

Inhaltsverzeichnis

1	Pakete	1
2	Arbeitsverzeichnis	3
3	Import von Datensätzen	3
3.1	Datenimport von Textdateien	4
3.2	Datenimport von Excel-Dateien	6
3.3	.RData - Das R-eigene Datenformat	7
4	Übersicht	8
4.1	Neue wichtige Konzepte	8
4.2	Neue wichtige Befehle, Argumente, Operatoren	8

1 Pakete

Wie im ersten Kapitel beschrieben, ist R quelloffen und kostenlos, was R in- und außerhalb der wissenschaftlichen *Community* sehr attraktiv macht. Prinzipiell kann jeder mit den nötigen Kenntnissen die Funktionalität von R erweitern und mit anderen teilen. Diese Erweiterungen werden Pakete genannt. Pakete beinhalten meist Funktionen. Sie können aber auch Objekte oder ganze Datensätze enthalten. Bei der Installation von R werden bereits eine Reihe von Paketen installiert (z.B., **base**, **stats** oder **graphics**), welche beim Start von R automatisch zur Verfügung stehen. Das gilt nur für die *R-Basispakete*; alle Pakete, die Sie manuell installieren, müssen mit jeder neuen Session auch manuell aktiviert werden. Im Folgenden lernen wir, wie R-Pakete installiert, aktiviert und deren Funktionen genutzt werden können.

1.0.1 Installation von Paketen

Es gibt viele Wege, um ein Paket in R zu installieren. Die klassische Variante ist über die *Console*. Der Name des gewünschten Pakets sollte bekannt sein und auch bei den Paketnamen ist die Unterscheidung von Groß- und Kleinbuchstaben wichtig. Außerdem dürfen Paketnamen keine Leerzeichen beinhalten. Die Installation zusätzlicher Pakete sieht in allgemeiner Form so aus:

```
install.packages("Paketname")
```


Ein praktisches Beispiel:

```
install.packages("car", dependencies = TRUE)
```

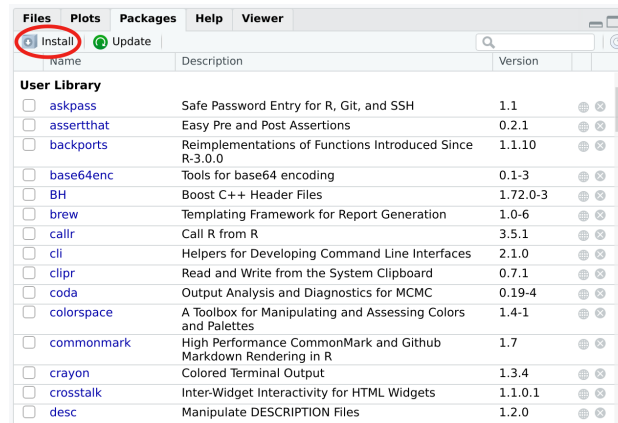
Als Erfolgsmeldung wartet man in der Konsole auf das schwarz gedruckte

```
package 'car' successfully unpacked and MD5 sums checked
```

Das Paket befindet sich nun in der sogenannten R Library des Computers und es kann darauf zugegriffen werden (s. u.). Mit dem Argument `dependencies = TRUE` geben wir R das Recht, weitere Pakete zu installieren, die für die Benutzung des Zielpakets benötigt werden.

Alternative 1: RStudio bietet einen weiteren Weg zur Installation von Paketen an. Über  **Install Packages** öffnet sich ein Fenster, in welches lediglich der Name des Paketes eingetippt werden muss.

Alternative 2: Ein letzter Weg ist über den **Packages**-Reiter im Fenster unten rechts. Hier kann über den **Install**-Button dasselbe Fenster wie in *Alternative 1* aufgerufen werden. Zusätzlich erhält man eine Übersicht über alle installierten Pakete.



