

**Corpus der Entscheidungen
des
Bundesfinanzhofs
(CE-BFH)**

CODEBOOK

Version 2023-10-15



DOI: [10.5281/zenodo.7691841](https://doi.org/10.5281/zenodo.7691841)

Titel	Corpus der Entscheidungen des Bundesfinanzhofs
Abkürzung	CE-BFH
Autor	Seán Fobbe
Version	2023-10-15
Download	https://doi.org/10.5281/zenodo.7691841
Lizenz	CC0 1.0 Universal

Zitiervorschlag

Seán Fobbe (2023). Corpus der Entscheidungen des Bundesfinanzhofs (CE-BFH). Version 2023-10-15. Zenodo. DOI: [10.5281/zenodo.7691841](https://doi.org/10.5281/zenodo.7691841).

Digital Object Identifier (DOI): Concept DOI und Version DOI

Soweit nicht anders angegeben ist die DOI immer eine »Version DOI« und bezieht sich nur auf eine bestimmte Version des Datensatzes. Sie verweist daher nur auf Version 2023-10-15. Für das Gesamtkonzept dieses Datensatzes steht eine »Concept DOI« zur Verfügung, die auf der Zenodo-Seite jeder Version unter »Cite all versions?« zu finden ist. Sie lautet [10.5281/zenodo.7691840](https://doi.org/10.5281/zenodo.7691840). Die »Concept DOI« verlinkt immer die aktuellste Version.

Urheberrecht

Der Datensatz und dieses Dokument sind unter einer **Creative Commons CC0 1.0 Universal (CC0 1.0) Public Domain Dedication Lizenz** veröffentlicht. Ich stelle den Datensatz und das Codebook vollständig gemeinfrei und verzichte weltweit auf alle damit verbundenen Urheberrechte, einschließlich aller ähnlichen Rechte, soweit dies gesetzlich möglich ist.

Sie können die Werke kopieren, modifizieren, verteilen und aufführen ohne um Erlaubnis bitten zu müssen, selbst für kommerzielle Zwecke. Patente und Markenschutzrechte bleiben von CC0 unberührt. CC0 hat auch keine Auswirkungen auf etwaige Datenschutz- oder Persönlichkeitsrechte. Jegliche Haftung für die Benutzung dieses Werkes ist ausgeschlossen, bis zu dem maximalen Umfang in dem dies gesetzlich möglich ist.

Wenn Sie diese Werke nutzen oder zitieren sollten Sie nicht den Eindruck erwecken, der Autor unterstütze ihre Nutzung.

Dies ist nur eine unverbindliche deutsche Zusammenfassung der Lizenz, den vollständigen und rechtsverbindlichen Lizenztext finden Sie hier: <https://creativecommons.org/publicdomain/zero/1.0/legalcode>

Disclaimer

Dieser Datensatz ist eine private wissenschaftliche Initiative und steht in keiner Verbindung zu Behörden, Gerichten oder anderen amtlichen Stellen der Bundesrepublik Deutschland.

Inhaltsverzeichnis

1	Einführung	4
2	Nutzung	5
2.1	CSV-Dateien	5
2.2	TXT-Dateien	5
3	Konstruktion	6
3.1	Beschreibung des Datensatzes	6
3.2	Datenquellen	6
3.3	Sammlung der Daten	6
3.4	Source Code und Compilation Report	6
3.5	Grenzen des Datensatzes	8
3.6	Urheberrechtsfreiheit von Rohdaten und Datensatz	8
3.7	Metadaten	8
3.7.1	Allgemein	8
3.7.2	Schema für die Dateinamen	8
3.7.3	Beispiel eines Dateinamens	8
3.8	Qualitätsprüfung	9
3.9	Grafische Darstellung	9
4	Varianten und Zielgruppen	10
5	Variablen	12
5.1	Datenstruktur	12
5.2	Allgemeine Hinweise	13
5.3	ID-Variablen	14
5.4	Text-Variablen	14
5.5	Thematische Variablen	16
5.6	Temporale Variablen	17
5.7	Meta-Variablen	18
6	Registerzeichen	20
7	Linguistische Kennzahlen	21
7.1	Erläuterung der Kennzahlen und Diagramme	21
7.2	Werte der Kennzahlen	21
7.3	Verteilung Zeichen	22
7.4	Verteilung Tokens	22
7.5	Verteilung Typen	23
7.6	Verteilung Sätze	23
8	Inhalt des Korpus	24
8.1	Zusammenfassung	24
8.2	Nach Typ der Entscheidung	24
8.3	Nach Spruchkörper (Aktenzeichen)	25
8.4	Nach Registerzeichen	26
8.5	Nach Entscheidungsjahr	27
8.6	Nach Eingangsjahr (ISO)	28

8.7 Nach Normen	30
9 Dateigrößen	33
10 Kryptographische Signaturen	35
10.1 Zwei-Phasen-Signatur	35
10.2 Persönliche GPG-Signatur	35
11 Changelog	36
11.1 Version 2023-10-15	36
12 Parameter für strenge Replikationen	37
Literaturverzeichnis	39

1 Einführung

Der **Bundesfinanzhof (BFH)** ist einer der fünf obersten Gerichtshöfe des Bundes und steht an der Spitze der Finanzgerichtsbarkeit der Bundesrepublik Deutschland (Art. 95 Abs. 1 GG, §§ 2, 10 f. FGO). Der BFH ist die höchste Instanz in Steuer- und Zollsachen und entscheidet über Revisionen gegen Urteile der Finanzgerichte, gegen urteilsgleiche Entscheidungen und Beschwerden gegen andere Entscheidungen der Finanzgerichte (§ 36 FGO).¹ Er wurde mit dem »Gesetz über den Bundesfinanzhof« vom 29. Juni 1950 errichtet und hat seinen Sitz in München (§ 2 FGO).

Im Jahr 2022 sind am Bundesfinanzhof 11 Senate eingerichtet, bestehend aus ca. 60 Richter:innen.² Zudem besteht ein Großer Senat, dem Richter:innen aller Senate angehören. Präsident:in und Vizepräsident:in vertreten das Gericht nach Außen. Entscheidungen der Senate ergehen in einer Besetzung von 5 Richter:innen, Beschlüsse außerhalb der mündlichen Verhandlung können von 3 Richter:innen getroffen werden (§ 10 Abs. 3 FGO). Der Große Senat besteht aus Gerichtspräsident:in und einer Richter:in jedes Senates, in dem der Vorsitz nicht von der Gerichtspräsident:in geführt wird (§ 11 Abs. 5 FGO), d.h. aktuell 11 Mitgliedern. Zu Beginn jeden Jahres bestimmt das Präsidium des Gerichts die thematische Geschäftsverteilung zwischen den Senaten.

Wieso dieser Datensatz? Die quantitative Analyse von juristischen Texten, insbesondere denen des BFH, ist in den deutschen Rechtswissenschaften ein noch junges und kaum bearbeitetes Feld.³ Zu einem nicht unerheblichen Teil liegt dies auch daran, dass die Anzahl an frei nutzbaren Datensätzen außerordentlich gering ist.

Die meisten hochwertigen Datensätze lagern (fast) unerreichbar in kommerziellen Datenbanken und sind wissenschaftlich gar nicht oder nur gegen Entgelt zu nutzen. Frei verfügbare Datenbanken wie *Opinio Iuris*⁴ und *openJur*⁵ verbieten ausdrücklich das maschinelle Auslesen der Rohdaten. Wissenschaftliche Initiativen wie der Juristische Referenzkorpus (JuReKo) sind nach jahrelanger Arbeit hinter verschlossenen Türen verschwunden.

In einem funktionierenden Rechtsstaat muss die Rechtsprechung öffentlich, transparent und nachvollziehbar sein. Im 21. Jahrhundert bedeutet dies auch, dass sie systematischer Überprüfung mittels quantitativen Analysen zugänglich sein muss. Der Erstellung und Aufbereitung des Datensatzes liegen daher die Prinzipien der allgemeinen Verfügbarkeit durch Urheberrechtsfreiheit, strenge Transparenz und vollständige wissenschaftliche Reproduzierbarkeit zugrunde. Die FAIR-Prinzipien (Findable, Accessible, Interoperable and Reusable) für freie wissenschaftliche Daten inspirieren sowohl die Konstruktion, als auch die Art der Publikation.⁶

¹ Der Finanzrechtsweg ist vergleichsweise kurz: in erster Instanz entscheiden die Finanzgerichte (§ 35 FGO) und es steht im Anschluss nur noch die Anrufung des BFH offen (§ 36 FGO).

² <https://www.bundesfinanzhof.de/de/gericht/organisation/>.

³ Besonders positive Ausnahmen finden sich unter: <https://www.quantitative-rechtswissenschaft.de/>
⁴ <https://opinioiuris.de/>

⁵ <https://openjur.de/>

⁶ Wilkinson, M., Dumontier, M., Aalbersberg, I. et al. The FAIR Guiding Principles for Scientific Data Management and Stewardship. *Sci Data* 3, 160018 (2016). <https://doi.org/10.1038/sdata.2016.18>

2 Nutzung

Die Daten sind in offenen, interoperablen und weit verbreiteten Formaten (CSV, TXT, PDF) veröffentlicht. Sie lassen sich grundsätzlich mit allen modernen Programmiersprachen (z.B. Python oder R), sowie mit grafischen Programmen nutzen.

Wichtig: Nicht vorhandene Werte sind sowohl in den Dateinamen als auch in der CSV-Datei mit "NA" codiert.

2.1 CSV-Dateien

Am einfachsten ist es die **CSV-Dateien** einzulesen. CSV⁷ ist ein einfaches und maschinell gut lesbares Tabellen-Format. In diesem Datensatz sind die Werte komma-separiert. Jede Spalte entspricht einer Variable, jede Zeile einer Entscheidung. Die Variablen sind unter Punkt 5 genauer erläutert.

Zum Einlesen empfehle ich für **R** das package **data.table** (via CRAN verfügbar). Dessen Funktion **fread()** ist etwa zehnmal so schnell wie die normale **read.csv()**-Funktion in Base-R. Sie erkennt auch den Datentyp von Variablen sicherer. Ein Beispiel:

```
library(data.table)
dt <- fread("filename.csv")
```

2.2 TXT-Dateien

Die **TXT-Dateien** inklusive Metadaten können zum Beispiel mit **R** und dem package **readtext** (via CRAN verfügbar) eingelesen werden. Ein Vorschlag:

```
library(readtext)
df <- readtext("./*.txt",
               docvarsfrom = "filenames",
               docvarnames = c("gericht",
                               "bfhe",
                               "datum",
                               "spruchkoerper_az",
                               "registerzeichen",
                               "eingangsnummer",
                               "eingangsjahr_az",
                               "zusatz_az",
                               "bfh_id",
                               "kollision"),
               dvsep = "_",
               encoding = "UTF-8")
```

⁷ Das CSV-Format ist in RFC 4180 definiert, siehe <https://tools.ietf.org/html/rfc4180>

3 Konstruktion

3.1 Beschreibung des Datensatzes

Dieser Datensatz ist eine digitale Zusammenstellung von möglichst allen begründeten Entscheidungen, die auf der amtlichen Internetpräsenz des Bundesfinanzhofs (BFH) am jeweiligen Stichtag veröffentlicht waren. Die Stichtage für jede Version entsprechen exakt der Versionsnummer.

Zusätzlich zu den aufbereiteten maschinenlesbaren Formaten (HTML und CSV) sind die PDF-Daten enthalten, damit Analyst:innen gegebenenfalls eine unabhängige Konvertierung vornehmen können. Die PDF-Rohdaten wurden inhaltlich nicht verändert und nur die Dateinamen angepasst, um die Lesbarkeit für Mensch und Maschine zu verbessern.

Speziell an Praktiker:innen richten sich die PDF-Sammlungen aller in der amtlichen Sammlung abgedruckten Entscheidungen (V-Entscheidungen).

3.2 Datenquellen

Datenquelle	Fundstelle
Primäre Datenquelle	https://www.bundesfinanzhof.de
Source Code	https://doi.org/10.5281/zenodo.7691843
Registerzeichen	https://doi.org/10.5281/zenodo.4569564

Die Tabelle der Registerzeichen und der ihnen zugeordneten Verfahrensarten stammt aus dem folgenden Datensatz: “Seán Fobbe (2021). Aktenzeichen der Bundesrepublik Deutschland (AZ-BRD). Version 1.0.1. Zenodo. DOI: 10.5281/zenodo.4569564.”

