ Variablen.
110     DROP TABLE IF EXISTS binary_values;
111     CREATE TEMPORARY TABLE binary_values (
112         id INTEGER,
113         value INTEGER
114     );
115
116     -- Füge die Werte der abhängigen Variable ein.
117     INSERT INTO binary_values
118     SELECT

```

```
119     row_number() OVER () AS id,
120     premium AS value
121 FROM sample
122 LIMIT number_datapoints;
123
124 -- Erstelle eine Relation für die alten und neuen Parameterwerte.
125 DROP TABLE IF EXISTS parameters;
126 CREATE TEMPORARY TABLE parameters (
127     variable VARCHAR(50),
128     old NUMERIC(65, 30),
129     new NUMERIC(65, 30)
130 );
131
132 -- Füge die Initialwerte der Parameter ein.
133 INSERT INTO parameters VALUES
134     ('alpha', 0, 0),
135     ('beta_money', 0, 0);
136
137 -- Erstelle eine Relation für die Werte der logistischen Funktion für alle
138     ↪ Datenpunkte.
139 DROP TABLE IF EXISTS logits;
140 CREATE TEMPORARY TABLE logits (
141     id INTEGER,
142     old NUMERIC(65, 30),
143     new NUMERIC(65, 30)
144 );
145
146 -- Befülle die Relation für die Werte der logistischen Funktion.
147 PERFORM calculate_logit();
148
149 -- Erstelle eine Relation für den Gradienten.
150 DROP TABLE IF EXISTS gradient;
151 CREATE TEMPORARY TABLE gradient (
152     variable VARCHAR(50),
153     value NUMERIC(65, 30)
154 );
155
156 -- Iteriere über die Anzahl der gewünschten Iterationen.
157 counter := 0;
158 WHILE counter < rounds AND step > 0.0000000000000000000000000001 LOOP
159
160     -- Berechne den Gradienten und die neuen Parameter mit der aktuellen
161     ↪ Schrittweite.
162     PERFORM calculate_gradient();
163     PERFORM calculate_new_parameters(step);
164     PERFORM calculate_logit();
```

```

163 -- Verringere die Schrittweite solange, bis die neuen Parameter ein
    ↳ besseres Ergebnis liefern als die alten.
165 WHILE NOT (
166     SELECT
167         SUM(LOG(bv.value * l.new + (1 - bv.value) * (1 - l.new))) >
168         SUM(LOG(bv.value * l.old + (1 - bv.value) * (1 - l.old)))
169     FROM logits l
170     JOIN binary_values bv ON bv.id = l.id
171 ) AND step > 0.000000000000000000000000000001 LOOP
172
173     step := step / 2;
174     PERFORM calculate_new_parameters(step);
175     PERFORM calculate_logit();
176
177 END LOOP;
178
179 -- Ersetze die alten Werte durch die neuen Werte.
180 UPDATE parameters
181 SET old = new;
182
183 UPDATE logits
184 SET old = new;
185
186 counter := counter + 1;
187
188 END LOOP;
189
190 -- Transformiere die Parameter linear, um den originalen Daten zu
    ↳ entsprechen.
191 UPDATE parameters
192 SET old = old / ((SELECT MAX(money) FROM sample) - (SELECT MIN(money) FROM
    ↳ sample))
193 WHERE parameters.variable = 'beta_money';
194
195 UPDATE parameters
196 SET old = old - (SELECT old FROM parameters p2 WHERE p2.variable =
    ↳ 'beta_money') * (SELECT MIN(money) FROM sample)
197 WHERE parameters.variable = 'alpha';
198
199 -- Gib eine Relation mit Parametername und zugehörigem Wert zurück.
200 RETURN QUERY
201 SELECT parameters.variable::VARCHAR(50), old AS value
202 FROM parameters;
203
204 -- Lösche die Relationen wieder.
```

```
205 DROP TABLE IF EXISTS datapoints;
206 DROP TABLE IF EXISTS binary_values;
207 DROP TABLE IF EXISTS parameters;
208 DROP TABLE IF EXISTS logits;
209 DROP TABLE IF EXISTS gradient;
210
211 END;
212 $$ LANGUAGE plpgsql;
```

## F. Python-Skripte für das Benchmarking

### F.1. Berechnung der Benchmarks

