

**Corpus der
Amtlichen Entscheidungssammlung
des
Bundesverfassungsgerichts
(C-BVerfGE)**

CODEBOOK

Version 2024-03-08



DOI: 10.5281/zenodo.10783177

Titel	Corpus der amtlichen Entscheidungssammlung des Bundesverfassungsgerichts
Abkürzung	C-BVerfGE
Autor	Seán Fobbe
Version	2024-03-08
Download	https://doi.org/10.5281/zenodo.10783177
Lizenz	CC0 1.0 Universal

Zitiervorschlag

Seán Fobbe (2024). Corpus der amtlichen Entscheidungssammlung des Bundesverfassungsgerichts (C-BVerfGE). Version 2024-03-08. Zenodo. DOI: 10.5281/zenodo.10783177.

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 2024-03-08. 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.3831111. 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	5
2	Nutzung	6
2.1	CSV-Dateien	6
2.2	TXT-Dateien	6
3	Konstruktion	7
3.1	Beschreibung des Datensatzes	7
3.2	Datenquellen	7
3.3	Sammlung der Daten	7
3.4	Source Code und Compilation Report	7
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	9
3.7.3	Beispiel eines Dateinamens	9
3.8	Qualitätsprüfung	9
3.9	Grafische Darstellung	9
4	Varianten und Zielgruppen	10
5	Variablen (Allgemein)	12
5.1	Hinweise	12
5.2	Erläuterungen der einzelnen Variablen	12
6	Variablen (Linguistische Annotationen)	18
6.1	Datenstruktur	18
6.2	Hinweise	18
6.3	Erläuterung der Variablen	19
7	Registerzeichen	20
8	Präsident:innen	21
8.1	Hinweise	21
8.2	Lebensdaten	21
8.3	Dienstalter und Lebensalter	21
9	Vize-Präsident:innen	22
9.1	Hinweise	22
9.2	Lebensdaten	22
9.3	Dienstalter und Lebensalter	23
10	Linguistische Kennzahlen	24
10.1	Erläuterung der Kennzahlen und Diagramme	24
10.2	Werte der Kennzahlen	24
10.3	Verteilung Zeichen	25
10.4	Verteilung Tokens	25
10.5	Verteilung Typen	26

10.6	Verteilung Sätze	26
11	Inhalt des Korpus	27
11.1	Zusammenfassung	27
11.2	Nach Typ der Entscheidung	27
11.3	Nach Typ des Spruchkörpers	28
11.4	Nach Spruchkörper (Aktenzeichen)	29
11.5	Nach Registerzeichen	30
11.6	Nach Präsident:in	32
11.7	Nach Vize-Präsident:in	33
11.8	Nach Entscheidungsjahr	34
11.9	Nach Eingangsjahr (ISO)	36
12	Dateigrößen	39
12.1	Verteilung PDF-Dateigrößen	39
12.2	Verteilung TXT-Dateigrößen	39
12.3	Gesamtgröße je ZIP-Archiv	40
13	Signaturprüfung	41
13.1	Allgemeines	41
13.2	Persönliche GPG-Signatur	41
14	Changelog	42
14.1	Version 2024-03-08	42
14.2	Version 2023-02-20	42
14.3	Version 2022-06-20	42
14.4	Version 2021-09-19	42
14.5	Version 2021-01-03	43
14.6	Version 1.1.0	43
14.7	Version 1.0.0	43
15	Parameter für strenge Replikationen	44
	Literaturverzeichnis	45

1 Einführung

Das **Bundesverfassungsgericht (BVerfG)** ist das höchste Gericht der Bundesrepublik Deutschland und ein Verfassungsorgan. Als »Hüter der Verfassung« ist es seit seiner Gründung im Jahr 1951 mit der Auslegung und Durchsetzung des Grundgesetzes betraut.

Seine Bedeutung im Verfassungsgefüge der Bundesrepublik Deutschland ist kaum zu überschätzen. So richtet es nicht nur über Streitigkeiten zwischen Verfassungsorganen und über Normenkontrollanträge, welche die Nichtigkeit von Gesetzen zur Folge haben können, sondern auch über »Verfassungsbeschwerden, die von jedermann mit der Behauptung erhoben werden können, durch die öffentliche Gewalt in einem seiner Grundrechte oder in einem seiner [grundrechtsgleichen Rechte] verletzt zu sein« (Art. 93 Abs. 4b GG). Das Instrument der Verfassungsbeschwerde ist in seiner Beliebtheit und Effektivität in der Geschichte Deutschlands beispiellos und von hoher wissenschaftlicher und praktischer Bedeutung. In nicht wenigen Verfahrensarten haben die Entscheidungen des BVerfG zudem Gesetzeskraft (§ 31 Abs. 2 BVerfGG).

Die quantitative Analyse von juristischen Texten, insbesondere denen des Bundesverfassungsgerichts, ist in den deutschen Rechtswissenschaften ein noch junges und kaum bearbeitetes Feld.¹ 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*² und *openJur*³ 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.⁴

¹ Besonders positive Ausnahmen finden sich unter: <https://www.quantitative-rechtswissenschaft.de/>

² <https://opiniojuris.de/>

³ <https://openjur.de/>

⁴ 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. CSV⁵ 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.

