

Schummelzettel: Textverarbeitung mit R und Sonderzeichen

<div>1. Pakete</div> <div>2.</div> <div>stringr: Für die meisten Textmanipulationen. readr: Für robusten Datenimport, wichtig für korrektes Encoding.</div> <div>2. Encoding</div> <div>UTF-8 als Standard verwenden: Sowohl in RStudio als auch beim Datenimport. Encoding beim Import angeben: Besonders wichtig bei Dateien von anderen Betriebssystemen.</div> <div><ul style="list-style-type: none">Mit locale Funktion im readr Paket:<code>library(readr)</code><code>text <- read_csv("datei.csv", locale = locale(encoding = "ISO-8859-1"))</code>Mit fileEncoding bei read.table für Windows:<code>read.table("Quelle.csv", sep="," , fileEncoding="UTF-8")</code></div> <div>Encoding überprüfen: <code>Encoding(text\$spalte)</code> Encoding ändern: <code>Encoding(text\$spalte) <- "latin1"</code> Probleme erkennen: Zeichen wie <fc> oder <e4> deuten auf Encoding-Fehler hin.</div>	<div>3. stringr Funktionen</div> <div>Groß-/Kleinschreibung:</div> <div><ul style="list-style-type: none"><code>str_to_upper(text)</code>: Alle Buchstaben groß.<code>str_to_lower(text)</code>: Alle Buchstaben klein.<code>str_to_title(text)</code>: Erster Buchstabe jedes Wortes groß (Title Case). Teile extrahieren/ersetzen:<code>str_sub(text, start, end)</code>: Extrahiert Teilstring anhand der Position.<code>str_detect(text, muster)</code>: Prüft auf Muster (TRUE/FALSE).<code>str_match(text, muster)</code>: Extrahiert Treffer des Musters.<code>str_extract(text, muster)</code>: Extrahiert den ersten Treffer des Musters.<code>str_replace(text, muster, ersatz)</code>: Ersetzt den ersten Treffer.<code>str_replace_all(text, muster, ersatz)</code>: Ersetzt alle Treffer.</div> <div>Länge und Bereinigung:</div> <div><ul style="list-style-type: none"><code>str_length(text)</code>: Gibt die Länge des Textes zurück.<code>str_pad(text, länge, seite, zeichen)</code>: Füllt Text auf.<code>str_trunc(text, länge)</code>: Kürzt Text auf eine bestimmte Länge.<code>str_squish(text)</code>: Entfernt überflüssige Leerzeichen.</div>
<div>4. Reguläre Ausdrücke (Regex)</div> <div>Definition: Muster zur Textsuche, -prüfung und -manipulation. Syntax:</div> <div><ul style="list-style-type: none">Literale Zeichen: Die meisten Zeichen matchen sich selbst.Sonderzeichen: Müssen mit Backslash \ escaped werden (z.B. <code>1 \+ 1 = 2</code>).Zeichenklassen:<ul style="list-style-type: none"><code>[ae]</code> matcht a oder e.<code>`</code> matcht eine Ziffer.<code>[^x]</code> matcht alles außer x.Kurzformen:<ul style="list-style-type: none"><code>\d</code>: Ziffer (äquivalent zu <code>`</code>)<code>\w</code>: Wortzeichen (Buchstaben, Zahlen, Unterstrich).<code>\s</code>: Whitespace (Leerzeichen, Tab, Zeilenumbruch).Quantifizierer:<ul style="list-style-type: none"><code>?</code>: Optional (0 oder 1 mal).<code>*</code>: 0 oder mehr mal.<code>+</code>: 1 oder mehr mal.<code>{n}</code>: Genau n mal.<code>{n,m}</code>: Zwischen n und m mal.Anker:<ul style="list-style-type: none"><code>^</code>: Anfang des Strings.<code>\$</code>: Ende des Strings.<code>\b</code>: Wortgrenze.Alternation: <code> </code> (entspricht "oder").Gruppierung: (muster) (erfasst Treffer in Gruppen).Lookaround:<ul style="list-style-type: none"><code>(?=muster)</code>: Positive Lookahead.<code>(?!muster)</code>: Negative Lookahead.<code>(?<=muster)</code>: Positive Lookbehind.<code>(?<!muster)</code>: Negative Lookbehind.</div>	<div>5. Sonderzeichen in Regex</div> <div>Umlaute: Können direkt verwendet werden, Encoding beachten! Häufige Probleme</div> <div><ul style="list-style-type: none">- Umlaute können als <code>`<fc>`</code> oder <code>`<e4>`</code> angezeigt werden- Dateien mit Umlauten können Probleme beim Einlesen verursachen</div> <div>Tipps zur Vermeidung</div> <div><ul style="list-style-type: none">- Immer UTF-8 verwenden für Deutsch- RStudio-Einstellungen auf UTF-8 setzen- Überprüfen der Datentypen beim Import- Pakete wie <code>`readr`</code> für robusteren Datenimport nutzen- Umlaute in Variablennamen und Dateinamen vermeiden</div> <div>Verarbeitung von Umlauten: Encoding-Einstellungen</div> <div><ul style="list-style-type: none">- <code>`Sys.setlocale()`</code> für globale Encoding-Einstellungen- Für Windows: Encoding bei Dateimporten explizit angeben</div> <div><pre>```r read.table("Quelle.csv", sep="," , fileEncoding="UTF-8") ```</pre></div> <div>Umlauten ersetzen</div> <div><ul style="list-style-type: none">- Mit <code>`stringr`</code> können Umlaute gleichzeitig ersetzt werden:</div> <div><pre>```r library(stringr) str_replace_all('üäö', c('ü' = 'ue', 'ä' = 'ae', 'ö' = 'oe')) ```</pre></div> <div>Andere Sonderzeichen:</div> <div><ul style="list-style-type: none">Mit Backslash escapen (z.B. <code>\.</code> für einen Punkt).In Zeichenklassen verwenden (z.B. <code>[.,!?</code> für Satzzeichen).</div> <div>regexr.com: Online-Editor zum Testen von Regex. rdocumentation.org: stringr-Dokumentation</div>