```
1 from time import time
2 from subprocess import call
3 import os
4 import sys
5 import csv
6 import json
7
8 # Erzeuge Funktion, um einen Fortschrittsbalken zu drucken.
9 def print_progressbar(iterations, progress):
10     sys.stdout.write("\033[100D")
11     for i in range(progress * 50 // iterations): sys.stdout.write("#")
12     for i in range(50 - (progress * 50 // iterations)):
13         ↪ sys.stdout.write(".")
14     sys.stdout.write(" || %i%% / 100%%" % (progress * 100 // iterations))
15     sys.stdout.flush()
16
17 # Erzeuge Funktion, die eine bestimmte Anzahl an Laufzeiten für eine
18 ↪ bestimmte Anzahl an Datenpunkten für eine bestimmte Art der Regression
19 ↪ in einer bestimmten Sprache berechnet.
20 def benchmark(type_regression, language, command, set_number_datapoints,
21 ↪ file, iterations):
22     print("Evaluating %s in %s..." % (type_regression, language))
23
24     # Iteriere über ein übergebenes Array mit den Anzahlen der zu
25     ↪ verwendenden Datenpunkte.
26     for number_datapoints in set_number_datapoints:
27         print("Use %s datapoints:" % (number_datapoints))
28
29         # Erzeuge eine Datei, in die die Ausgaben aus stdout und stderr
30         ↪ geschrieben werden.
31         logfile = open("logs/%s-%s-%i.txt" % (language,
32 ↪ type_regression.replace(" ", "-"), number_datapoints), "a")
33
34         # Füge die Anzahl der Datenpunkte in den Kommandozeilen-Befehl ein.
35         if command.count("%") == 2:
36             exec_command = command % (number_datapoints, 8 / number_datapoints)
37         else:
```

```
31     exec_command = command % number_datapoints
32
33     # Iteriere über die Anzahl der gewünschten Iterationen.
34     for i in range(iterations):
35         # Drucke den Fortschrittsbalken
36         print_progessbar(iterations, i)
37
38         # Speichere die Startzeit.
39         start_time = time()
40
41         # Führe den Kommandozeilen-Befehl aus.
42         call(exec_command, stdout = logfile, stderr = logfile, shell = True)
43
44         # Speichere die Endzeit.
45         end_time = time()
46
47         # Füge eine neue Zeile in der Benchmark-Datei ein.
48         file.write("%s,%s,%i,%s\n" % (
49             language,
50             type_regression,
51             number_datapoints,
52             (end_time - start_time)
53         ))
54
55         # Schließe die Log-Datei.
56         logfile.close()
57
58         # Drucke 100% im Fortschrittbalken und eine leere Zeile.
59         print_progessbar(iterations, iterations)
60         print("")
61     return
62
63 def main(argv):
64     # Bestimme default-Werte.
65     iterations = 100
66     set_number_datapoints = [10, 100, 1000, 10000, 100000]
67     run_r = False
68     run_tensorflow = False
69     run_mysql = False
70     run_postgresql = False
71     run_simple_linear_regression = False
72     run_multiple_linear_regression = False
73     run_logistic_regression = False
74
75     # Durchlaufe die übergebenen Argumente.
76     for arg in argv:
```

```

77     # Fixiere die Anzahl der Datenpunkte.
78     if "--datapoints=" in arg: set_number_datapoints =
    ↪     [int(arg.split("=")[1])]
79
80     # Setze die Anzahl der Iterationen.
81     if "--iterations=" in arg: iterations = int(arg.split("=")[1])
82
83     # Berechne Benchmarks für R.
84     if arg in ["--r", "-r"]: run_r = True
85
86     # Verwende TensorFlow.
87     if arg in ["--tensorflow", "-t"]: run_tensorflow = True
88
89     # Berechne Benchmarks für MySQL.
90     if arg in ["--mysql", "-m"]: run_mysql = True
91
92     # Berechne Benchmarks für PostgreSQL.
93     if arg in ["--postgresql", "-p"]: run_postgresql = True
94
95     # Berechne Benchmarks für einfache lineare Regression.
96     if arg in ["--simple-linear", "-slr"]: run_simple_linear_regression =
    ↪     True
97
98     # Berechne Benchmarks für multiple lineare Regression.
99     if arg in ["--multiple-linear", "-mlr"]:
    ↪     run_multiple_linear_regression = True
100
101     # Berechne Benchmarks für logistische Regression.
102     if arg in ["--logistic", "-lr"]: run_logistic_regression = True
103
104     # Wenn keine Art der Regression und keine Sprache spezifiziert wurde,
    ↪     berechne alles.
105     if not ((
106         run_r or
107         run_tensorflow or
108         run_mysql or
109         run_postgresql
110     ) and (
111         run_simple_linear_regression or
112         run_multiple_linear_regression or
113         run_logistic_regression
114     )):
115         run_r = True
116         run_tensorflow = True
117         run_mysql = True
118         run_postgresql = True

```