Hier empfehle ich für **R** dringend 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 Vorschlag:

```
library(data.table)
dt.bverfg <- 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.bverfg <- readtext("*.txt",
  docvarsfrom = "filenames",
  docvarnames = c("gericht",
    "datum",
    "spruchkoerper_typ",
    "spruchkoerper_az",
    "registerzeichen",
    "eingangsnummer",
    "eingangsjahr_az",
    "kollision",
    "name",
    "band",
    "seite"),
  dvsep = "_",
  encoding = "UTF-8")
```

⁵ 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 vielen Entscheidungen, die in der amtlichen Entscheidungssammlung des Bundesverfassungsgerichts (BVerfGE) veröffentlicht sind. Er enthält alle Entscheidungen, die auf der offiziellen Webseite des Bundesverfassungsgerichts am jeweiligen Stichtag in der Auflistung der Entscheidungen der BVerfGE verlinkt waren. Die Stichtage für jede Version sind in der Versionsnummer festgehalten und für frühere Versionen im Changelog dokumentiert.

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

3.2 Datenquellen

Datenquelle	Fundstelle
Primäre Datenquelle	https://www.bundesverfassungsgericht.de
Source Code	https://doi.org/10.5281/zenodo.10783178
Entscheidungsnamen	https://doi.org/10.5281/zenodo.10783178
BVerfGE-Fundstellen	https://doi.org/10.5281/zenodo.10783178
Personendaten	https://doi.org/10.5281/zenodo.4568682
Registerzeichen	https://doi.org/10.5281/zenodo.4569564

Die Personendaten stammen aus folgendem Datensatz: »Seán Fobbe and Tilko Swalve (2021). Presidents and Vice-Presidents of the Federal Courts of Germany (PVP-FCG). Version 2021-04-08. Zenodo. DOI: 10.5281/zenodo.4568682«.

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.10783178>

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 vollständigen Source Code, 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 des Erstellungs-Prozesses interessieren, lesen Sie diesen bitte zuerst.

3.5 Grenzen des Datensatzes

Nutzer sollten folgende wichtige Grenzen beachten:

1. Der Datensatz enthält nur das, was das Gericht auch tatsächlich veröffentlicht, nämlich begründete Entscheidungen, die auch in der BVerfGE abgedruckt wurden (*publication bias*).
2. Es kann aufgrund technischer Grenzen bzw. Fehler sein, dass manche — im Grunde verfügbare — Entscheidungen nicht oder nicht korrekt abgerufen werden (*automation bias*).
3. Es werden nur PDF- und HTML-Dateien abgerufen (*file type bias*). Manche Entscheidungen sind nur als HTML verfügbar. Die Metadaten der Entscheidungen ohne PDF-Datei werden explizit im Compilation Report dokumentiert.
4. Erst ab dem Jahr 1998 sind begründete Entscheidungen des BVerfG einigermaßen vollständig veröffentlicht, auch wenn frühere Entscheidungen vereinzelt auf der Webseite verfügbar sind (*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«).

Der HTML-Quelltext wurde — wie in jeder HTML-Datei selbst dokumentiert ist — mit dem »Government Site Builder« der Bundesverwaltung erstellt, d.h. computergeneriert. Durch Maschinen generierte Texte sind keine »persönliche geistige Schöpfung« iSv § 2 Abs. 2 UrhG und daher urheberrechtlich nicht geschützt. Den verbleibenden Text-Bestandteilen (z.B. Buttons) fehlt es mindestens an der Schöpfungshöhe. Bilder oder andere Texte als Entscheidungstexte werden nicht abgerufen.

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 wurden weitgehend aus den Hyperlinks zur jeweiligen Datei und dem HTML-Quelltext extrahiert. Hinzugefügt wurden von mir eine Reihe weitere Variablen, sowie Unter- und Trennstriche um die Maschinenlesbarkeit zu erleichtern. Der volle Satz

an Metadaten ist nur in den CSV-Dateien enthalten. Alle hinzugefügten Metadaten sind zusammen mit dem Source Code vollständig maschinenlesbar dokumentiert und liegen entweder im CSV-Format vor oder sind direkt im Source Code enthalten.

Die Dateinamen der PDF- und TXT-Dateien enthalten Gerichtsname, Datum (Langform nach ISO-8601, d.h. YYYY-MM-DD), den Typ des Spruchkörpers, das offizielle Aktenzeichen, eine Kollisions-ID, den Namen der Entscheidung, sowie die BVerfGE-Fundstelle (Band und Seite).

3.7.2 Schema für die Dateinamen

```
[gericht]_[datum]_[spruchkoerper_typ]_[spruchkoerper_az]_[registerzeichen]_[  
eingangsnummer]_[eingangsjahr_az]_[kollision]_[name]_[band]_[seite]
```

3.7.3 Beispiel eines Dateinamens

```
BVerfG_1997-07-08_S_1_BvR_1243_95_NA_Partielehrer_96_152.txt
```

3.8 Qualitätsprüfung

Die Typen der Variablen wurden mit *regular expressions* strikt validiert. Die möglichen Werte der jeweiligen Variablen wurden zudem durch Frequenztabellen und Visualisierungen auf ihre Plausibilität geprüft. Insgesamt werden zusammen mit jeder Kompilierung Dutzende Tests zur Qualitätsprüfung durchgeführt. Alle Ergebnisse der Qualitätsprüfungen sind aggregiert im Compilation Report zusammen mit dem Source Code und einzeln im Archiv »ANALYSE« zusammen mit dem Datensatz veröffentlicht.

3.9 Grafische Darstellung

Die Robenfarbe der Bundesverfassungsrichter ist »scharlachrot«. Der Hex-Wert hierfür ist #ca2129. Das ist besonders bei der Erstellung thematisch passender Graphen hilfreich. Alle im Compilation Report und diesem Codebook präsentierten Graphen sind in diesem scharlachrot 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 zudem 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 Bundesverfassungsgericht 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.
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.
CSV_Annotiert	Legal Tech/Quantitative Forschung. Alle Entscheidungen in tokenisierter Form mit linguistischen Annotationen. Beachten Sie bitte die besondere Variablen-Struktur unter Punkt 6. Jede Spalte entspricht einer Variable, jede Zeile einem Token.