3.3 Sammlung der Daten

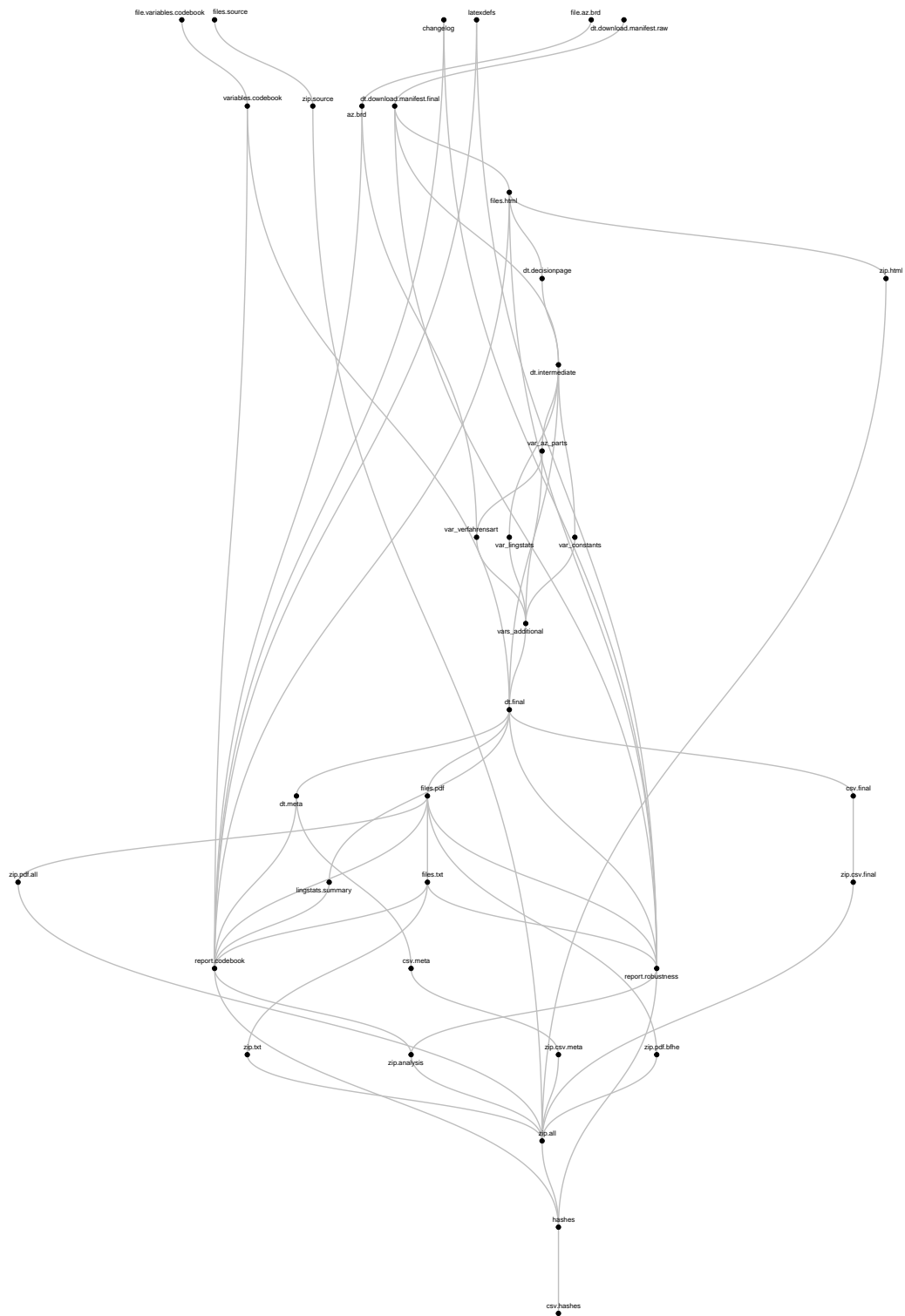
Die Daten wurden unter Beachtung des Robot Exclusion Standard (RES) gesammelt. Der Abruf geschieht ausschließlich über TLS-verschlüsselte Verbindungen. Die Entscheidungen sind laut dem Gericht anonymisiert, aber ungekürzt.

3.4 Source Code und Compilation Report

Der gesamte Source Code — sowohl für die Erstellung des Datensatzes, als auch für dieses Codebook — ist öffentlich einsehbar und dauerhaft erreichbar im wissenschaftlichen Archiv des CERN unter dieser Adresse hinterlegt: <https://doi.org/10.5281/zenodo.7691843>

Mit jeder Kompilierung des vollständigen Datensatzes wird auch ein umfangreicher **Compilation Report** in einem attraktiv designten PDF-Format erstellt (ähnlich diesem Codebook). Der Compilation Report enthält den kommentierten Source Code für die Daten-Pipeline, dokumentiert relevante Rechenergebnisse, gibt sekundengenaue Zeitstempel an und ist mit einem klickbaren Inhaltsverzeichnis versehen. Er ist zusammen mit dem Source Code hinterlegt. Wenn Sie sich für Details der Herstellung interessieren, lesen Sie diesen bitte zuerst.

CE-BFH | Version 2023-10-15 | Vollständiger Prozess der Datensatz-Kompilierung



Fobbe | DOI: 10.5281/zenodo.7691841

Abbildung 1: Der vollständige Prozess der Datensatz-Kompilierung.

3.5 Grenzen des Datensatzes

Nutzer:innen sollten folgende wichtige Grenzen beachten:

- Der Datensatz enthält nur das, was das Gericht auch tatsächlich veröffentlicht, nämlich begründete Entscheidungen (*publication bias*).
- Es kann aufgrund technischer Grenzen bzw. Fehler sein, dass manche — im Grunde verfügbare — Entscheidungen nicht oder nicht korrekt abgerufen werden (*automation bias*).
- Es werden HTML- und PDF-Dateien abgerufen (*file type bias*). Der Text der Entscheidungen in der CSV-Datei stammt aus den HTML-Dateien. Die PDF-Dateien sind beigefügt, falls bei der Extraktion Fehler auftreten sollten, sind in der Regel aber für den Gebrauch in der traditionellen Rechtswissenschaft und -praxis gedacht.
- Erst ab dem 1. Januar 2010 sind begründete Entscheidungen des Bundesfinanzhofs einigermaßen vollständig veröffentlicht (*temporal bias*). Die Frequenztabellen geben hierzu genauer Auskunft.

3.6 Urheberrechtsfreiheit von Rohdaten und Datensatz

An den Entscheidungstexten und amtlichen Leitsätzen besteht gem. § 5 Abs. 1 UrhG kein Urheberrecht, da sie amtliche Werke sind. § 5 UrhG ist auf amtliche Datenbanken analog anzuwenden (BGH, Beschluss vom 28.09.2006, I ZR 261/03, »Sächsischer Ausschreibungsdienst«).

Alle eigenen Beiträge (z.B. durch Zusammenstellung und Anpassung der Metadaten) und damit den gesamten Datensatz stelle ich gemäß einer *CC0 1.0 Universal Public Domain Lizenz* vollständig urheberrechtsfrei.

3.7 Metadaten

3.7.1 Allgemein

Die Metadaten in den Dateinamen sind größtenteils unverändert von den jeweiligen Datenbankeinträgen aus der amtlichen Datenbank des Bundesfinanzhofs entnommen. Berechnet und hinzugefügt wurden durch den Autor des Datensatzes eine Reihe weitere Variablen, sowie in den Dateinamen der PDF/TXT-Dateien Unter- und Trennstriche, um die Maschinenlesbarkeit zu erleichtern. Der volle Satz an Metadaten ist nur in den CSV-Dateien enthalten. Alle hinzugefügten Metadaten sind vollständig maschinenlesbar dokumentiert. Sie sind entweder im Source Code enthalten, mit dem Source Code zusammen dokumentiert oder über dauerhaft stabile Identifikatoren (z.B. DOI) zitiert.

Die Dateinamen der PDF- und TXT-Dateien enthalten Gerichtsname, die Bezeichnung als V- oder NV-entscheidung, Datum, das offizielle Aktenzeichen, einen Zusatz zum Aktenzeichen und die vom BFH in der Datenbank genutzte einzigartige ID.

3.7.2 Schema für die Dateinamen

```
[gericht]_[bfhe]_[datum]_[spruchkoerper_az]_  
[registerzeichen]_[eingangsnummer]_[eingangsjahr_az]_[bfh_id]
```

3.7.3 Beispiel eines Dateinamens

```
BFH_V_2023-07-11_X_R_17_22_STRE202310190.pdf
```

3.8 Qualitätsprüfung

Die Inhalte der Variablen wurden strikt validiert. Die möglichen Werte der jeweiligen Variablen wurden zudem durch Frequenztabellen und Visualisierungen auf ihre Plausibilität geprüft. Insgesamt werden zusammen mit jeder Kompilierung über 30 automatisierte Tests zur Qualitätsprüfung durchgeführt. Alle Ergebnisse der Qualitätsprüfungen sind aggregiert im Robustness Checks Report, im Compilation Report und einzeln im Archiv »ANALYSE« zusammen mit dem Datensatz veröffentlicht.

3.9 Grafische Darstellung

Die Robenfarbe der Richter:innen des Bundesfinanzhofs ist »karmesinrot«. Der Hex-Wert hierfür ist vermutlich #7e0731. Das ist besonders bei der Erstellung thematisch passender Diagrammen hilfreich. Alle im Compilation Report und diesem Codebook präsentierten Diagramme sind in diesem karmesinrot gehalten.

4 Varianten und Zielgruppen

Dieser Datensatz ist in verschiedenen Varianten verfügbar, die sich an unterschiedliche Zielgruppen richten. Zielgruppe sind nicht nur quantitativ forschende Rechtswissenschaftler:innen, sondern auch traditionell arbeitende Jurist:innen. Idealerweise müssen quantitative Methoden ohnehin immer durch qualitative Interpretation, Theoriebildung und kritische Auseinandersetzung verstärkt werden (*mixed methods approach*).

Lehrende werden von den vorbereiteten Tabellen und Diagrammen besonders profitieren, die bei der Erläuterung der Charakteristika der Daten hilfreich sein können und Zeit im universitären Alltag sparen. Alle Tabellen und Diagramme liegen auch als separate Dateien vor, um sie einfach z.B. in Präsentations-Folien oder Handreichungen zu integrieren.

Variante	Zielgruppe und Beschreibung
PDF	Traditionelle juristische Forschung. Die PDF-Dokumente wie sie vom BFH auf der amtlichen Webseite bereitgestellt werden, jedoch verbessert durch semantisch hochwertige Dateinamen, die der leichteren Auffindbarkeit von Entscheidungen dienen. Die Dateinamen sind so konzipiert, dass sie auch für die traditionelle qualitative juristische Arbeit einen erheblichen Mehrwert bieten. Im Vergleich zu den CSV-Dateien enthalten die Dateinamen nur einen reduzierten Umfang an Metadaten, um Kompatibilitätsprobleme zu vermeiden und die Lesbarkeit zu verbessern. Neben dem vollen Datensatz ist für Praktiker:innen auch eine Variante aufbereitet, die nur <i>V-Entscheidungen</i> der amtlichen Sammlung enthalten.
CSV_Datensatz	Legal Tech/Quantitative Forschung. Diese CSV-Datei ist die für statistische Analysen empfohlene Variante des Datensatzes. Sie enthält den Volltext aller Entscheidungen, sowie alle in diesem Codebook beschriebenen Metadaten. Jede Spalte entspricht einer Variable, jede Zeile einer Entscheidung.
CSV_Metadaten	Legal Tech/Quantitative Forschung. Wie die vorige CSV-Variante, nur ohne die Entscheidungstexte. Sinnvoll für Analyst:innen, die sich nur für die Metadaten interessieren und Speicherplatz sparen wollen. Jede Spalte entspricht einer Variable, jede Zeile einer Entscheidung.
TXT	Subsidiär für alle Zielgruppen. Diese Variante enthält die vollständigen, aus den PDF-Dateien extrahierten Entscheidungstexte, aber nur einen reduzierten Umfang an Metadaten, der dem der PDF-Dateien entspricht. Die TXT-Dateien sind optisch an das Layout der PDF-Dateien angelehnt. Geeignet für qualitativ arbeitende Forscher:innen, die nur wenig Speicherplatz oder eine langsame Internetverbindung zur Verfügung haben oder für quantitativ arbeitende Forscher:innen, die beim Einlesen der CSV-Dateien Probleme haben.