```
119     run_simple_linear_regression = True
120     run_multiple_linear_regression = True
121     run_logistic_regression = True
122
123     # Erzeuge eine Datei für die berechneten Laufzeiten und schreibe die
124     ↪ erste Zeile.
125     file = open("benchmarks-%i.csv" % (time()), "w")
126     file.write("language,type,datapoints,time\n")
127
128     # Berechne einfache lineare Regression in R.
129     if run_r and run_simple_linear_regression:
130         benchmark(
131             "simple linear regression",
132             "r",
133             "Rscript r/simpleLinearRegression.R %i -",
134             set_number_datapoints,
135             file,
136             iterations
137         )
138
139     # Berechne multiple lineare Regression in R.
140     if run_r and run_multiple_linear_regression:
141         benchmark(
142             "multiple linear regression",
143             "r",
144             "Rscript r/multipleLinearRegression.R %i -",
145             set_number_datapoints,
146             file,
147             iterations
148         )
149
150     # Berechne logistische Regression in R.
151     if run_r and run_logistic_regression:
152         benchmark(
153             "logistic regression",
154             "r",
155             "Rscript r/logisticRegression.R %i -",
156             set_number_datapoints,
157             file,
158             iterations
159         )
160
161     # Berechne einfache lineare Regression in TensorFlow.
162     if run_tensorflow and run_simple_linear_regression:
163         benchmark(
164             "simple linear regression",
```



```

164         "tensorflow",
165         "python3 tensorflow/simpleLinearRegression.py %i -",
166         set_number_datapoints,
167         file,
168         iterations
169     )
170
171     # Berechne multiple lineare Regression in TensorFlow.
172     if run_tensorflow and run_multiple_linear_regression:
173         benchmark(
174             "multiple linear regression",
175             "tensorflow",
176             "python3 tensorflow/multipleLinearRegression.py %i -",
177             set_number_datapoints,
178             file,
179             iterations
180         )
181
182     # Berechne logistische Regression in TensorFlow.
183     if run_tensorflow and run_logistic_regression:
184         benchmark(
185             "logistic regression",
186             "tensorflow",
187             "python3 tensorflow/logisticRegression.py %i -",
188             set_number_datapoints,
189             file,
190             iterations
191         )
192
193     # Berechne einfache lineare Regression in MySQL.
194     if run_mysql and run_simple_linear_regression:
195         benchmark(
196             "simple linear regression",
197             "mysql",
198             "echo \"CALL regression.simple_linear_regression(%i)\" | " + "mysql
199             ↪ -u %s -p%s" % (
200                 json.loads(os.environ["MYSQL_CONFIG"])["user"],
201                 json.loads(os.environ["MYSQL_CONFIG"])["password"]
202             ),
203             set_number_datapoints,
204             file,
205             iterations
206         )
207
208     # Berechne multiple lineare Regression in MySQL.
209     if run_mysql and run_multiple_linear_regression:

```

```
209     benchmark(
210         "multiple linear regression",
211         "mysql",
212         "echo \"CALL regression.multiple_linear_regression(%i)\" | " +
213         ↪ "mysql -u %s -p%s" % (
214             json.loads(os.environ["MYSQL_CONFIG"])["user"],
215             json.loads(os.environ["MYSQL_CONFIG"])["password"]
216         ),
217         set_number_datapoints,
218         file,
219         iterations
220     )
221
222     # Berechne logistische Regression in MySQL.
223     if run_mysql and run_logistic_regression:
224         benchmark(
225             "logistic regression",
226             "mysql",
227             "echo \"CALL regression.logistic_regression(%i, 1000, %f)\" | " +
228             ↪ "mysql -u %s -p%s" % (
229                 json.loads(os.environ["MYSQL_CONFIG"])["user"],
230                 json.loads(os.environ["MYSQL_CONFIG"])["password"]
231             ),
232             set_number_datapoints,
233             file,
234             iterations
235         )
236
237     # Berechne einfache lineare Regression in PostgreSQL.
238     if run_postgresql and run_simple_linear_regression:
239         benchmark(
240             "simple linear regression",
241             "postgresql",
242             "echo \"SELECT simple_linear_regression(%i)\" | psql regression",
243             set_number_datapoints,
244             file,
245             iterations
246         )
247
248     # Berechne multiple lineare Regression in PostgreSQL.
249     if run_postgresql and run_multiple_linear_regression:
250         benchmark(
251             "multiple linear regression",
252             "postgresql",
253             "echo \"SELECT multiple_linear_regression(%i)\" | psql regression",
254             set_number_datapoints,
```