Häufige Probleme:

- Zu Beginn wird vom CRAN (CRAN steht für The Comprehensive R Archive Network. R sucht nach dem gewünschten Paket auf CRAN und lädt es herunter, vorausgesetzt es existiert) manchmal nach dem bevorzugten *Mirror* gefragt. Hier einfach *Germany* und eine möglichst nahe Stadt (z.B., Göttingen) auswählen, dort steht dann der Server, auf dem die Dateien liegen.
- Pakete veranlassen oft die Installation weiterer Pakete, auf die sie angewiesen sind. Bei Versionswechseln (ca. alle 3-6 Monate) kann es sein, dass nicht alle Pakete gleich aktuell sind. Bei der Frage, ob aus neueren source-Files kompiliert werden soll, **NEIN** angeben.
- Falls Pakete aktueller sind als die eigene R Version, wird eine meist unproblematische Warnmeldung ausgegeben: **Warnung: R Paket wurde unter R version X.Y.Z erstellt.**

1.0.2 Verwendung von Paketen

Nach der erfolgreichen Installation ist ein Paket nicht automatisch in R geladen. Bevor wir also die zusätzliche Funktionalität unseres neuen Pakets nutzen können, müssen wir das Paket für uns zugänglich machen. Das geht wie folgt:

```
library(Paketname)
```

```
library(car)
```

Nun können Sie alle Funktionen des Pakets **car** in R nutzen. Einige dieser Funktionen werden Sie im Laufe des Semesters kennenlernen.

1.0.3 Aktualisieren von Paketen

Wie R werden auch Pakete mit neuen Funktionen ausgestattet oder *Bugs* behoben (Bugs¹ sind im Informatikerjargon Fehler in einem Programm, die meist durch den geschriebenen Code verursacht wurden). Aktualisieren kann man Pakete mit dem `update.packages()`-Befehl. Wenn der Befehl ohne weitere Argumente ausgeführt wird, wird für jedes Paket explizit gefragt, ob es aktualisiert werden soll. Falls man diesen Schritt überspringen will, hilft das Argument `ask`.

```
update.packages(ask = FALSE) # Alle Pakete werden ohne explizite Nachfrage aktualisiert
```

- Achtung! Das kann abhängig von der Anzahl der installierten Pakete sehr lange dauern.
 - Falls man nur ein bestimmtes Paket aktualisieren möchte, kann man einfach wieder mit `install.packages` arbeiten

```
install.packages("Paket")
```

¹Herkunft des Wortes Bug: Das Wort Bug kommt aus dem Englischen für Insekt oder Käfer. Die Legende besagt, dass der erste Computerfehler entstand, weil eine Motte in einen Schalter geflogen ist, was beim Computer zu Fehlverhalten führte.

2 Arbeitsverzeichnis

- Das Arbeitsverzeichnis (*working directory*, wd) ist der Dateiordner, auf den R standardmäßig zugreift, wenn Objekte gespeichert oder geöffnet werden sollen.
- Die Definition des eigenen Arbeitsverzeichnisses zu Beginn jeder Sitzung verkürzt die Befehle für Einlesen, Speichern und Öffnen von Daten.
- Wenn RStudio über Doppelklick auf ein Skript geöffnet wird, ist der Ordner, in dem das Skript liegt, das Arbeitsverzeichnis.
- Das aktuelle Arbeitsverzeichnis kann man mithilfe der Funktion `getwd()` abfragen:

```
getwd()
```

```
[1] "C:/Users/Icke/Documents"
```

- Hier ist der Dateipfad des aktuellen Arbeitsverzeichnisses angezeigt, also dem Ordner (hier: Documents) mit Überordnern und Laufwerk auf den R automatisch zugreift, wenn wir Dateien öffnen oder speichern wollen.
- Das Arbeitsverzeichnis sollte zu Beginn jeder R-Sitzung (nach jedem Öffnen des Programms) gesetzt werden. Das Arbeitsverzeichnis kann auch während einer Sitzung geändert werden.
- Im Gegensatz zu Windows-Dateipfaden beachte man die Schreibweise mit Forward-Slash(!) als Trenner
- Das Arbeitsverzeichnis lässt sich ändern über das RStudio Hauptmenü Session Set Working Directory Choose Directory ...
 - Dialogfenster öffnet sich
 - zum gewünschten Ordner navigieren
 - *Open* klicken
 - Der Befehl `setwd("Ordnerpfad")` wird in die Console geschrieben und ausgeführt.

```
setwd("X:/Pfad/zum/Ordner")
```

Falls Sie den Befehl händisch in die Konsole eintippen möchten, vergessen Sie nicht den Ordnerpfad in Anführungszeichen zu setzen (" " oder ' ').

```
setwd("Ordnerpfad")
```

```
getwd()
```

- Mit `getwd()` können wir überprüfen, ob das Anlegen des Arbeitsverzeichnisses geklappt hat.
 - Das Anlegen des Arbeitsverzeichnisses war erfolgreich, falls `getwd()` den in `setwd()` angegebenen Pfad ausgibt
- Die Ausgabe von `setwd()` sollte aus der *Console* in das Skript kopiert werden, damit das Skript beim erneuten Öffnen weiterhin fehlerfrei läuft (falls in diesem auf Dateien aus Ordnern auf dem Computer zugegriffen wird).