Variante	Zielgruppe und Beschreibung
HTML	Subsidiär für Legal Tech/Quantitative Forschung. Diese Variante enthält die ursprünglichen HTML-Dateien, wie sie auf der Webseite des BFH präsentiert werden. Nur sinnvoll, falls Probleme bei der Nutzung der CSV-Dateien auftreten.
ANALYSE	Alle Lehrenden und Forschenden. Dieses Archiv enthält alle während dem Kompilierungs- und Prüfprozess erstellten Tabellen (CSV) und Diagramme (PDF, PNG) im Original. Sie sind inhaltsgleich mit den in diesem Codebook verwendeten Tabellen und Diagrammen. Das PDF-Format eignet sich besonders für die Verwendung in gedruckten Publikationen, das PNG-Format besonders für die Darstellung im Internet. Analyst:innen mit fortgeschrittenen Kenntnissen in R können auch auf den Source Code zurückgreifen. Empfohlen für Nutzer:innen die einzelne Inhalte aus dem Codebook für andere Zwecke (z.B. Präsentationen, eigene Publikationen) weiterverwenden möchten.

5 Variablen

5.1 Datenstruktur

```
## Classes 'data.table' and 'data.frame':  10310 obs. of  33 variables:
## $ aktenzeichen      : chr  "VII S 34/09" "IV R 43/07" "VII B 118/09" "VII
  B 165/09" ...
## $ bfh_id           : chr  "STRE201050294" "STRE201050190" "STRE201050120
  " "STRE201050238" ...
## $ doc_id           : chr  "BFH_NV_2010-01-04_VII_S_34_9_STRE201050294" "
  BFH_NV_2010-01-05_IV_R_43_7_STRE201050190" "BFH_NV_2010-01-07_VII_B_118_9_
  STRE201050120" "BFH_NV_2010-01-07_VII_B_165_9_STRE201050238" ...
## $ ecll             : chr  NA NA NA NA ...
## $ gericht          : chr  "BFH" "BFH" "BFH" "BFH" ...
## $ text_leitsatz    : chr  "1. NV: Ob bei Verpachtung eines Betriebes der
  Pächter oder der Verpächter Milcherzeuger ist, bedarf einer umfas"| __
  truncated__ "NV: Der Gesellschafter einer zweigliedrigen GbR ist nach § 48
  Abs. 1 Nr. 3 FGO klagebefugt und damit notwendig "| __truncated__ "NV:
  Restschuldbefreiung erlangt der Insolvenzschuldner nicht mit dem Ablauf der
  sog. Wohlverhaltensphase, sonde"| __truncated__ "1. NV: Die für die
  Erstattung nach Art. 901 Abs. 2 ZKDVO erforderlichen Nachweise können nicht
  allein mit den i"| __truncated__ ...
## $ url_html         : chr  "https://www.bundesfinanzhof.de/de/
  entscheidung/entscheidungen-online/detail/STRE201050294/" "https://www.
  bundesfinanzhof.de/de/entscheidung/entscheidungen-online/detail/STRE201050190
  /" "https://www.bundesfinanzhof.de/de/entscheidung/entscheidungen-online/
  detail/STRE201050120/" "https://www.bundesfinanzhof.de/de/entscheidung/
  entscheidungen-online/detail/STRE201050238/" ...
## $ url_pdf          : chr  "https://www.bundesfinanzhof.de/de/
  entscheidung/entscheidungen-online/detail/pdf/STRE201050294?type=1646225765"
  "https://www.bundesfinanzhof.de/de/entscheidung/entscheidungen-online/detail/
  pdf/STRE201050190?type=1646225765" "https://www.bundesfinanzhof.de/de/
  entscheidung/entscheidungen-online/detail/pdf/STRE201050120?type=1646225765"
  "https://www.bundesfinanzhof.de/de/entscheidung/entscheidungen-online/detail/
  pdf/STRE201050238?type=1646225765" ...
## $ zeichen          : num  14086 8666 4027 6056 11204 ...
## $ tokens           : int  2118 1482 619 1027 1899 2112 1523 3447 721 897
  ...
## $ typen            : int  649 494 290 394 583 764 540 967 302 369 ...
## $ saetze           : int  47 99 29 70 87 75 78 164 25 49 ...
## $ adv              : logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
## $ bfhe             : chr  "NV" "NV" "NV" "NV" ...
## $ normen           : chr  "FGO § 133a|FGO § 118 Abs 2|FGO § 126 Abs 3 S
  1 Nr 2|EGV 1788/2003 Art 5 Buchst c|GG Art 103 Abs 1|FGO § 96 Abs 1 S 1" "FGO
  § 48 Abs 1 Nr 3|FGO § 60 Abs 3" "Ins0 § 294 Abs 3|Ins0 § 300|Ins0 § 301" "ZK
  Art 238|ZK Art 239|ZKDV Art 901 Abs 2|ZKDV Art 902 Abs 1|ZKDV Art 904 Buchst
  a|EWGV 2454/93 Art 901 Abs 2|EW"| __truncated__ ...
## $ pkh              : logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
## $ registerzeichen  : chr  "S" "R" "B" "B" ...
## $ spruchkoerper_az : chr  "VII" "IV" "VII" "VII" ...
## $ spruchkoerper_db : chr  "VII. Senat" "IV. Senat" "VII. Senat" "VII.
  Senat" ...
## $ titel            : chr  "Tatsachenfeststellung und rechtliche Wü
  rdigung bei Ermittlung des Milcherzeugers - Keine Bindung an rechtliche "| __
```

```

truncated_ "Notwendige Beiladung des aus einer zweigliedrigen GbR
ausgeschiedenen Gesellschafters bei Streit über Einkunftsart der GbR" "Kein
insolvenzrechtliches Aufrechnungsverbot zwischen Aufhebung des
Insolvenzverfahrens und Erteilung der Restschuldbefreiung" "Zoll: Erstattung
der Einfuhrabgaben bei Wiederausfuhr der Ware" ...
## $ verfahrensart      : chr  "Sonstige Verfahren" "Revision" "Beschwerden"
"Beschwerden" ...
## $ vorinstanz        : chr  "vorgehend BFH , 25. Mai 2009, Az: VII R
28/08" "vorgehend Finanzgericht Berlin-Brandenburg , 13. Juni 2007, Az: 15 K
3202/04 B" "vorgehend Finanzgericht des Landes Sachsen-Anhalt , 17. März
2009, Az: 2 K 1682/08" "vorgehend FG Hamburg, 26. Mai 2009, Az: 4 K 58/07"
...
## $ datum            : Date, format: "2010-01-04" "2010-01-05" ...
## $ entscheidungsjahr : int   2010 2010 2010 2010 2010 2010 2010 2010 2010
2010 ...
## $ eingangsjahr_az   : int   9 7 9 9 9 7 8 7 9 9 ...
## $ eingangsjahr_iso  : num   2009 2007 2009 2009 2009 ...
## $ eingangsnummer    : int   34 43 118 165 99 34 159 24 113 110 ...
## $ veroeffentlichung : Date, format: "2010-05-19" "2010-04-14" ...
## $ veroeffentlichungsjahr: int   2010 2010 2010 2010 2010 2010 2010 2010 2010
2010 ...
## $ doi_concept       : chr   "10.5281/zenodo.7691840" "10.5281/zenodo
.7691840" "10.5281/zenodo.7691840" "10.5281/zenodo.7691840" ...
## $ doi_version       : chr   "10.5281/zenodo.7691841" "10.5281/zenodo
.7691841" "10.5281/zenodo.7691841" "10.5281/zenodo.7691841" ...
## $ lizenz            : chr   "Creative Commons Zero 1.0 Universal" "
Creative Commons Zero 1.0 Universal" "Creative Commons Zero 1.0 Universal" "
Creative Commons Zero 1.0 Universal" ...
## $ version           : chr   "2023-10-15" "2023-10-15" "2023-10-15"
"2023-10-15" ...
## - attr(*, ".internal.selfref")=<externalptr>

```

5.2 Allgemeine Hinweise

- **Doppelte Codierung der Spruchkörper** — Für viele Urteile sind die Spruchkörper doppelt enthalten, einmal aus der Datenbank (Variable »spruchkoerper_db«), einmal durch das Aktenzeichen (Variable »spruchkoerper_az«).
- **Fehlende Werte** sind immer mit »NA« codiert.
- **Strings** können grundsätzlich alle in UTF-8 definierten Zeichen (insbesondere Buchstaben, Zahlen und Sonderzeichen) enthalten.
- Die **Reihenfolge** der Variablen entspricht der im CSV-Datensatz. Der Datensatz wird automatisiert darauf getestet, ob alle Variablen im Datensatz auch in diesem Codebook dokumentiert sind.

5.3 ID-Variablen

ID-Variablen stellen verschiedene Identifikatoren für die Entscheidung zur Verfügung, beispielsweise Aktenzeichen oder ECLI.

Variable	Type	Description
aktenzeichen	String	Das amtliche Aktenzeichen im Format [senatsnummer] [registerzeichen] [eingangsnummer] / [eingangsjahr] [ggf. zusatz_az]. Quelle: BFH-Datenbank.
bfh_id	String	Die vom BFH vergebene ID der Entscheidung in der amtlichen Datenbank. Die ID ist einzigartig. Quelle: BFH-Datenbank.
doc_id	String	(Nur CSV) Dateiname der PDF- und TXT-Dateien, ohne Dateierweiterung. Quelle: Kompositum verschiedener Variablen des Datensatzes.
ecli	String	(Nur CSV-Datei) Der European Case Law Identifier (ECLI) der Entscheidung. Die ECLI ist einzigartig. Die ECLI ist vor allem dann hilfreich, wenn dieser Datensatz mit anderen Datensätzen zusammengeführt und Duplikate vermieden werden sollen. Nicht für alle Entscheidungen vorhanden. Quelle: BFH-Datenbank, aus den HTML-Dateien extrahiert
gericht	String	Name des Gerichts. Es ist nur der Wert »BFH« vergeben. Dies ist der ECLI-Code für »Bundesfinanzhof«. Diese Variable dient vor allem zur einfachen und transparenten Verbindung der Daten mit anderen Datensätzen. Quelle: Autor des Datensatzes.

5.4 Text-Variablen

Text-Variablen enthalten den Volltext der Entscheidung, Teilstücke davon (z.B. Leitsätze), den Umfang des Volltextes (Zeichen, Tokens, Typen, Sätze) und dessen Quelle (URLs zu Volltexten).

Variable	Type	Description
text	String	(Nur CSV) Volltext der Entscheidung. Achtung: wenige Entscheidungen (ca. 20) haben keinen Text, weil es Parallelentscheidungen sind. Quelle: BFH-Datenbank, aus den HTML-Dateien extrahiert

(continued)

Variable	Type	Description
text_leitsatz	String	(Nur CSV) Text der Leitsätze der Entscheidung. Quelle: BFH-Datenbank, aus den HTML-Dateien extrahiert.
url_html	String	(Nur CSV) Link zum Volltext der Entscheidung als HTML in der amtlichen Datenbank des Gerichts. Quelle: BFH-Datenbank.
url_pdf	String	(Nur CSV) Link zum Volltext der Entscheidung als PDF in der amtlichen Datenbank des Gerichts. Quelle: BFH-Datenbank.
zeichen	Integer	(Nur CSV) Die Anzahl Zeichen eines Dokumentes. Quelle: Mit R berechnet.
tokens	Integer	(Nur CSV) Die Anzahl Tokens (beliebige Zeichenfolge getrennt durch whitespace) eines Dokumentes. Diese Zahl kann je nach Tokenizer und verwendeten Einstellungen erheblich schwanken. Für diese Berechnung wurde eine reine Tokenisierung ohne Entfernung von Inhalten durchgeführt. Benutzen Sie diesen Wert eher als Anhaltspunkt für die Größenordnung denn als exakte Aussage und führen sie ggf. mit ihrer eigenen Software eine Kontroll-Rechnung durch. Quelle: Mit R berechnet
typen	Integer	(Nur CSV) Die Anzahl <i>einzigartiger</i> Tokens (beliebige Zeichenfolge getrennt durch whitespace) eines Dokumentes. Diese Zahl kann je nach Tokenizer und verwendeten Einstellungen erheblich schwanken. Für diese Berechnung wurde eine reine Tokenisierung und Typenzählung ohne Entfernung von Inhalten durchgeführt. Benutzen Sie diesen Wert eher als Anhaltspunkt für die Größenordnung denn als exakte Aussage und führen sie ggf. mit ihrer eigenen Software eine Kontroll-Rechnung durch. Quelle: mit R berechnet.