```

253     file,
254     iterations
255 )
256
257 # Berechne logistische Regression in PostgreSQL.
258 if run_postgresql and run_logistic_regression:
259     benchmark(
260         "logistic regression",
261         "postgresql",
262         "echo \"SELECT logistic_regression(%i, 1000, %f)\" | psql
        ↪ regression",
263         set_number_datapoints,
264         file,
265         iterations
266     )
267
268 # Schließe die Benchmark-Datei.
269 file.close()
270 return
271
272 # Führe die main-Funktion aus.
273 if __name__ == "__main__":
274     main(sys.argv)

```

## F.2. Auswertung der Benchmarks

```

1  import sys
2  import os.path as p
3  import csv
4  import matplotlib.pyplot as plt
5
6  # Erzeuge Funktion, die die Benchmarks aus der csv-Datei einliest und
  ↪ gruppiert.
7  def calculate_benchmarks():
8      # Bestimme den Dateipfad der benchmark-Datei und öffne diese.
9      filename = p.abspath(p.join(p.dirname(p.realpath(__file__)),
        ↪ "benchmarks.csv"))
10     csvfile = open(filename, newline="")
11     csvreader = csv.reader(csvfile, delimiter=",", quotechar="|")
12
13     # Erstelle ein dictionary, in dem die Laufzeiten aggregiert werden
    ↪ sollen.
14     count = {
15         "simple linear regression": {
16             "10": {

```

```
17     "r": {"count": 0, "time": 0},
18     "tensorflow": {"count": 0, "time": 0},
19     "mysql": {"count": 0, "time": 0},
20     "postgresql": {"count": 0, "time": 0}
21 },
22 "100": {
23     "r": {"count": 0, "time": 0},
24     "tensorflow": {"count": 0, "time": 0},
25     "mysql": {"count": 0, "time": 0},
26     "postgresql": {"count": 0, "time": 0}
27 },
28 "1000": {
29     "r": {"count": 0, "time": 0},
30     "tensorflow": {"count": 0, "time": 0},
31     "mysql": {"count": 0, "time": 0},
32     "postgresql": {"count": 0, "time": 0}
33 },
34 "10000": {
35     "r": {"count": 0, "time": 0},
36     "tensorflow": {"count": 0, "time": 0},
37     "mysql": {"count": 0, "time": 0},
38     "postgresql": {"count": 0, "time": 0}
39 },
40 "100000": {
41     "r": {"count": 0, "time": 0},
42     "tensorflow": {"count": 0, "time": 0},
43     "mysql": {"count": 0, "time": 0},
44     "postgresql": {"count": 0, "time": 0}
45 }
46 },
47 "multiple linear regression": {
48     "10": {
49         "r": {"count": 0, "time": 0},
50         "tensorflow": {"count": 0, "time": 0},
51         "mysql": {"count": 0, "time": 0},
52         "postgresql": {"count": 0, "time": 0}
53     },
54     "100": {
55         "r": {"count": 0, "time": 0},
56         "tensorflow": {"count": 0, "time": 0},
57         "mysql": {"count": 0, "time": 0},
58         "postgresql": {"count": 0, "time": 0}
59     },
60     "1000": {
61         "r": {"count": 0, "time": 0},
62         "tensorflow": {"count": 0, "time": 0},
```

```

63     "mysql": {"count": 0, "time": 0},
64     "postgresql": {"count": 0, "time": 0}
65 },
66 "10000": {
67     "r": {"count": 0, "time": 0},
68     "tensorflow": {"count": 0, "time": 0},
69     "mysql": {"count": 0, "time": 0},
70     "postgresql": {"count": 0, "time": 0}
71 },
72 "100000": {
73     "r": {"count": 0, "time": 0},
74     "tensorflow": {"count": 0, "time": 0},
75     "mysql": {"count": 0, "time": 0},
76     "postgresql": {"count": 0, "time": 0}
77 }
78 },
79 "logistic regression": {
80     "10": {
81         "r": {"count": 0, "time": 0},
82         "tensorflow": {"count": 0, "time": 0},
83         "mysql": {"count": 0, "time": 0},
84         "postgresql": {"count": 0, "time": 0}
85     },
86     "100": {
87         "r": {"count": 0, "time": 0},
88         "tensorflow": {"count": 0, "time": 0},
89         "mysql": {"count": 0, "time": 0},
90         "postgresql": {"count": 0, "time": 0}
91     },
92     "1000": {
93         "r": {"count": 0, "time": 0},
94         "tensorflow": {"count": 0, "time": 0},
95         "mysql": {"count": 0, "time": 0},
96         "postgresql": {"count": 0, "time": 0}
97     },
98     "10000": {
99         "r": {"count": 0, "time": 0},
100        "tensorflow": {"count": 0, "time": 0},
101        "mysql": {"count": 0, "time": 0},
102        "postgresql": {"count": 0, "time": 0}
103    },
104    "100000": {
105        "r": {"count": 0, "time": 0},
106        "tensorflow": {"count": 0, "time": 0},
107        "mysql": {"count": 0, "time": 0},
108        "postgresql": {"count": 0, "time": 0}

```