3 Import von Datensätzen

Wir arbeiten in den folgenden Sitzungen mit Datensätzen aus psychologischen Befragungen. Diese Daten liegen meist im Excelformat (.xlsx-Datei) oder einem Textdatenformat vor (z. B. .dat oder .txt).

Das in R verwendete Datenformat nennt sich *RData* (oder abgekürzt *.rda*). Um Daten im *RData* Format speichern zu können, müssen wir sie zunächst aus den anderen Dateiformaten in R **importieren**. Wenn Daten im *RData* Format vorliegen, können sie hingegen einfach **geöffnet** werden.

Das korrekte Importieren ist extrem wichtig um sicherzustellen, dass man mit den richtigen Daten weiterarbeitet. Gleichzeitig ist es manchmal relativ komplex. Die gute Nachricht ist, dass man es meist nur einmal machen muss, denn danach lässt es sich einfach mit der *RData* Datei weiterarbeiten. Die Dateien **erstis.xlsx**

(Excel-Datei) und **erstis.dat** (Textdatei) enthalten Daten von Studierenden der FU aus früheren Semestern und stammen aus dem Buch *R für Einsteiger* (Luhmann, 2015)². Das beigelegte *Codebook* gibt Aufschluss über die Variablenbenennungen und Ausprägungen.

Wir üben den Datenimport im Folgenden in zwei Varianten:

- Einlesen einer Text-Datei über die allgemeine `read.table()`-Funktion
- Einlesen der Excel-Datei mit Hilfe des RStudio-Menüs und Zusatzpaketen

3.1 Datenimport von Textdateien

- Textdatei-Formate sind universell und können in R mit der `read.table()`-Funktion importiert werden.
- Um die Pfadangabe zu sparen, sollte das Arbeitsverzeichnis auf den Ordner gelegt werden, der den Datensatz (als Textdatei) enthält.

In der R-Hilfe für die Funktion `read.table()` finden wir u. a. folgende Voreinstellungen:

```
read.table(file,
            header = FALSE,
            sep = "",
            dec = ".",
            na.strings = "NA")
```

- Diese Voreinstellungen müssen auf den Datensatz angepasst werden!
- Bei `file` muss der Name der Textdatei (inklusive Endung) angegeben werden.
- Mit diesem Befehl wird der Inhalt der Datendatei eingelesen. Um damit im Anschluss weiterarbeiten zu können, müssen wir den Inhalt in einem Objekt ablegen, z. B.

```
Daten <- read.table(...)
```

Die wichtigsten Parameter bei der Verwendung von `read.table` sind: `read.table(file, header, sep, dec, na.strings)`

Diese bedeuten der Reihe nach:

1. **file** gibt den Namen der Datei an, welche den einzulesenden Datensatz enthält. Es wird dabei im Arbeitsverzeichnis nach dieser Datei gesucht, sofern kein vollständiger Pfad (über eine URL) angegeben wird.
2. **header** ist ein Wahrheitswert, kann also die Werte `TRUE` oder `FALSE` übergeben bekommen. Ist dieser Wert `TRUE`, so wird die erste Zeile des Datensatzes zur Benennung der Variablen verwendet. Die Angabe dieses Werts ist nicht notwendig. Wird er nicht angegeben, so wird dieser Wert per Voreinstellung auf `FALSE` gesetzt.
3. **sep** gibt das Zeichen an, welches im jeweiligen Datensatz den Trenner zwischen dem Ende eines Wertes und dem Anfang eines anderen Wertes darstellt. Wird diese Angabe nicht gemacht, wird davon ausgegangen, dass *white space* (ein Leerzeichen) verschiedene Werte voneinander trennt.
4. **dec** gibt das Zeichen an, welches im jeweiligen Datensatz als Dezimaltrennzeichen dient. Wird diese Angabe nicht gemacht, wird davon ausgegangen, dass der Punkt als Dezimaltrennzeichen verwendet wird.
5. **na.strings** ermöglicht die Angabe benutzerdefinierter fehlender Werte. Per Voreinstellung werden diese mit `NA` gekennzeichnet.