(continued)

Variable	Type	Description
saetze	Integer	(Nur CSV) Die Anzahl Sätze. Die Definition entspricht in etwa dem üblichen Verständnis eines Satzes. Die Regeln für die Bestimmung von Satzanfang und Satzende sind im Detail allerdings sehr komplex und in »Unicode Standard: Annex No 29« beschrieben. Diese Zahl kann je nach Software und verwendeten Einstellungen erheblich schwanken. Für diese Berechnung wurde eine reine Zählung ohne Entfernung von Inhalten durchgeführt. Benutzen Sie diesen Wert eher als Anhaltspunkt für die Größenordnung denn als exakte Aussage und führen sie ggf. mit ihrer eigenen Software eine Kontroll-Rechnung durch. Quelle: Mit R berechnet

5.5 Thematische Variablen

Thematische Variablen geben Auskunft über eine grobe thematische Zuordnung der Entscheidung, beispielsweise zu Registerzeichen, Verfahrensart, Normen, Vorinstanz.

Variable	Type	Description
adv	Logical	Ob es sich um eine Entscheidung zur Aufhebung der Vollziehung (AdV) handelt. Entweder TRUE oder FALSE. Quelle: REGEX-Suche nach "AdV" im Aktenzeichen.
bfhe	String	(Nur CSV) Ob die Entscheidung in der amtlichen Sammlung veröffentlicht wird (»V«) oder nicht (»NV«). Der BFH spricht hier auch von V-Entscheidungen und NV-Entscheidungen. Quelle: BFH-Datenbank.
normen	String	(Nur CSV) Die rechtlichen Normen, die von der Entscheidung betroffen sind. Normen beginnen jeweils mit dem Gesetzesnamen, gefolgt von der genauen Fundstelle. Mehrere Normen sind durch einen vertikale Balken (» «) getrennt. Quelle: BFH-Datenbank.
pkh	Logical	(Nur CSV) Ob es sich um eine Entscheidung zur Prozesskostenhilfe (PKH) handelt. Entweder TRUE oder FALSE. Quelle: REGEX-Suche nach "PKH" im Aktenzeichen.

(continued)

Variable	Type	Description
registerzeichen	String	Das amtliche Registerzeichen. Eine Erläuterung der Abkürzungen findet sich im Abschnitt 6. Quelle: Mit REGEX aus Variable "aktenzeichen" extrahiert.
spruchkoerper_az	String	Der im Aktenzeichen angegebene Spruchkörper. Die Senate sind mit römischen Ziffern nummeriert. Der Große Senat ist mit »GrS« gekennzeichnet. Quelle: Mit REGEX aus Aktenzeichen extrahiert.
spruchkoerper_db	String	(Nur CSV) Der Spruchkörper, wie er in der amtlichen Datenbank des Gerichts eingetragen ist. Quelle: BFH-Datenbank
titel	String	(Nur CSV) Der Titel der Entscheidung. Enthält eine kurze thematische und rechtliche Einordnung. Quelle: BFH-Datenbank
verfahrensart	String	(Nur CSV) Die Verfahrensart, auf die das Registerzeichen hinweist. Siehe auch Abschnitt 6. Quelle: Abgleich von Registerzeichen mit AZ-BRD-Datensatz.
vorinstanz	String	(Nur CSV) Die Vorinstanz des Verfahrens. Quelle: BFH-Datenbank.

5.6 Temporale Variablen

Temporale Variablen bieten Informationen zu wichtigen Zeitpunkten, wie Verkündung der Entscheidung, Veröffentlichung der Entscheidung oder Eingang des Verfahrens.

Variable	Type	Description
datum	Date	Datum der Entscheidung im Format YYYY-MM-DD (ISO-8601). Quelle: BFH-Datenbank.
entscheidungsjahr	Integer	Jahr der Entscheidung im Format YYYY (ISO-8601). Quelle: Berechnet aus Variable "datum".
eingangsjahr_az	Integer	Eingangsjahr laut Aktenzeichen. Das Jahr in dem das Verfahren beim Gericht anhängig wurde. Das Format ist eine zweistellige Jahreszahl (YY). Quelle: Mit REGEX aus Variable "aktenzeichen" extrahiert.

(continued)

Variable	Type	Description
eingangsjahr_iso	Integer	(Nur CSV) Eingangsjahr im Format YYYY-MM-DD (ISO-8601). Quelle: Aus Variable “eingangsjahr_az” berechnet.
eingangsnummer	Integer	Eingangsnummer. Verfahren des gleichen Eingangsjahres erhalten vom Gericht eine fortlaufende Nummer (Ordinalzahl) in der Reihenfolge ihres Eingangs. Quelle: Mit REGEX aus Variable “aktenzeichen” extrahiert.
veroeffentlichung	Date	(Nur CSV) Das Datum der Veröffentlichung der Entscheidung im Format YYYY-MM-DD (ISO-8601). Quelle: BFH-Datenbank.
veroeffentlichungsjahr	Integer	(Nur CSV) Das Jahr der Veröffentlichung der Entscheidung im Format YYYY (ISO-8601). Quelle: Aus Variable “veroeffentlichung” berechnet.

5.7 Meta-Variablen

Meta-Variablen beziehen sich auf den Datensatz selbst. Sie dokumentieren Versionsnummer, verschiedene DOIs und die Lizenz des Datensatzes. Streng genommen sind sie innerhalb des Datensatzes Konstanten (weil der Inhalt immer gleich ist) und nur im Vergleich zwischen Datensätzen echte Variablen.

Variable	Type	Description
doi_concept	String	(Nur CSV) Der Digital Object Identifier (DOI) des Gesamtkonzeptes des Datensatzes. Dieser ist langzeit-stabil (persistent). Über diese DOI kann via www.doi.org immer die aktuellste Version des Datensatzes abgerufen werden. Prinzip F1 der FAIR-Data Prinzipien («data are assigned globally unique and persistent identifiers») empfiehlt die Dokumentation jeder Messung mit einem persistenten Identifikator. Selbst wenn die CSV-Dateien ohne Kontext weitergegeben werden kann ihre Herkunft so immer zweifelsfrei und maschinenlesbar bestimmt werden. Quelle: Vom Autor hinzugefügt.

(continued)

Variable	Type	Description
doi_version	String	(Nur CSV) Der Digital Object Identifier (DOI) der konkreten Version des Datensatzes. Dieser ist langzeit-stabil (persistent). Über diese DOI kann via www.doi.org immer diese konkrete Version des Datensatzes abgerufen werden. Prinzip F1 der FAIR-Data Prinzipien («data are assigned globally unique and persistent identifiers») empfiehlt die Dokumentation jeder Messung mit einem persistenten Identifikator. Selbst wenn die CSV-Dateien ohne Kontext weitergegeben werden kann ihre Herkunft so immer zweifelsfrei und maschinenlesbar bestimmt werden. Quelle: Vom Autor hinzugefügt.
lizenz	String	Die Lizenz für den Gesamtdatensatz. In diesem Datensatz immer »Creative Commons Zero 1.0 Universal«. Quelle: Vom Autor hinzugefügt.
version	Date	(Nur CSV) Die Versionsnummer des Datensatzes im Format YYYY-MM-DD (Langform nach ISO-8601). Die Versionsnummer entspricht immer dem Datum an dem der Datensatz erstellt und die Daten von der Webseite des Gerichts abgerufen wurden. Quelle: Vom Autor hinzugefügt.

6 Registerzeichen

Die Tabelle der Registerzeichen und der ihnen zugeordneten Verfahrensarten stammt aus dem folgenden Datensatz: »Seán Fobbe (2021). Aktenzeichen der Bundesrepublik Deutschland (AZ-BRD). Version 1.0.1. Zenodo. DOI: 10.5281/zenodo.4569564.«

Die im Datensatz enthaltenen Registerzeichen wurden jeweils um die runden Klammern bereinigt, um Probleme bei der Nutzung unter Windows zu vermeiden.

Die Bedeutung des Registerzeichens »ER-S« ist mir nicht klar, aber möglicherweise ist es eine Kombination aus den Registerzeichen »E«, »R« und »S«. Ich arbeite an der Aufklärung.

Registerzeichen	Verfahrensart
AR	Allgemeines Register: Vorverfahren oder sonstige Verfahrensarten
B	Beschwerden
E	Erinnerung zu Kosten und Streitwert, Erinnerung in Kostenfestsetzungsverfahren
GrS	Großer Senat
K	Entschädigungsklagen wegen überlanger Verfahrensdauer
PKH	Prozesskostenhilfe
R	Revision
S	Sonstige Verfahren
ER-S	Unklar

7 Linguistische Kennzahlen

7.1 Erläuterung der Kennzahlen und Diagramme

Zur besseren Einschätzung des inhaltlichen Umfangs des Korpus dokumentiere ich an dieser Stelle die Verteilung der Werte für einige klassische linguistische Kennzahlen:

Kennzahl	Definition
Zeichen	Zeichen entsprechen grob den <i>Graphemen</i> , den kleinsten funktionalen Einheiten in einem Schriftsystem. Beispiel: das Wort »RichterIn« besteht aus 9 Zeichen.
Tokens	Eine beliebige Zeichenfolge, getrennt durch whitespace-Zeichen, d.h. ein Token entspricht in der Regel einem »Wort«, kann aber auch Zahlen, Sonderzeichen oder sinnlose Zeichenfolgen enthalten, weil es rein syntaktisch berechnet wird.
Typen	Einzigartige Tokens. Beispiel: wenn das Token »Finanzrecht« zehnmal in einer Entscheidung vorhanden ist, wird es als ein Typ gezählt.
Sätze	Entsprechen in etwa dem üblichen Verständnis eines Satzes. Die Regeln für die Bestimmung von Satzanfang und Satzende sind im Detail aber sehr komplex und in »Unicode Standard: Annex No 29« beschrieben.

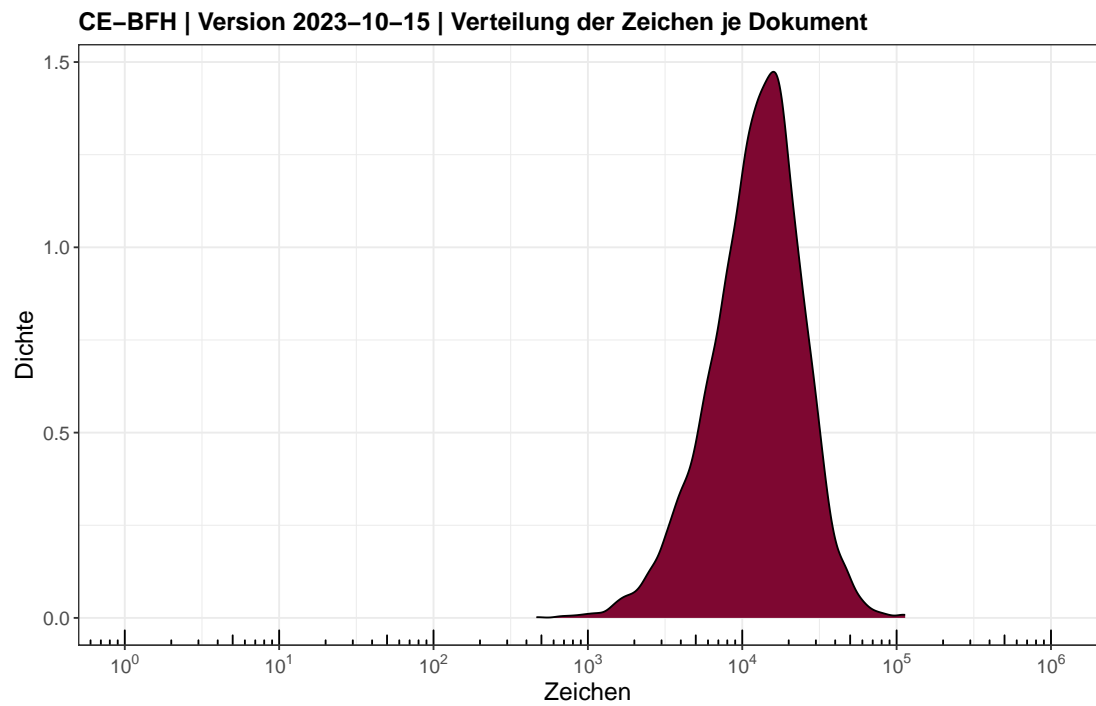
Es handelt sich bei den Diagrammen jeweils um »Density Charts«, die sich besonders dafür eignen die Schwerpunkte von Variablen mit stark schwankenden numerischen Werten zu visualisieren. Die Interpretation ist denkbar einfach: je höher die Kurve, desto dichter sind in diesem Bereich die Werte der Variable. Der Wert der y-Achse kann außer Acht gelassen werden, wichtig sind nur die relativen Flächenverhältnisse und die x-Achse.