Variante	Zielgruppe und Beschreibung
CSV_Segmentiert	Legal Tech/Quantitative Forschung. Experimentell! Alle Entscheidungen in segmentierter Form, d.h. sie sind in einzelne Text-Abschnitte unterteilt (z.B. Leitsätze, Entscheidungsformel, Begründung, Unterschriften). Manche Teile einer Entscheidung sind bewusst nicht enthalten (z.B. lange Zitate aus Gesetzen), weil diese nicht die eigentliche Aktivität des Gerichts wiedergeben. Die Nummerierung der Leitsätze und Absätze der Begründung sollte in der Regel (aber nicht immer!) der originalen Nummerierung in der PDF-Fassung entsprechen. Diese Fassung wurde aus den HTML-Dateien gewonnen und ist noch experimentell.
HTML	Legal Tech/Quantitative Forschung. Die HTML-Dokumente wie sie vom Bundesverfassungsgericht auf der amtlichen Webseite bereitgestellt werden, mit originalen Dateinamen.
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 qualitative Forscher:innen, die nur wenig Speicherplatz oder eine langsame Internetverbindung zur Verfügung haben oder für quantitative Forscher:innen, die beim Einlesen der CSV-Dateien Probleme haben.
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 (Allgemein)

5.1 Hinweise

- Fehlende Werte sind immer mit »NA« codiert
- Strings können grundsätzlich alle in UTF-8 definierten Zeichen (insbesondere Buchstaben, Zahlen und Sonderzeichen) enthalten.

5.2 Erläuterungen der einzelnen Variablen

Variable	Typ	Erläuterung
doc_id	String	(Nur CSV-Datei) Der Name der extrahierten TXT-Datei.
text	String	(Nur CSV-Datei) Der vollständige Inhalt der Entscheidung, so wie er in der von www.bundesverfassungsgericht.de heruntergeladenen PDF-Datei dokumentiert ist. In der segmentierten Variante stammt der Text aus der HTML-Dateien.
segment	String	(Nur segmentierte Variante) Das Segment der Entscheidung. Bezieht sich auf die Variable »text«. Segmentarten sind »leitsatz« (Leitsätze), »gegenstand« (Entscheidungsgegenstand), »formel« (Entscheidungsformel), »tenor« (Tenor), »gruende« (Entscheidungsgründe, ggf. mit Anmerkung ob Sondervotum) und »unterschriften« (Unterschriften der Richter:innen). Die Erkennung von Sondervoten ist noch fehleranfällig. Einzelne Segmente sind mit einer Kombination aus Art und Ordinalzahl definiert, z.B. »gruende-133-sondervotum«.
gericht	String	In diesem Datensatz ist nur der Wert »BVerfG« vergeben. Dies ist der ECLI-Gerichtscode für »Bundesverfassungsgericht«. Diese Variable dient vor allem zur einfachen und transparenten Verbindung der Daten mit anderen Datensätzen.
datum	Datum (ISO)	Das Datum der Entscheidung im Format YYYY-MM-DD (Langform nach ISO-8601). Die Langform ist für Menschen einfacher lesbar und wird maschinell auch öfter automatisch als Datumsformat erkannt.
entscheidung_typ	String	(Nur CSV-Datei) Der Typ der Entscheidung. Es sind die Werte »B« (Beschluss) und »U« (Urteil) vergeben. Wurde durch <i>regular expressions</i> aus der Variable »zitiervorschlag« berechnet.

Variable	Typ	Erläuterung
spruchkoerper_typ	String	Der Typ des Spruchkörpers. Es sind die Werte »K« (Kammer), »S« (Senat), »P« (Plenum) und »B« (Beschwerdekammer gem. § 97c BVerfGG) vergeben.
spruchkoerper_az	Natürliche Zahl	Der im Aktenzeichen angegebene Spruchkörper. Es sind nur die Werte »1« und »2« vergeben. Die Werte stehen für den 1. oder 2. Senat des Gerichts. Für Verzögerungsentscheidungen der Beschwerdekammer ist der Wert »NA«. Achtung: Um die Entscheidungen eines bestimmten Senats zu analysieren reicht es nicht, die Variable »spruchkoerper_az« zu nutzen, es muss zusätzlich noch die Variable »spruchkoerper_typ« auf »S« gesetzt werden, weil ansonsten noch mit dem Senat assoziierte Entscheidungen seiner Kammern und des Plenums mit ausgewählt werden.
registerzeichen	String	Das amtliche Registerzeichen. Es gibt die Verfahrensart an, in der die Entscheidung ergangen ist. Eine Erläuterung der Registerzeichen findet sich unter Punkt 7.
verfahrensart	String	Die ausführliche Beschreibung der Verfahrensart, die dem Registerzeichen zugeordnet ist. Eine Erläuterung der Registerzeichen und der zugehörigen Verfahrensarten findet sich unter Punkt 7.
eingangsnummer	Natürliche Zahl	Verfahren des gleichen Eingangsjahres erhalten vom Gericht eine Nummer in der Reihenfolge ihres Eingangs. Die Zahl ist in den Dateinamen mit führenden Nullen (falls <1000) codiert.
eingangsjahr_az	Natürliche Zahl	Das im Aktenzeichen angegebene Jahr in dem das Verfahren beim Gericht anhängig wurde. Das Format ist eine zweistellige Jahreszahl (YY).
eingangsjahr_iso	Natürliche Zahl	(Nur CSV-Datei) Das nach ISO-8601 codierte Jahr in dem das Verfahren beim Bundesverfassungsgericht anhängig wurde. Das Format ist eine vierstellige Jahreszahl (YYYY), um eine maschinenlesbare und eindeutige Jahreszahl für den Eingang zur Verfügung zu stellen. Wurde aus der Variable »eingangsjahr_az« durch den Autor des Datensatzes berechnet, unter der Annahme, dass Jahreszahlen über 50 dem 20. Jahrhundert zugeordnet sind und andere Jahreszahlen dem 21. Jahrhundert.