²Luhmann, M. 2015. "R für Einsteiger. Einführung in Die Statistiksoftware für Die Sozialwissenschaften. Beltz." Weinheim, Basel.

erstis.dat - Notepad

File Edit Format View Help

Dateiname: **erstis.dat**

Variablenname: **code**, **gruppe**, **geschl**

Tabstopp als Spaltentrenner

Fehlende Werte

	code	gruppe	geschl	gebjahr	alter	abi	kinder	job	berlin	wohntort.alt	uni1	uni2	uni3	uni4	uni5	uni6
1	1	1	1	1970	38	1992	1	2	1	2	1	1	1	1	1	0
2	2	3	2	1967	41	1987	1	1	1	3	1	1	1	1	1	0
3	1	2	1	1989	19	2007	2	1	1	4	0	0	1	1	1	0
4	2	2	2	1976	32	1993	2	2	1	3	0	1	1	0	0	0
5	4	-9	-9	1969	39	1987	1	2	1	3	0	0	0	0	0	0
6	4	1	1	1976	32	1995	1	1	1	3	0	1	1	1	1	0
7	3	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9
8	2	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9
9	3	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9
10	4	2	1	1985	23	2005	2	2	1	3	0	1	1	1	1	1
11	3	2	1	1983	25	2002	2	1	1	3	1	1	1	1	1	0
12	2	1	1	1985	23	2005	2	2	1	3	1	1	1	1	1	1
13	2	1	1	1985	23	2005	2	2	1	3	1	1	1	1	0	0
14	2	1	1	1988	20	2008	2	1	2	2	0	1	1	1	1	0
15	3	1	1	1980	28	2000	2	2	1	3	1	1	1	1	1	1
16	3	1	1	1983	25	-9	2	1	-9	1	1	0	0	1	1	0
17	3	2	1	1983	25	-9	2	1	-9	1	0	0	1	1	0	0
18	1	1	1	1977	31	1996	2	1	1	1	1	1	1	1	0	0
19	1	2	1	1976	32	1994	2	2	1	-9	0	0	0	0	0	0
20	4	1	1	1985	23	2006	2	2	1	1	1	1	1	1	1	0
21	2	1	1	1989	19	2008	2	2	1	3	0	1	1	1	1	0
22	3	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9
23	4	2	1	1990	18	2008	2	2	1	1	0	1	1	1	0	0
24	3	1	1	1978	30	1998	2	2	2	3	0	1	0	1	1	1
25	1	1	1	1967	41	1980	1	2	2	2	0	1	1	0	1	0
26	1	2	1	1965	43	1980	1	2	2	2	0	0	0	1	1	0
27	1	1	1	1983	25	2003	2	2	1	1	0	1	1	0	1	0
28	2	1	1	1987	21	2007	2	2	2	4	0	0	1	1	1	0
29	4	1	1	1988	20	2008	2	1	1	3	0	0	0	1	0	0
30	1	1	1	1987	21	2007	2	1	1	2	0	1	1	1	1	0
31	4	1	1	1987	21	2007	2	1	1	4	0	1	1	1	1	0
32	1	2	1	1972	36	1990	1	2	1	3	1	1	1	1	1	0
33	4	1	1	1986	22	2006	2	1	1	3	0	1	0	0	1	1
34	3	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9	-9
35	4	2	1	1988	20	2008	2	2	1	3	0	1	1	0	1	0
36	1	1	1	1990	18	2008	2	2	1	1	0	0	1	1	1	0
37	1	1	1	1989	19	2008	2	2	1	1	0	1	1	1	1	0
38	1	2	1	1989	19	2008	2	2	1	2	0	0	1	1	1	0
39	3	1	1	1989	19	2008	2	2	1	4	1	1	1	1	1	0
40	1	1	1	1983	25	2002	2	2	1	1	0	0	1	1	1	0
41	2	2	1	1988	20	2007	2	1	1	3	0	1	1	1	1	0
42	2	2	1	1978	30	1998	1	1	1	3	0	1	1	1	1	1
43	3	1	1	1988	20	2007	2	1	1	3	0	1	0	1	1	0
44	2	1	1	1988	20	2008	2	2	1	3	0	1	0	0	1	0
45	3	1	1	1989	19	2007	2	2	1	3	0	0	1	0	0	1
46	2	1	1	1977	31	1998	1	2	1	3	1	0	1	0	1	0
47	4	1	1	1985	23	2005	2	2	1	3	1	1	1	1	1	1
48	1	1	1	1983	25	2003	2	1	1	1	1	1	1	1	1	0
49	2	1	1	1980	28	2007	1	2	1	3	1	0	0	1	1	0
50	4	1	1	1989	19	2008	2	2	1	1	0	1	1	1	1	0
51	2	2	1	1985	23	2005	1	1	1	-9	-9	0	1	0	1	1
52	1	1	1	1985	23	2005	2	2	1	1	0	1	1	1	1	1