Vorsicht bei der Interpretation: Die x-Achse ist logarithmisch skaliert, d.h. in 10er-Potenzen und damit nicht-linear. Die kleinen Achsen-Markierungen zwischen den Schritten der Exponenten sind eine visuelle Hilfestellung um diese nicht-Linearität zu verstehen.

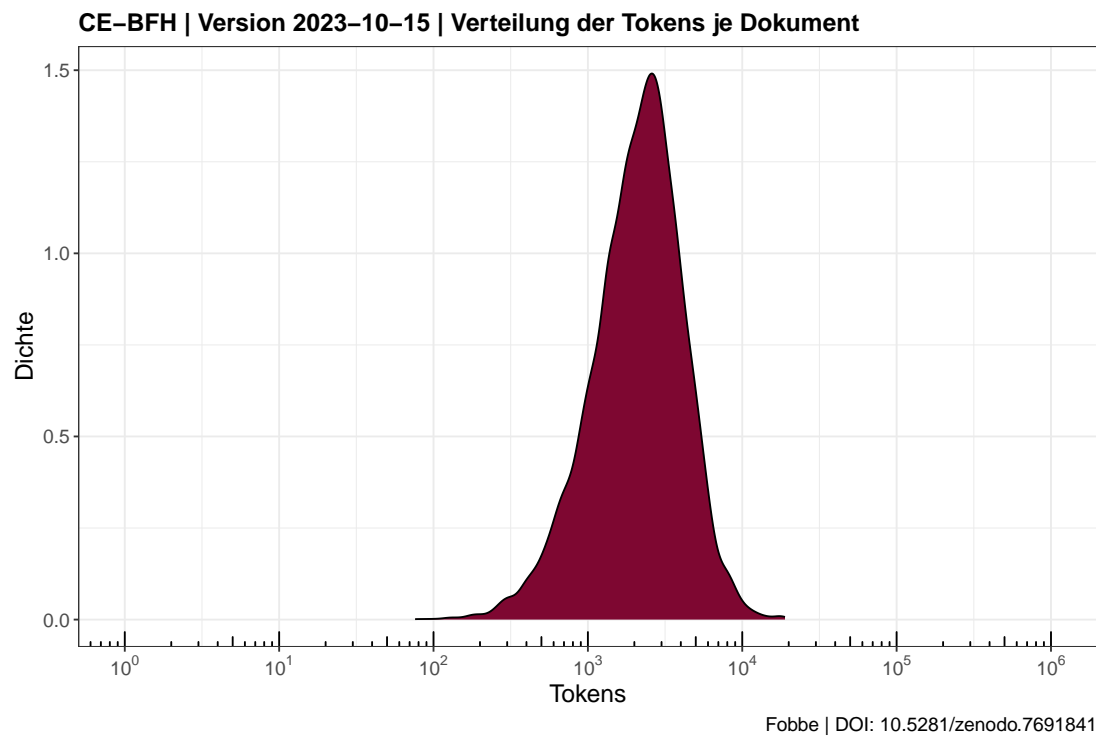
7.2 Werte der Kennzahlen

Variable	Summe	Min	Quart1	Median	Mittel	Quart3	Max
zeichen	154,700,545	0	8,037	12,948	15,004.90	19,220.5	113,616
tokens	26,230,089	0	1,366	2,201	2,544.14	3,268.0	18,918
typen	269,062	0	484	682	725.87	904.0	3,251
saetze	1,125,917	0	61	95	109.21	140.0	790

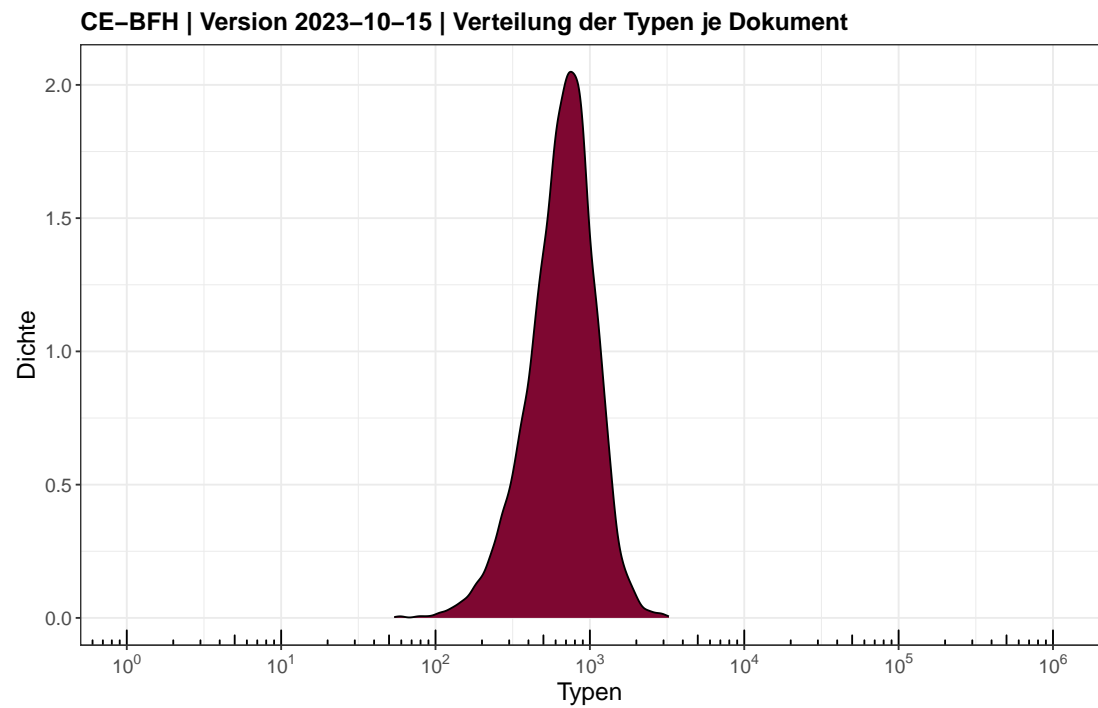
7.3 Verteilung Zeichen



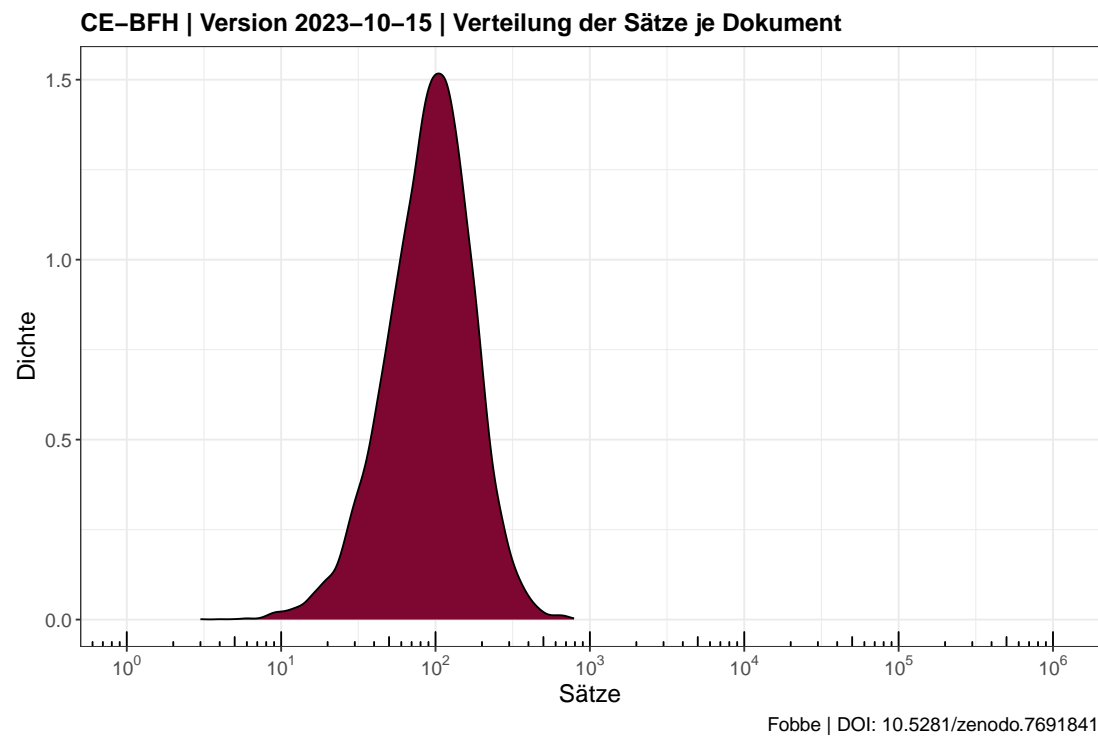
7.4 Verteilung Tokens



7.5 Verteilung Typen



7.6 Verteilung Sätze

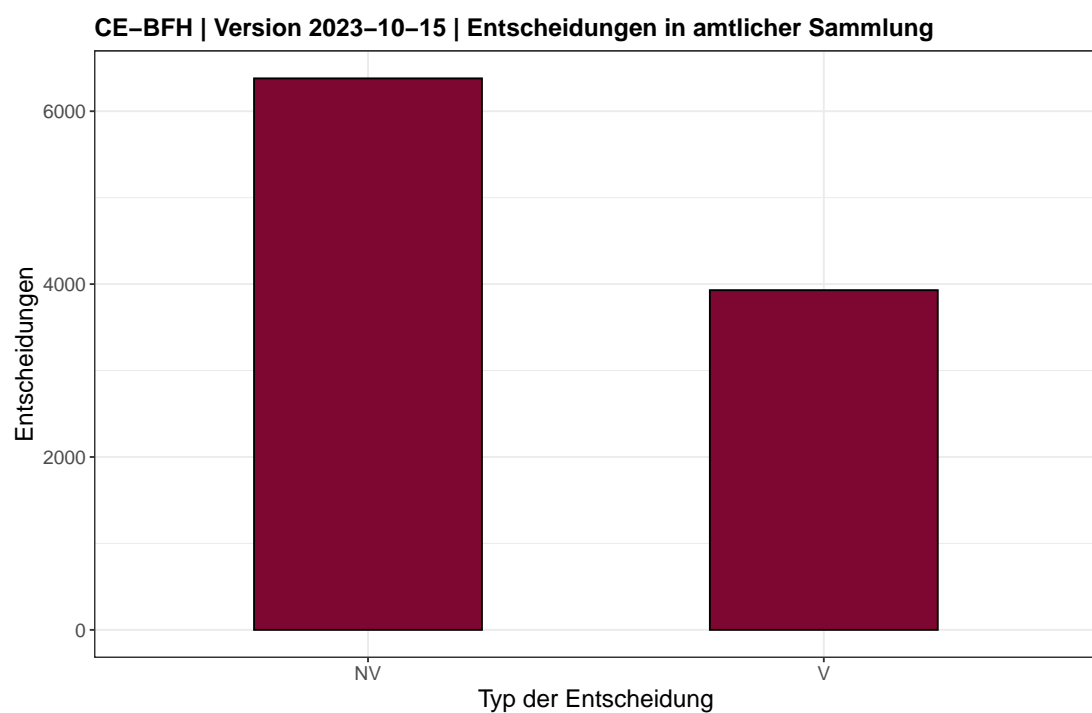


8 Inhalt des Korpus

8.1 Zusammenfassung

Variable	Anzahl	Min	Quart1	Median	Mean	Quart3	Max
entscheidungsjahr	14	2010	2012	2014	2014.89	2018	2023
eingangsjahr_iso	22	2001	2010	2013	2013.50	2016	2023
eingangsnummer	266	1	17	36	50.64	66	281