```

148     "|" + " %s | %s | %s | %s | %s | %s " %
        ↳ tuple(table["tensorflow"]) + "|",
149     "|" + "-" * 89 + "|",
150     "|" + " %s | %s | %s | %s | %s | %s " %
        ↳ tuple(table["mysql"]) + "|",
151     "|" + "-" * 89 + "|",
152     "|" + " %s | %s | %s | %s | %s | %s " %
        ↳ tuple(table["postgresql"]) + "|",
153     "|" + "-" * 89 + "|"
154 ]
155
156 # Drucke die Tabelle und eine Leerzeile danach.
157 print(str.join("\n", print_table))
158 print("\n")
159
160 # Gib das dictionary mit allen aggregierten Laufzeiten zurück.
161 return count
162
163 # Erzeuge Funktion, um die Laufzeiten für eine bestimmte Art der
        ↳ Regression zu plotten.
164 def plot(benchmarks, regression_type, plot_title):
165     # Definiere ein Array mit den Anzahlen der Datenpunkte für den Plot.
166     x = [10, 100, 1000, 10000, 100000]
167
168     # Erzeuge ein dictionary mit leeren Arrays für die durchschnittlichen
        ↳ Laufzeiten.
169     values = {
170         "r": [],
171         "tensorflow": [],
172         "mysql": [],
173         "postgresql": []
174     }
175
176     # Durchlaufe alle Sprachen.
177     for language in ["r", "tensorflow", "mysql", "postgresql"]:
178         # Durchlaufe die Anzahl der Datenpunkte.
179         for number_datapoints in x:
180             # Füge die durchschnittliche Laufzeit in das Array ein (falls
                ↳ Laufzeiten vorhanden sind).
181             if
                ↳ benchmarks[regression_type][str(number_datapoints)][language]["count"]
                ↳ > 0:
182                 values[language].append(
183                     ↳ benchmarks[regression_type][str(number_datapoints)][language]["time"]
                    ↳ /

```

```
184         ↪ benchmarks[regression_type][str(number_datapoints)][language]["count"]
185     )
186     else:
187         values[language].append(None)
188
189     # Erzeuge den Plot.
190     plt.loglog(x, values["r"], "r-", label="R")
191     plt.loglog(x, values["tensorflow"], "y-", label="TensorFlow")
192     plt.loglog(x, values["mysql"], "b-", label="MySQL")
193     plt.loglog(x, values["postgresql"], "g-", label="PostgreSQL")
194     plt.title(plot_title)
195     plt.legend()
196     plt.xlabel("Anzahl der Datenpunkte")
197     plt.ylabel("Laufzeit in Sekunden")
198     plt.show()
199
200 def main(argv):
201     # Erzeuge default-Wert, ob Plots erstellt werden sollen.
202     print_plots = True
203
204     # Überschreibe default-Wert, falls ein entsprechendes Argument übergeben
205     ↪ wurde.
206     if len(argv) == 2:
207         if argv[1] == "-":
208             print_plots = False
209
210     # Berechne und drucke die Benchmarks.
211     benchmarks = calculate_benchmarks()
212
213     # Plote die Benchmarks (falls gewünscht).
214     if print_plots:
215         plot(benchmarks, "simple linear regression", "Einfache lineare
216             ↪ Regression")
217         plot(benchmarks, "multiple linear regression", "Multiple lineare
218             ↪ Regression")
219         plot(benchmarks, "logistic regression", "Logistische Regression")
220
221     # Führe die main-Funktion aus.
222     if __name__ == "__main__":
223         main(sys.argv)
```



# Bibliografie

- [1] Khan Hamid Ahmad Farooq. An efficient and simple algorithm for matrix inversion. 2010.
- [2] Ludwig Fahrmeir. *Statistik*. Springer Spektrum, 2016.
- [3] Trevor John Hastie John Mckinley Chambers. *Statistical models in S*. Wadsworth & Brooks/Cole, 1992.