Ln 1, Col 1 100% Windows (CRLF) UTF-8

Abbildung 1: Darstellung einer Textdatei

Struktur einer .txt-Datei	Argumente für <code>read.table()</code>	
Variablenamen in 1. Zeile	Nein Ja	<code>header = FALSE</code> (default) <code>header = TRUE</code>
Spaltentrennung	Leerzeichen Tabstopp Semikolon	<code>sep = ""</code> (default) <code>sep = "\t"</code> <code>sep = ";"</code>
Dezimaltrennzeichen in Textdatei	Punkt Komma	<code>dec = "."</code> (default) <code>dec = ","</code>
Kodierung fehlender Werte	NA -99 -99 und -9	<code>na.strings = "NA"</code> (default) <code>na.strings = "-99"</code> <code>na.strings = c("-99", "-9")</code>

3.2 Datenimport von Excel-Dateien

Für den Import von Excel-Dateien gibt es verschiedene Möglichkeiten. In RStudio hat man mit **File** **Import Dataset** **From Excel** eine integrierte Bedienoberfläche für den Import von Excel-Dateien. Dafür muss man aber das Paket `readxl` installieren und aktivieren:

```
install.packages("readxl")
library(readxl)
```

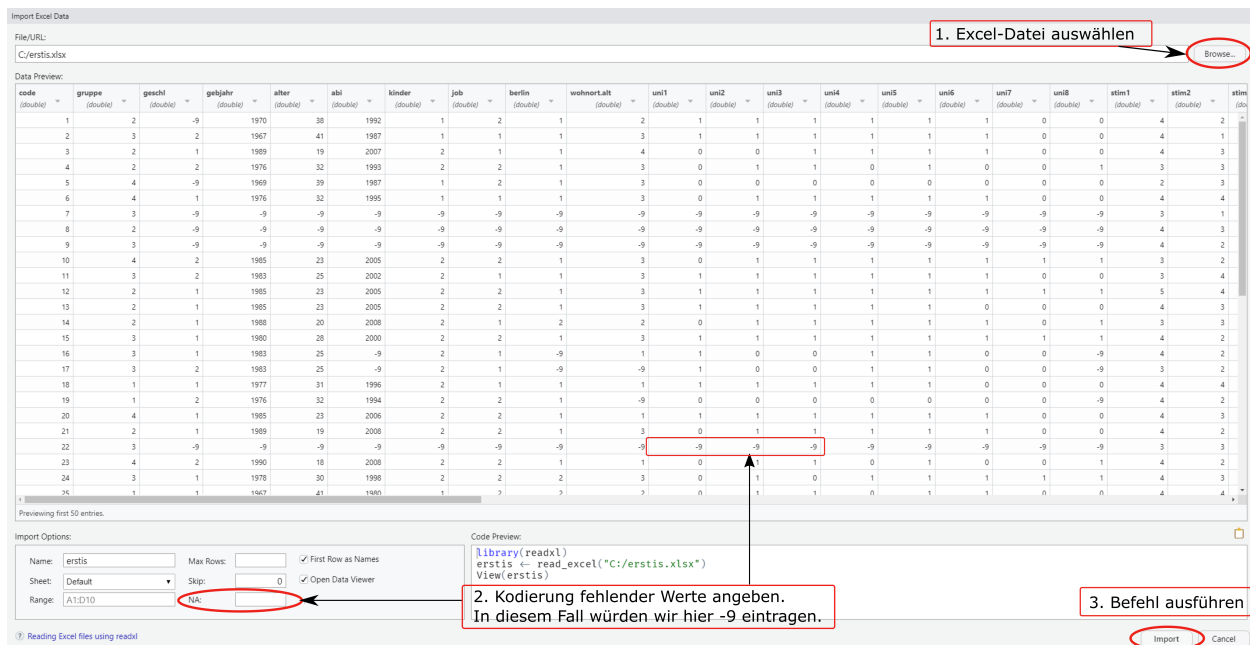


Abbildung 2: Datenimport in RStudio

Über **Browse** können die gewünschten Daten auf dem Computer lokalisiert und in die Vorschau geladen werden (Abbildung 2). Falls die Vorschau die Daten richtig widerspiegelt, können, über **Import**, die Daten in R importiert werden.

Ein alternatives Paket ist `gdata`, das ebenfalls für den Import von Excel-Dateien programmiert wurde. Hier fehlt die graphische Bedienoberfläche, allerdings lassen sich komplexere Datenimporte durchführen. Der Name der Funktion und die wichtigsten Parameter sind hier aufgelistet:

```
read.xls(xls,
        sheet = 1,
        na.strings = "NA")
```

Diese bedeuten der Reihe nach:

1. **xls** gibt den Namen der Datei ein, welche den einzulesenden Datensatz enthält. Es wird dabei im Arbeitsverzeichnis nach dieser Datei gesucht, sofern kein vollständiger Pfad (über eine **URL**) angegeben wird.