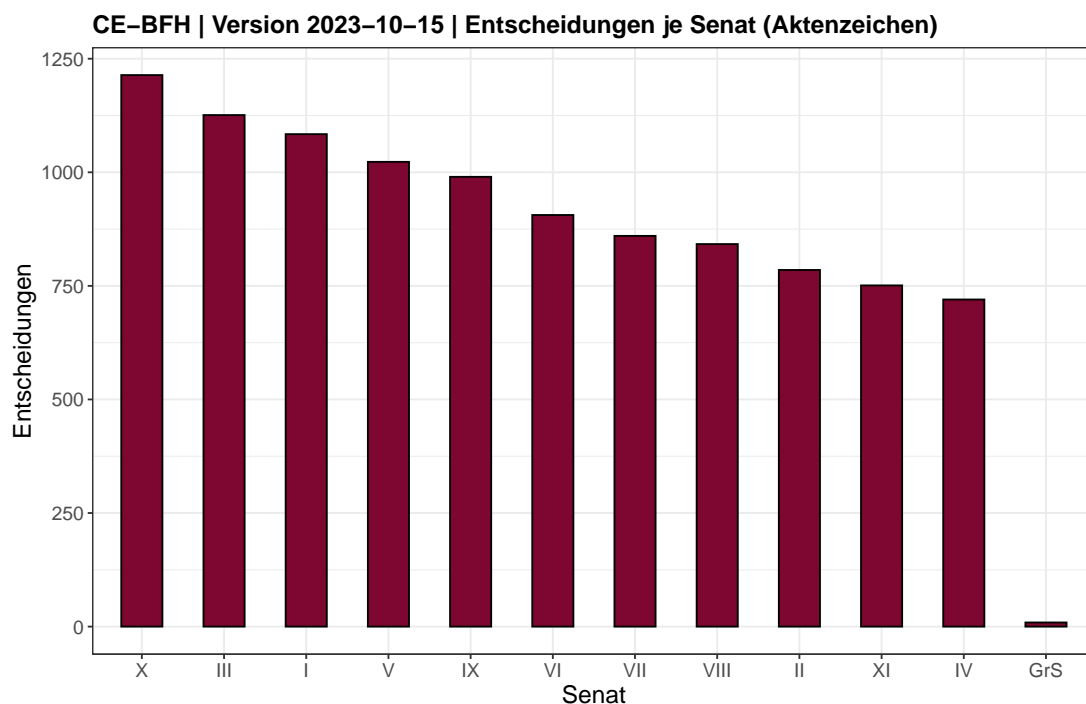
8.2 Nach Typ der Entscheidung



Fobbe | DOI: 10.5281/zenodo.7691841

Typ	Entscheidungen	% Gesamt	% Kumulativ
NV	6380	61.88	61.88
V	3930	38.12	100.00
Total	10310	100.00	100.00

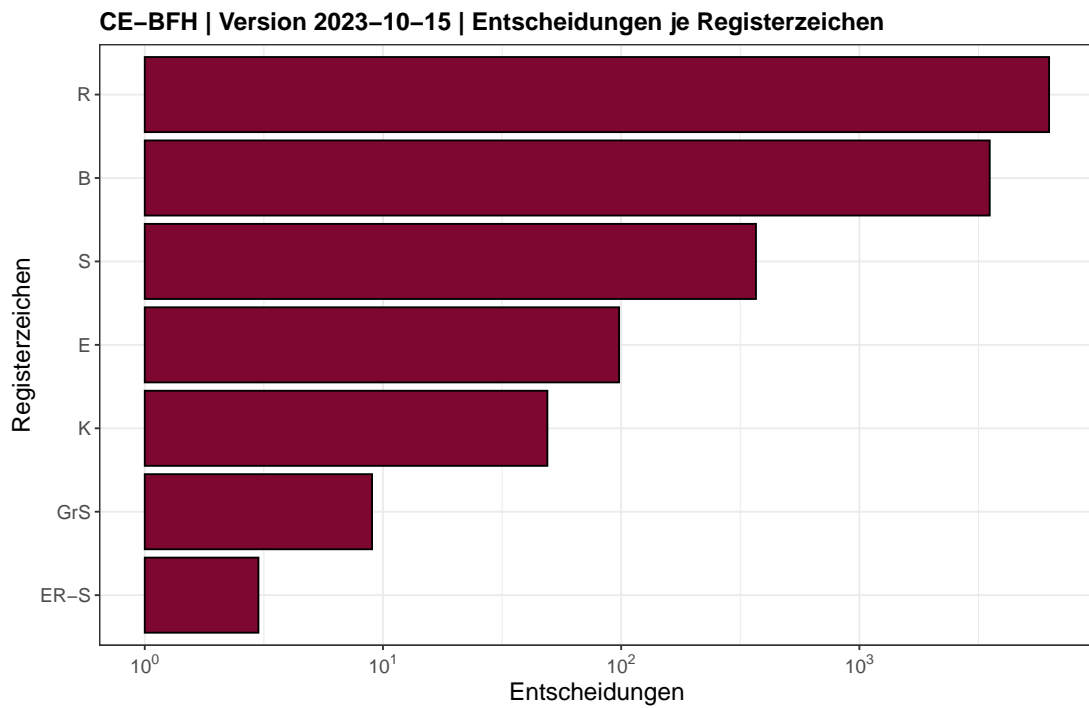
8.3 Nach Spruchkörper (Aktenzeichen)



Fobbe | DOI: 10.5281/zenodo.7691841

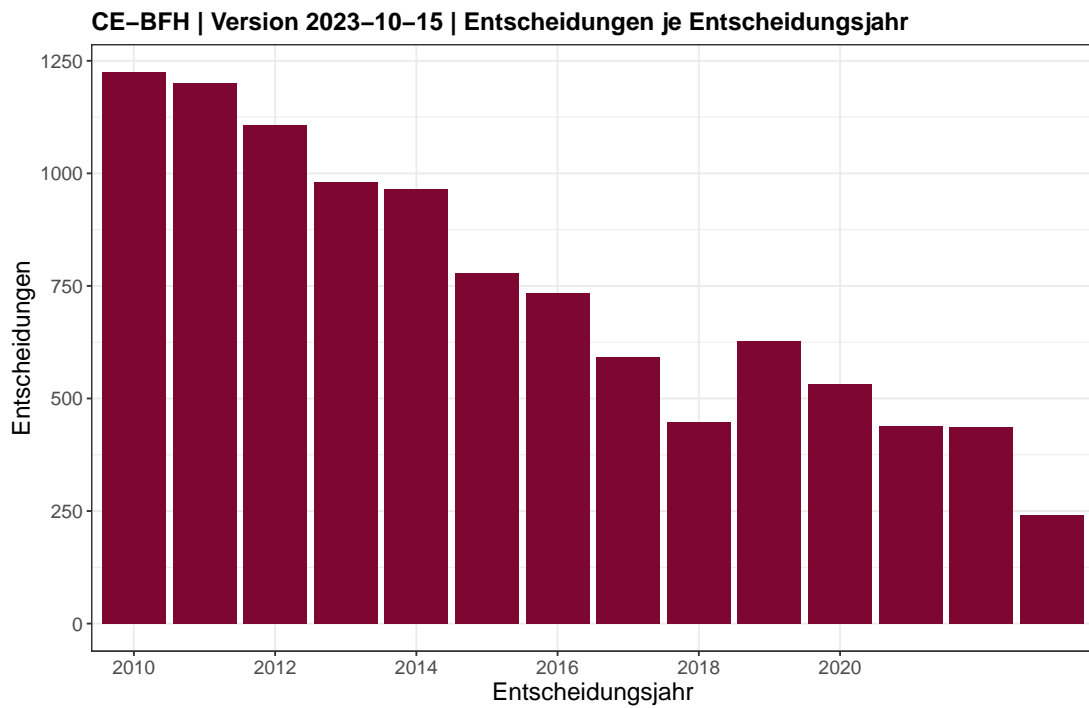
Senat	Entscheidungen	% Gesamt	% Kumulativ
GrS	9	0.09	0.09
I	1084	10.51	10.60
II	785	7.61	18.22
III	1126	10.92	29.14
IV	720	6.98	36.12
IX	990	9.60	45.72
V	1023	9.92	55.65
VI	906	8.79	64.43
VII	860	8.34	72.77
VIII	842	8.17	80.94
X	1214	11.77	92.72
XI	751	7.28	100.00
Total	10310	100.00	100.00

8.4 Nach Registerzeichen



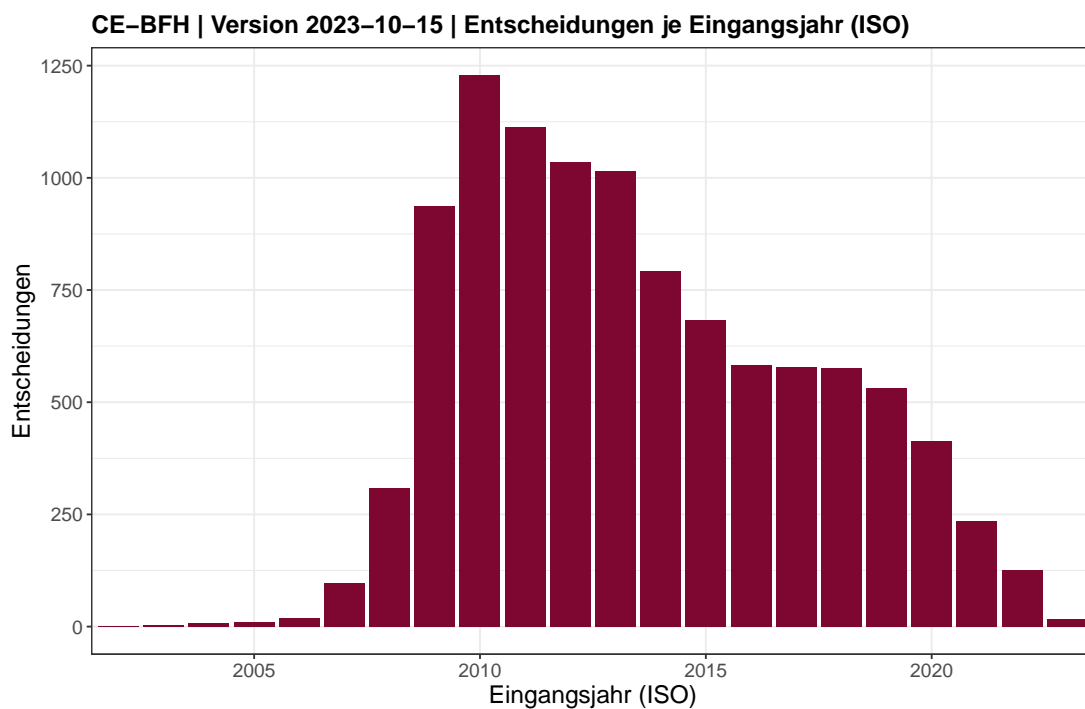
Registerzeichen	Entscheidungen	% Gesamt	% Kumulativ
B	3525	34.19	34.19
E	98	0.95	35.14
ER-S	3	0.03	35.17
GrS	9	0.09	35.26
K	49	0.48	35.73
R	6258	60.70	96.43
S	368	3.57	100.00
Total	10310	100.00	100.00

8.5 Nach Entscheidungsjahr



Entscheidungsjahr	Entscheidungen	% Gesamt	% Kumulativ
2010	1225	11.88	11.88
2011	1201	11.65	23.53
2012	1107	10.74	34.27
2013	981	9.52	43.78
2014	966	9.37	53.15
2015	778	7.55	60.70
2016	735	7.13	67.83
2017	591	5.73	73.56
2018	448	4.35	77.90
2019	628	6.09	84.00
2020	532	5.16	89.16
2021	439	4.26	93.41
2022	437	4.24	97.65
2023	242	2.35	100.00
Total	10310	100.00	100.00