Variable	Typ	Erläuterung
entscheidungsjahr	Natürliche Zahl	(Nur CSV-Datei) Das Jahr in dem die Entscheidung ergangen ist. Das Format ist eine vierstellige Jahreszahl (YYYY). Wurde aus der Variable »datum« durch den Autor des Datensatzes berechnet.
kollision	String	In wenigen Fällen sind am gleichen Tag mehrere Entscheidungen zum gleichen Aktenzeichen ergangen. Diese werden ab der zweiten Entscheidung pro Tag durch eine Kollisions-ID mit einem Kleinbuchstaben ausdifferenziert. Für die erste Entscheidung ist der Wert der Variable »NA«, also nicht vorhanden. Die zweite Entscheidung ist mit »a« identifiziert, die dritte mit »b« und so fort. In der offiziellen Beschreibung der ECLI-Ordinalzahl wird diese Variable als »Kollisionsnummer« bezeichnet. Buchstaben sind allerdings keine Nummern, daher die abweichende Bezeichnung.
name	String	Der Name der Entscheidung. Für viele Entscheidungen aus der amtlichen Sammlung sind bekannte Namen vorhanden, diese wurden benutzt soweit möglich und auffindbar. Für weniger bekannte Entscheidungen wurde ein möglichst informativer Name vom Autor vergeben. Die konkrete Darstellung (ohne Leerzeichen, mit Bindestrichen usw.) ist Gründen der maschinellen Lesbarkeit geschuldet.
band	Natürliche Zahl	Der Band der amtlichen Sammlung in dem die Entscheidung veröffentlicht ist.
seite	Natürliche Zahl	Die genaue Fundstelle (Seitenzahl) der Entscheidung im jeweiligen Band der amtlichen Sammlung. Nur sinnvoll nutzbar im Zusammenspiel mit der Variable »band«.
aktenzeichen	String	(Nur CSV-Datei) Das amtliche Aktenzeichen. Die Variable wurde aus den Variablen »spruchkoerper_az«, »registerzeichen«, »eingangsnummer« und »eingangsjahr_az« durch den Autor des Datensatzes berechnet. Im Falle mehrere verbundener Verfahren mit einer einheitlichen Entscheidung ist dies das Aktenzeichen des Pilotverfahrens.

Variable	Typ	Erläuterung
aktenzeichen_alle	String	(Nur CSV-Datei) Alle Aktenzeichen der von der Entscheidung betroffenen Verfahren, falls es sich um verbunden Verfahren mit einheitlicher Entscheidung handelt. Ansonsten ist der Wert dieser Variable identisch mit der Variable »aktenzeichen«.
ecli	String	(Nur CSV-Datei) Der European Case Law Identifier (ECLI) der Entscheidung. Jeder Entscheidung ist eine einzigartige ECLI zugewiesen, ggf. mit Kollisions-ID. Die ECLI ist vor allem dann hilfreich, wenn dieser Datensatz mit anderen Datensätzen zusammengeführt werden und Doppelungen vermieden werden sollen. Alle inhaltlichen Bestandteile der ECLI sind in diesem Datensatz zusätzlich auch anderen und besser verständlichen Variablen zugewiesen. Nutzen Sie bevorzugt diese anderen Variablen, statt Informationen aus der ECLI zu extrahieren. Die Variable wurde aus den Variablen »entscheidungsjahr«, »spruchkoerper_typ«, »datum«, »kollision«, »spruchkoerper_az«, »registerzeichen«, »eingangsnummer« und »eingangsjahr_az« durch den Autor des Datensatzes berechnet.
zitiervorschlag	String	(Nur CSV-Datei) Der vom BVerfG vorgegebene Zitiervorschlag.
kurzbeschreibung	String	(Nur CSV-Datei) Kurzbeschreibung des Inhalts des Verfahrens wie auf der Website des BVerfG angegeben.
pressemittteilung	String	(Nur CSV-Datei) Nummer und Datum der zugehörigen Pressemitteilung, falls vorhanden. Ansonsten »NA«.
praesi	String	(Nur CSV-Datei) Der Nachname des oder der Präsident:in in dessen/deren Amtszeit das Datum der Entscheidung fällt.
v_praesi	String	(Nur CSV-Datei) Der Nachname des oder der Vize-Präsident:in in dessen/deren Amtszeit das Datum der Entscheidung fällt.
richter	String	(Nur CSV-Datei) Die Nachnamen der Richter:innen, die die Entscheidung unterschrieben haben. Ggf. mit Angabe falls die Person verhindert war. Die einzelnen Namen sind jeweils durch vertikale Striche (» «) voneinander getrennt.
zeichen	Natürliche Zahl	(Nur CSV-Datei) Die Anzahl Zeichen eines Dokumentes.

Variable	Typ	Erläuterung
tokens	Natürliche Zahl	(Nur CSV-Datei) 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.
typen	Natürliche Zahl	(Nur CSV-Datei) 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.
saetze	Natürliche Zahl	(Nur CSV-Datei) 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.
version	Datum (ISO)	(Nur CSV-Datei) 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.