2. **sheet** gibt die Nummer der Tabelle oder des Tabellen-Sheets an, welches eingelesen werden soll. Wird die Nummer nicht angegeben, so wird die erste Tabelle eingelesen.
3. **na.strings** ermöglicht die Angabe benutzerdefinierter fehlender Werte (wie in `read.table`).
Ein Beispiel: Folgender Befehl liest die zweite Tabelle in der Excel-Datei **daten.xls** ein und speichert sie unter dem Namen **daten_excel**:

```
daten_excel <- read.xls("daten_excel ", sheet = 2, na.strings = "-9")
```

3.3 .RData - Das R-eigene Datenformat

Wenn Objekte, die in R erstellt oder verändert wurden, im Arbeitsverzeichnis abgespeichert werden sollen, bietet sich der **save**-Befehl an.

```
save(Objektname, file = "Dateiname.RData")
```

Der gewünschte Dateiname muss in Anführungsstrichen und mit der Dateiendung ***.RData** angegeben werden. Generell macht es Sinn, den Objektnamen auch als Dateinamen zu verwenden.

Um ein ***.RData**-Objekt in R zu laden, benutzt man den **load**-Befehl.

```
load("Dateiname.RData", verbose = TRUE)
```

Beispiel:

```
objekt <- c(1, 3, 5)
save(objekt, file = "objekt.RData")
load("objekt.RData", verbose = TRUE)
```

Das optionale **verbose**-Argument bewirkt, dass der Objektnamen in die Konsole geschrieben wird.

4 Übersicht

4.1 Neue wichtige Konzepte

- Pakete
- Import von Daten
- Arbeitsverzeichnis

4.2 Neue wichtige Befehle, Argumente, Operatoren

Funktion	Verwendung
<code>install.packages("paket")</code>	Installiert Paket paket in die lokale R Library
<code>library(paket)</code>	Aktiviert Paket paket zur Verwendung
<code>update.packages()</code>	Aktualisiert installierte Pakete
<code>read.table("Daten")</code>	Liest eine Textdatei in R ein
<code>read.xls("Daten")</code>	List eine Excel-Datei in R ein
<code>getwd()</code>	Gibt das aktuelle Arbeitsverzeichnis an
<code>setwd("Pfad")</code>	Legt den Pfad des Arbeitsverzeichnisses fest
<code>load("Daten")</code>	List eine .RData in R ein
<code>save("Dateiname.RData")</code>	Objekte aus R im Arbeitsverzeichnis speichern