8.6 Nach Eingangsjahr (ISO)

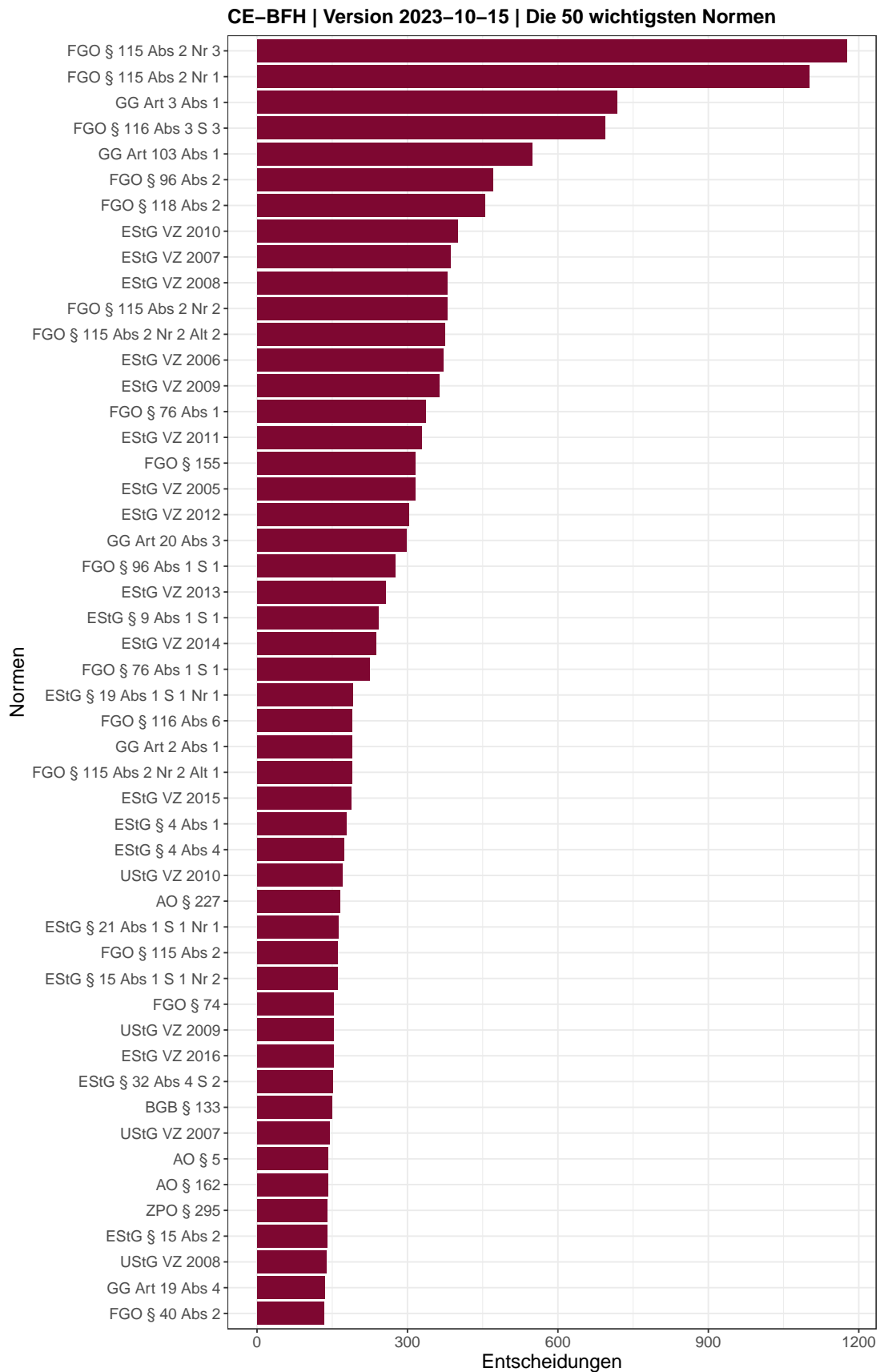


Eingangsjahr	Entscheidungen	% Gesamt	% Kumulativ
2001	1	0.01	0.01
2003	3	0.03	0.04
2004	7	0.07	0.11
2005	10	0.10	0.20
2006	19	0.18	0.39
2007	98	0.95	1.34
2008	308	2.99	4.33
2009	938	9.10	13.42
2010	1229	11.92	25.34
2011	1113	10.80	36.14
2012	1035	10.04	46.18
2013	1014	9.84	56.01
2014	793	7.69	63.71
2015	684	6.63	70.34
2016	583	5.65	75.99
2017	579	5.62	81.61
2018	575	5.58	87.19

(continued)

Eingangsjahr	Entscheidungen	% Gesamt	% Kumulativ
2019	531	5.15	92.34
2020	414	4.02	96.35
2021	234	2.27	98.62
2022	126	1.22	99.84
2023	16	0.16	100.00
Total	10310	100.00	100.00

8.7 Nach Normen



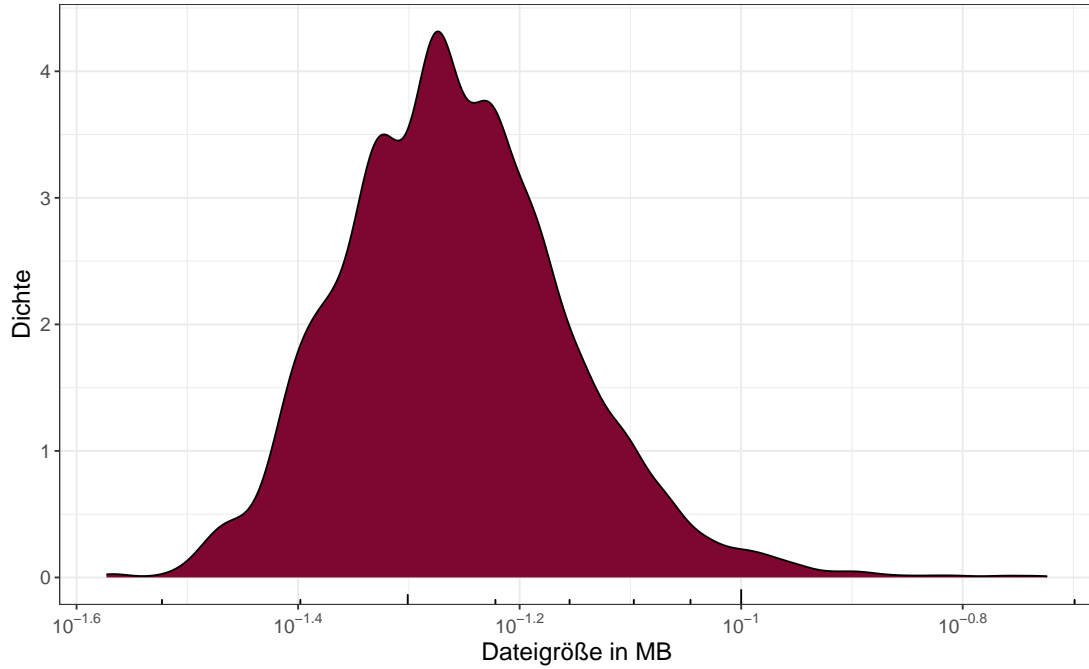
Normen	Entscheidungen	% Gesamt
FGO § 115 Abs 2 Nr 3	1176	1.69
FGO § 115 Abs 2 Nr 1	1101	1.58
GG Art 3 Abs 1	717	1.03
FGO § 116 Abs 3 S 3	694	1.00
GG Art 103 Abs 1	548	0.79
FGO § 96 Abs 2	470	0.67
FGO § 118 Abs 2	454	0.65
EStG VZ 2010	400	0.57
EStG VZ 2007	386	0.55
EStG VZ 2008	380	0.55
FGO § 115 Abs 2 Nr 2	379	0.54
FGO § 115 Abs 2 Nr 2 Alt 2	374	0.54
EStG VZ 2006	371	0.53
EStG VZ 2009	363	0.52
FGO § 76 Abs 1	336	0.48
EStG VZ 2011	328	0.47
EStG VZ 2005	316	0.45
FGO § 155	316	0.45
EStG VZ 2012	302	0.43
GG Art 20 Abs 3	298	0.43
FGO § 96 Abs 1 S 1	275	0.39
EStG VZ 2013	257	0.37
EStG § 9 Abs 1 S 1	242	0.35
EStG VZ 2014	237	0.34
FGO § 76 Abs 1 S 1	224	0.32
EStG § 19 Abs 1 S 1 Nr 1	191	0.27
FGO § 116 Abs 6	190	0.27
FGO § 115 Abs 2 Nr 2 Alt 1	189	0.27
GG Art 2 Abs 1	189	0.27
EStG VZ 2015	188	0.27
EStG § 4 Abs 1	178	0.26
EStG § 4 Abs 4	174	0.25
UStG VZ 2010	170	0.24
AO § 227	165	0.24

(continued)

Normen	Entscheidungen	% Gesamt
EStG § 21 Abs 1 S 1 Nr 1	162	0.23
EStG § 15 Abs 1 S 1 Nr 2	161	0.23
FGO § 115 Abs 2	161	0.23
FGO § 74	153	0.22
EStG VZ 2016	152	0.22
UStG VZ 2009	152	0.22
EStG § 32 Abs 4 S 2	151	0.22
BGB § 133	149	0.21
UStG VZ 2007	145	0.21
AO § 5	142	0.20
AO § 162	141	0.20
ZPO § 295	140	0.20
EStG § 15 Abs 2	139	0.20
UStG VZ 2008	138	0.20
GG Art 19 Abs 4	135	0.19
FGO § 40 Abs 2	134	0.19

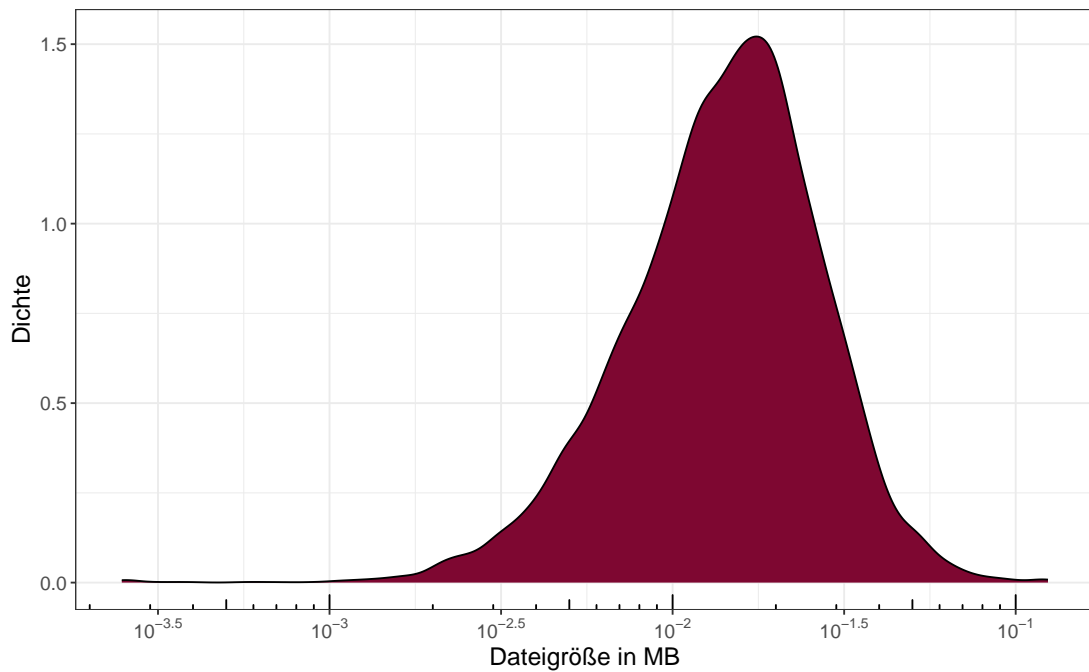
9 Dateigrößen

CE-BFH | Version 2023-10-15 | Verteilung der Dateigrößen (PDF)