Variable	Typ	Erläuterung
doi_concept	String	(Nur CSV-Datei) 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.
doi_version	String	(Nur CSV-Datei) 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.
lizenz	String	Die Lizenz für den Gesamtdatensatz. In diesem Datensatz immer »Creative Commons Zero 1.0 Universal«.

6 Variablen (Linguistische Annotationen)

6.1 Datenstruktur

```
## Classes 'data.table' and 'data.frame': 10692952 obs. of 12 variables:
## $ doc_id : chr "BVerfG_1951-09-09_S_2_BvQ_0001_51_NA_Südweststaat-
Volksabstimmung-EA_1_1.txt" "BVerfG_1951-09-09_S_2_BvQ_0001_51_NA_Sü
dweststaat-Volksabstimmung-EA_1_1.txt" "BVerfG_1951-09-09_S_2_BvQ_0001_51_NA_
Südweststaat-Volksabstimmung-EA_1_1.txt" "BVerfG_1951-09-09_S_2_BvQ_0001_51_
NA_Südweststaat-Volksabstimmung-EA_1_1.txt" ...
## $ sentence_id : int 1 1 1 1 2 2 2 2 3 3 ...
## $ token_id : int 1 2 3 4 1 2 3 4 1 2 ...
## $ token : chr "BUNDESVERFASSUNGSGERICHT" "\n\n" "-" "2" ...
## $ lemma : chr "BUNDESVERFASSUNGSGERICHT" "\n\n" "--" "2" ...
## $ pos : chr "PROPN" "SPACE" "PUNCT" "NUM" ...
## $ tag : chr "NE" "_SP" "$(" "CARD" ...
## $ head_token_id: int 1 1 1 1 1 1 4 2 3 1 ...
## $ dep_rel : chr "ROOT" "dep" "punct" "punct" ...
## $ entity : chr "MISC_B" "MISC_I" "MISC_I" "MISC_I" ...
## $ nounphrase : chr "beg_root" "" "" "" ...
## $ whitespace : logi FALSE FALSE TRUE TRUE TRUE FALSE ...
## - attr(*, ".internal.selfref")=<externalptr>
```

6.2 Hinweise

Diese Variante des Datensatzes beruht nur auf den Variablen »doc_id« und »text« des regulären Datensatzes, die tokenisiert und mittels der Software »spacy«⁶ mit linguistischen Annotationen versehen wurden.

Die Metadaten des Gesamtdatensatzes sind nicht in der linguistische annotierten Fassung enthalten, weil die Bereitstellung von Metadaten für jedes Token die Dateigröße und damit auch den RAM-Bedarf für Analysen gewaltig steigern würde. Um anhand der Metadaten Teilmengen der linguistischen Annotationen zu bilden, gehen sie bitte wie folgt vor:

1. CSV-Datei mit Metadaten einlesen
2. Anhand der Metadaten die gewünschte Teilmenge der Dokumente bilden
3. CSV-Datei mit Linguistischen Annotationen einlesen
4. Die Werte der Variable »doc_id« der Teilmenge nutzen um aus den annotierten Daten nur diejenigen herauszufiltern, deren Variable »doc_id« mit der Teilmenge übereinstimmt

⁶ Die den Annotationen zugrundeliegende Software ist *spacy*, die hier verfügbar ist <https://spacy.io/>. Diese wird in R mittels des *spacyr* packages integriert: <https://spacyr.quanteda.io/>.

6.3 Erläuterung der Variablen

Variable	Typ	Erläuterung
doc_id	String	Der Dateiname des Dokumentes, aus dem die Tokens stammen. Identische Werte wie im Hauptdatensatz. Geeignet um Metadaten mit den linguistischen Annotationen zu verbinden.
sentence_id	Natürliche Zahl	Die Ordinalzahl des Satzes in dem Dokument, dem das Token zugeordnet ist.
token_id	Natürliche Zahl	Die Nummer des Tokens in einem Dokument.
token	String	Einzelne Tokens in der Reihenfolge ihres Vorkommens im jeweiligen Dokument.
lemma	String	Die lemmatisierte Form des Tokens.
pos	String	Grobe Annotation des einzelnen Tokens nach dem universal dependency POS tagset, siehe auch http://universaldependencies.org/u/pos/ .
tag	String	Feine Annotation des einzelnen Tokens mit dem »de_core_news_sm«-Modell von spacy. Für eine detaillierte Erläuterung der Annotationen siehe: https://spacy.io/models/de
head_token_id	Natürliche Zahl	Das führende Token.
dep_rel	String	Die <i>dependency relation</i> zum head_token.
entity	String	Erkennung von benannten Entitäten (Named Entity Recognition).
nounphrase	String	Erkennung von Nominalphrasen.
whitespace	Logisch	Ob es sich bei dem Token um Whitespace handelt oder nicht.

7 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.”

Registerzeichen	Verfahrensart
AR	Allgemeines Register: Vorverfahren oder sonstige Verfahrensarten
BvA	Verwirkung von Grundrechten
BvB	Verfassungswidrigkeit von Parteien
BvC	Wahlprüfungsverfahren
BvD	Anklage des Bundespräsidenten
BvE	Organstreitverfahren
BvF	Abstrakte Normenkontrolle
BvG	Bund-Länder-Streitigkeiten
BvH	Andere Streitigkeiten zwischen Bund und Ländern
BvJ	Anklage von Richtern des Bundesverfassungsgerichts
BvK	Landesverfassungsstreitigkeiten
BvL	Konkrete Normenkontrolle
BvM	Feststellung der Anwendbarkeit einer Regel des Völkerge- wohnheitsrechts
BvN	Divergenzvorlagen eines Landesverfassungsgerichts zur Aus- legung des Grundgesetzes
BvO	Fortgeltung vorkonstitutionellen Rechts als Bundesrecht
BvP	Sonstige durch Bundesrecht zugewiesene Verfahren
BvQ	Einstweilige Anordnungen
BvR	Verfassungsbeschwerden; Kommunalverfassungsbeschwerden
BvT	Sonstige Verfahren
PBvS	Beendigung des Richteramtes am Bundesverfassungsgericht
PBvU	Plenarentscheidungen
PBvV	Rechtsgutachten
PKH	Prozesskostenhilfe
Vz	Verzögerungsbeschwerde

8 Präsident:innen

8.1 Hinweise

- Die Personaldaten stammen aus folgendem Datensatz: »Seán Fobbe and Tilko Swalve (2021). Presidents and Vice-Presidents of the Federal Courts of Germany (PVP-FCG). Version 2021-04-08. Zenodo. DOI: 10.5281/zenodo.4568682«.
- Das Datum bezieht sich jeweils auf das Amt als Präsident:in, nicht auf die Amtszeit als Richter:in.

8.2 Lebensdaten

Nachname	Vorname	Amtsantritt	Amtsende	Geboren	Gestorben
Höpker-Aschoff	Hermann	1951-09-07	1954-01-15	1883-01-31	1954-01-15
VACANCY-1	VACANCY-1	1954-01-16	1954-03-22	NA	NA
Wintrich	Josef	1954-03-23	1958-10-19	1891-02-15	1958-10-19
VACANCY-2	VACANCY-2	1958-10-20	1959-01-07	NA	NA
Müller	Gebhard	1959-01-08	1971-12-07	1900-04-17	1990-08-07
Benda	Ernst	1971-12-08	1983-12-19	1925-01-15	2009-03-02
Zeidler	Wolfgang	1983-12-20	1987-11-15	1924-09-02	1987-12-31
Herzog	Roman	1987-11-16	1994-06-30	1934-04-05	2017-01-10
VACANCY-3	VACANCY-3	1994-07-01	1994-09-13	NA	NA
Limbach	Jutta	1994-09-14	2002-04-09	1934-03-27	2016-09-10
Papier	Hans-Jürgen	2002-04-10	2010-03-15	1943-07-06	NA
Voßkuhle	Andreas	2010-03-16	2020-06-21	1963-12-21	NA
Harbarth	Stephan	2020-06-22	NA	1971-12-19	NA

8.3 Dienstalter und Lebensalter

Nachname	Vorname	Alter (Amtsantritt)	Alter (Amtsende)	Alter (Tod)
Höpker-Aschoff	Hermann	68	70	70
Wintrich	Josef	63	67	67
Müller	Gebhard	58	71	90
Benda	Ernst	46	58	84
Zeidler	Wolfgang	59	63	63
Herzog	Roman	53	60	82
Limbach	Jutta	60	68	82
Papier	Hans-Jürgen	58	66	NA
Voßkuhle	Andreas	46	56	NA
Harbarth	Stephan	48	NA	NA

9 Vize-Präsident:innen

9.1 Hinweise

- Die Personaldaten stammen aus folgendem Datensatz: »Seán Fobbe and Tilko Swalve (2021). Presidents and Vice-Presidents of the Federal Courts of Germany (PVP-FCG). Version 2021-04-08. Zenodo. DOI: 10.5281/zenodo.4568682«.
- Das Datum bezieht sich jeweils auf das Amt als Vize-Präsident:in, nicht auf die Amtszeit als Richter:in.

9.2 Lebensdaten

Nachname	Vorname	Amtsantritt	Amtsende	Geboren	Gestorben
Katz	Rudolf	1951-09-07	1961-07-23	1895-11-23	1961-07-23
Wagner	Friedrich Wilhelm	1961-12-19	1967-10-18	1894-02-28	1971-03-27
Seuffert	Walter	1967-10-18	1975-11-07	1907-02-04	1989-12-28
Zeidler	Wolfgang	1975-11-07	1983-12-20	1924-09-02	1987-12-31
Herzog	Roman	1983-12-20	1987-11-16	1934-04-05	2017-01-10
Mahrenholz	Ernst Gottfried	1987-11-16	1994-03-24	1929-06-18	2021-01-28
Limbach	Jutta	1994-03-24	1994-09-14	1934-03-27	2016-09-10
Henschel	Johann Friedrich	1994-09-29	1995-10-13	1931-06-10	2007-03-19
Seidl	Otto	1995-10-13	1998-02-27	1931-12-11	NA
Papier	Hans-Jürgen	1998-02-27	2002-04-10	1943-07-06	NA
Hassemer	Winfried	2002-04-10	2008-05-07	1940-02-17	2014-01-09
Voßkuhle	Andreas	2008-05-07	2010-03-16	1963-12-21	NA
Kirchhof	Ferdinand	2010-03-16	2018-11-30	1950-06-21	NA
Harbarth	Stephan	2018-11-30	2020-06-22	1971-12-19	NA
König	Doris	2020-06-22	NA	1957-06-25	NA

9.3 Dienstalter und Lebensalter

Nachname	Vorname	Alter (Amtsantritt)	Alter (Amtsende)	Alter (Tod)
Katz	Rudolf	55	65	65
Wagner	Friedrich Wilhelm	67	73	77
Seuffert	Walter	60	68	82
Zeidler	Wolfgang	51	59	63
Herzog	Roman	49	53	82
Mahrenholz	Ernst Gottfried	58	64	91
Limbach	Jutta	59	60	82
Henschel	Johann Friedrich	63	64	75
Seidl	Otto	63	66	NA
Papier	Hans-Jürgen	54	58	NA
Hassemer	Winfried	62	68	73
Voßkuhle	Andreas	44	46	NA
Kirchhof	Ferdinand	59	68	NA
Harbarth	Stephan	46	48	NA
König	Doris	62	NA	NA

10 Linguistische Kennzahlen

10.1 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 gelegentlich auch sinnlose Zeichenfolgen enthalten, weil es rein syntaktisch berechnet wird.
Typen	Einzigartige Tokens. Beispiel: wenn das Token »Verfassungsrecht« 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.