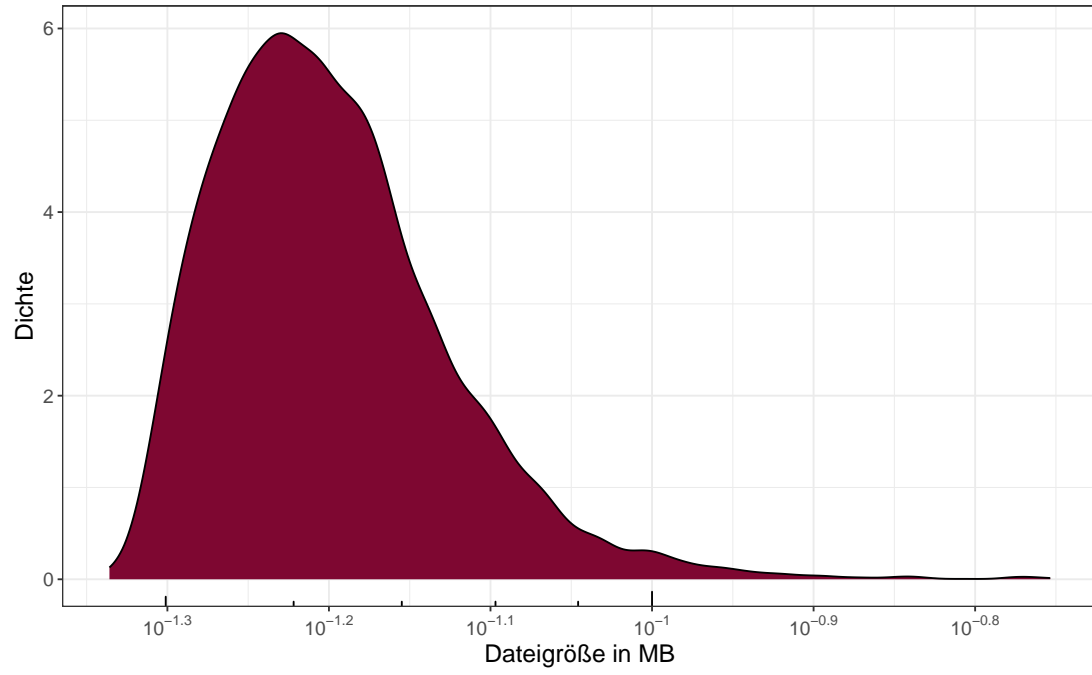
Fobbe | DOI: 10.5281/zenodo.7691841

CE-BFH | Version 2023-10-15 | Verteilung der Dateigrößen (TXT)



Fobbe | DOI: 10.5281/zenodo.7691841

CE-BFH | Version 2023-10-15 | Verteilung der Dateigrößen (HTML)



Fobbe | DOI: 10.5281/zenodo.7691841

10 Kryptographische Signaturen

10.1 Zwei-Phasen-Signatur

Die Integrität und Echtheit der einzelnen Archive des Datensatzes sind durch eine Zwei-Phasen-Signatur sichergestellt.

In **Phase I** werden während der Kompilierung für jedes ZIP-Archiv, das Codebook und die Robustness Checks Hash-Werte in zwei verschiedenen Verfahren (SHA2-256 und SHA3-512) berechnet und in einer CSV-Datei dokumentiert.

In **Phase II** werden diese CSV-Datei und der Compilation Report mit meinem persönlichen geheimen GPG-Schlüssel signiert. Dieses Verfahren stellt sicher, dass die Kompilierung von jedermann durchgeführt werden kann, insbesondere im Rahmen von Replikationen, die persönliche Gewähr für Ergebnisse aber dennoch vorhanden bleibt.

10.2 Persönliche GPG-Signatur

Die während der Kompilierung des Datensatzes erstellte CSV-Datei mit den Hash-Prüfsummen und der Compilation Report sind mit meiner persönlichen GPG-Signatur versehen. Der mit dieser Version korrespondierende Public Key ist sowohl mit dem Datensatz als auch mit dem Source Code hinterlegt. Er hat folgende Kenndaten:

Name: Sean Fobbe (fobbe-data@posteo.de)

Fingerabdruck: FE6F B888 F0E5 656C 1D25 3B9A 50C4 1384 F44A 4E42

11 Changelog

11.1 Version 2023-10-15

- Erstveröffentlichung

12 Parameter für strenge Replikationen

```
## [1] "OpenSSL 3.0.2 15 Mar 2022 (Library: OpenSSL 3.0.2 15 Mar 2022)"
```

```
## R version 4.2.2 (2022-10-31)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 22.04.2 LTS
##
## Matrix products: default
## BLAS: /usr/lib/x86_64-linux-gnu/openblas-pthread/libblas.so.3
## LAPACK: /usr/lib/x86_64-linux-gnu/openblas-pthread/libopenblas-p-r0.3.20.so
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8 LC_NUMERIC=C
## [3] LC_TIME=en_US.UTF-8 LC_COLLATE=en_US.UTF-8
## [5] LC_MONETARY=en_US.UTF-8 LC_MESSAGES=en_US.UTF-8
## [7] LC_PAPER=en_US.UTF-8 LC_NAME=C
## [9] LC_ADDRESS=C LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] stats graphics grDevices utils datasets methods base
##
## other attached packages:
## [1] future.apply_1.10.0 future_1.32.0 quanteda_3.2.4
## [4] readtext_0.81 data.table_1.14.8 scales_1.2.1
## [7] ggraph_2.1.0 ggplot2_3.4.1 pdftools_3.3.3
## [10] kableExtra_1.3.4 knitr_1.42 rvest_1.0.3
## [13] httr_1.4.5 mgsub_1.7.3 zip_2.2.2
## [16] fs_1.6.1 testthat_3.1.7 RcppTOML_0.2.2
## [19] tarchetypes_0.7.5 targets_0.14.3
##
## loaded via a namespace (and not attached):
## [1] webshot_0.5.4 rprojroot_2.0.3 future.callr_0.8.1
## [4] tools_4.2.2 backports_1.4.1 utf8_1.2.3
## [7] R6_2.5.1 colorspace_2.1-0 withr_2.5.0
## [10] tidyselct_1.2.0 gridExtra_2.3 processx_3.8.0
## [13] compiler_4.2.2 cli_3.6.0 xml2_1.3.3
## [16] desc_1.4.2 labeling_0.4.2 stringfish_0.15.7
## [19] callr_3.7.3 askpass_1.1 systemfonts_1.0.4
## [22] stringr_1.5.0 digest_0.6.31 rmarkdown_2.20
## [25] svglite_2.1.1 pkgconfig_2.0.3 htmltools_0.5.4
## [28] parallelly_1.34.0 fastmap_1.1.1 rlang_1.0.6
## [31] rstudioapi_0.14 farver_2.1.1 generics_0.1.3
## [34] RApiSerialize_0.1.2 dplyr_1.1.0 magrittr_2.0.3
## [37] Matrix_1.5-1 waldo_0.4.0 Rcpp_1.0.10
## [40] munsell_0.5.0 fansi_1.0.4 viridis_0.6.2
## [43] lifecycle_1.0.3 furrr_0.3.1 stringi_1.7.12
## [46] yaml_2.3.7 MASS_7.3-58.1 brio_1.1.3
## [49] grid_4.2.2 parallel_4.2.2 listenv_0.9.0
## [52] ggrepel_0.9.3 lattice_0.20-45 graphlayouts_0.8.4
## [55] ps_1.7.2 pillar_1.8.1 igraph_1.4.1
## [58] base64url_1.4 codetools_0.2-18 stopwords_2.3
```

```
## [61] pkgload_1.3.2      fastmatch_1.1-3    glue_1.6.2
## [64] evaluate_0.20      qpdf_1.3.0         RcppParallel_5.1.7
## [67] vctrs_0.5.2        tweenr_2.0.2       gtable_0.3.1
## [70] purrr_1.0.1        polyclip_1.10-4    tidyr_1.3.0
## [73] qs_0.25.5          xfun_0.37          ggforce_0.4.1
## [76] tidygraph_1.2.3    viridisLite_0.4.1  tibble_3.2.0
## [79] tinytex_0.44       globals_0.16.2
```

Literaturverzeichnis

- Allaire, JJ, Yihui Xie, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron Atkins, Hadley Wickham, Joe Cheng, Winston Chang, and Richard Iannone. 2023. *Rmarkdown: Dynamic Documents for R*.
- Bengtsson, Henrik. 2021. “A Unifying Framework for Parallel and Distributed Processing in R Using Futures.” *The R Journal* 13 (2): 208–27. <https://doi.org/10.32614/RJ-2021-048>.
- . 2022. *Future.apply: Apply Function to Elements in Parallel Using Futures*.
- . 2023. *Future: Unified Parallel and Distributed Processing in R for Everyone*.
- Benoit, Kenneth, Kohei Watanabe, Haiyan Wang, Paul Nulty, Adam Obeng, Stefan Müller, and Akitaka Matsuo. 2018. “Quanteda: An R Package for the Quantitative Analysis of Textual Data.” *Journal of Open Source Software* 3 (30): 774. <https://doi.org/10.21105/joss.00774>.
- Benoit, Kenneth, Kohei Watanabe, Haiyan Wang, Paul Nulty, Adam Obeng, Stefan Müller, Akitaka Matsuo, and William Lowe. 2022. *Quanteda: Quantitative Analysis of Textual Data*. <https://quanteda.io>.
- Csardi, Gabor, and Tamas Nepusz. 2006. “The Igraph Software Package for Complex Network Research.” *InterJournal Complex Systems*: 1695. <https://igraph.org>.
- Csárdi, Gábor, Kuba Podgórski, and Rich Geldreich. 2022. *Zip: Cross-Platform Zip Compression*. <https://github.com/r-lib/zip#readme>.
- Dowle, Matt, and Arun Srinivasan. 2023. *Data.table: Extension of ‘Data.frame’*.
- Eddelbuettel, Dirk. 2023. *RcppTOML: Rcpp Bindings to Parser for “Tom’s Obvious Markup Language”*. <http://dirk.eddelbuettel.com/code/rcpp.toml.html>.
- file., See AUTHORS. 2023. *Igraph: Network Analysis and Visualization*.
- Gagolewski, Marek. 2022. “stringi: Fast and Portable Character String Processing in R.” *Journal of Statistical Software* 103 (2): 1–59. <https://doi.org/10.18637/jss.v103.i02>.
- Gagolewski, Marek, Bartek Tartanus, others; Unicode, Inc., and others. 2023. *Stringi: Fast and Portable Character String Processing Facilities*.
- Landau, William Michael. 2021a. *Tarchetypes: Archetypes for Targets*.
- . 2021b. “The Targets R Package: A Dynamic Make-Like Function-Oriented Pipeline Toolkit for Reproducibility and High-Performance Computing.” *Journal of Open Source Software* 6 (57): 2959. <https://doi.org/10.21105/joss.02959>.
- . 2023a. *Tarchetypes: Archetypes for Targets*.
- . 2023b. *Targets: Dynamic Function-Oriented Make-Like Declarative Pipelines*.
- Ooms, Jeroen. 2023. *Pdftools: Text Extraction, Rendering and Converting of Pdf Documents*.
- Pedersen, Thomas Lin. 2022. *Ggraph: An Implementation of Grammar of Graphics for Graphs and Networks*.
- Ushey, Kevin. 2023. *Renv: Project Environments*. <https://rstudio.github.io/renv/>.
- Wickham, Hadley. 2022. *Rvest: Easily Harvest (Scrape) Web Pages*.

- Xie, Yihui. 2014. “Knitr: A Comprehensive Tool for Reproducible Research in R.” In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. Chapman; Hall/CRC.
- . 2015. *Dynamic Documents with R and Knitr*. 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC. <https://yihui.org/knitr/>.
- . 2023. *Knitr: A General-Purpose Package for Dynamic Report Generation in R*. <https://yihui.org/knitr/>.
- Xie, Yihui, J. J. Allaire, and Garrett Grolemond. 2018. *R Markdown: The Definitive Guide*. Boca Raton, Florida: Chapman; Hall/CRC. <https://bookdown.org/yihui/rmarkdown>.
- Xie, Yihui, Christophe Dervieux, and Emily Riederer. 2020. *R Markdown Cookbook*. Boca Raton, Florida: Chapman; Hall/CRC. <https://bookdown.org/yihui/rmarkdown-cookbook>.
- Zhu, Hao. 2021. *KableExtra: Construct Complex Table with Kable and Pipe Syntax*.