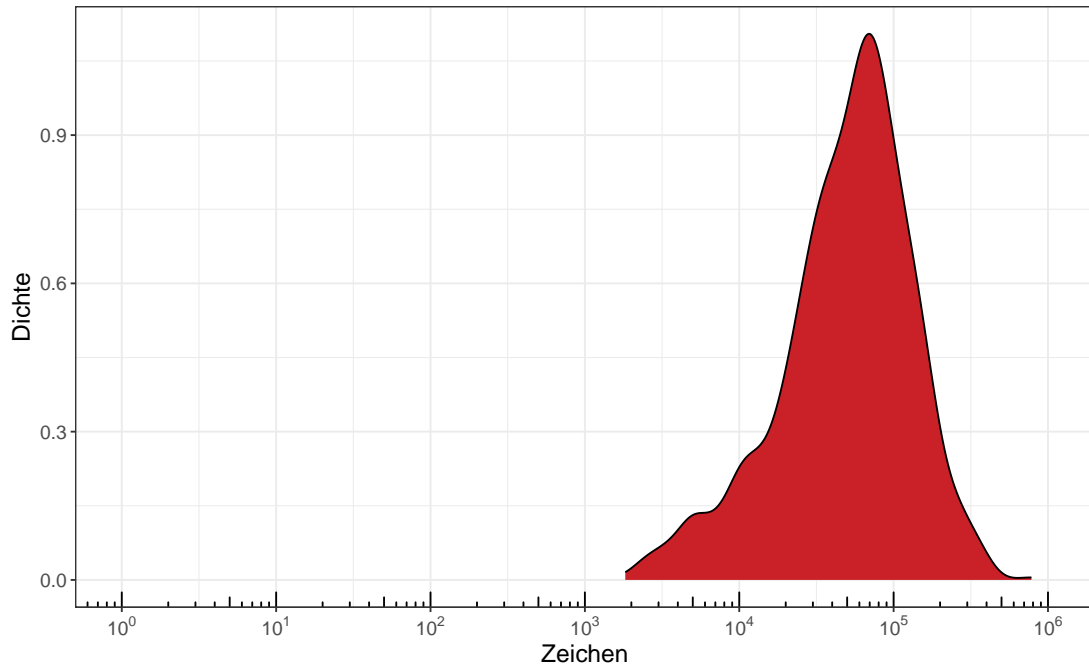
Vorsicht bei der Interpretation: Die x-Achse ist logarithmisch skaliert, d.h. in 10er-Potenzen und damit nicht-linear. Die kleinen Achsen-Markierungen zwischen den Schritten der Exponenten sind eine visuelle Hilfestellung um diese nicht-Linearität zu verstehen.

10.2 Werte der Kennzahlen

Kennzahl	Summe	Min	Quart1	Median	Mittel	Quart3	Max
zeichen	66,748,077	1,827	29,121.00	56,657.5	72,394.88	93,920.50	781,226
tokens	10,312,081	178	4,326.75	8,694.5	11,184.47	14,497.25	115,540
typen	177,109	92	1,171.50	1,908.0	2,072.66	2,674.50	13,491
saetze	558,006	8	250.25	469.5	605.21	783.00	5,282

10.3 Verteilung Zeichen

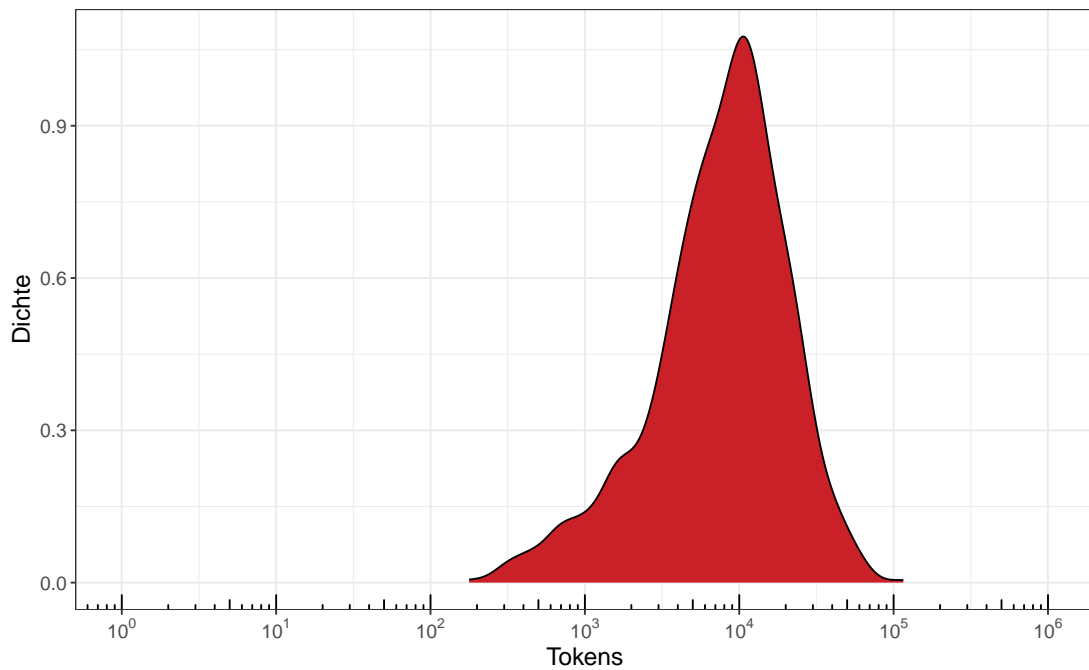
C-BVerfGE | Version 2024-03-08 | Verteilung der Zeichen je Dokument



Fobbe | DOI: 10.5281/zenodo.10783177

10.4 Verteilung Tokens

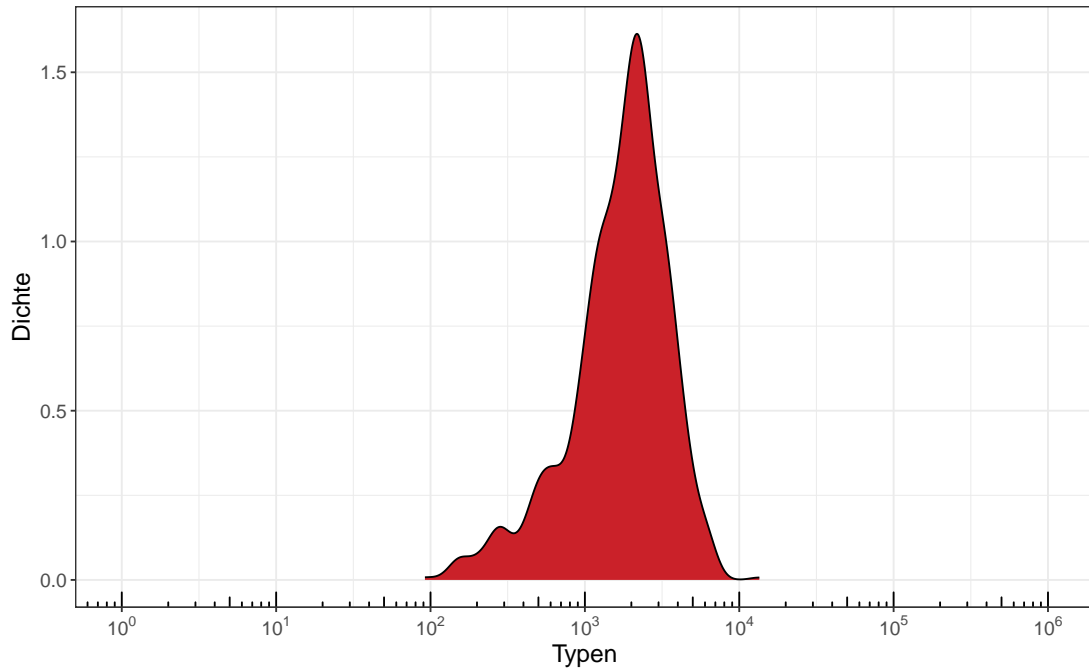
C-BVerfGE | Version 2024-03-08 | Verteilung der Tokens je Dokument



Fobbe | DOI: 10.5281/zenodo.10783177

10.5 Verteilung Typen

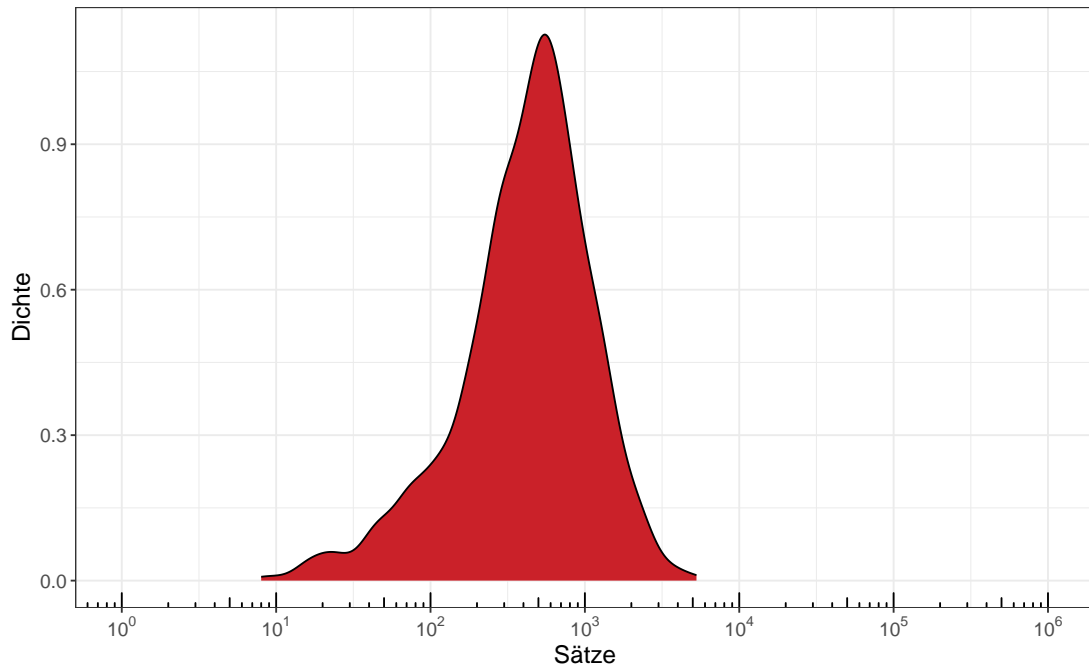
C-BVerfGE | Version 2024-03-08 | Verteilung der Typen je Dokument



Fobbe | DOI: 10.5281/zenodo.10783177

10.6 Verteilung Sätze

C-BVerfGE | Version 2024-03-08 | Verteilung der Sätze je Dokument



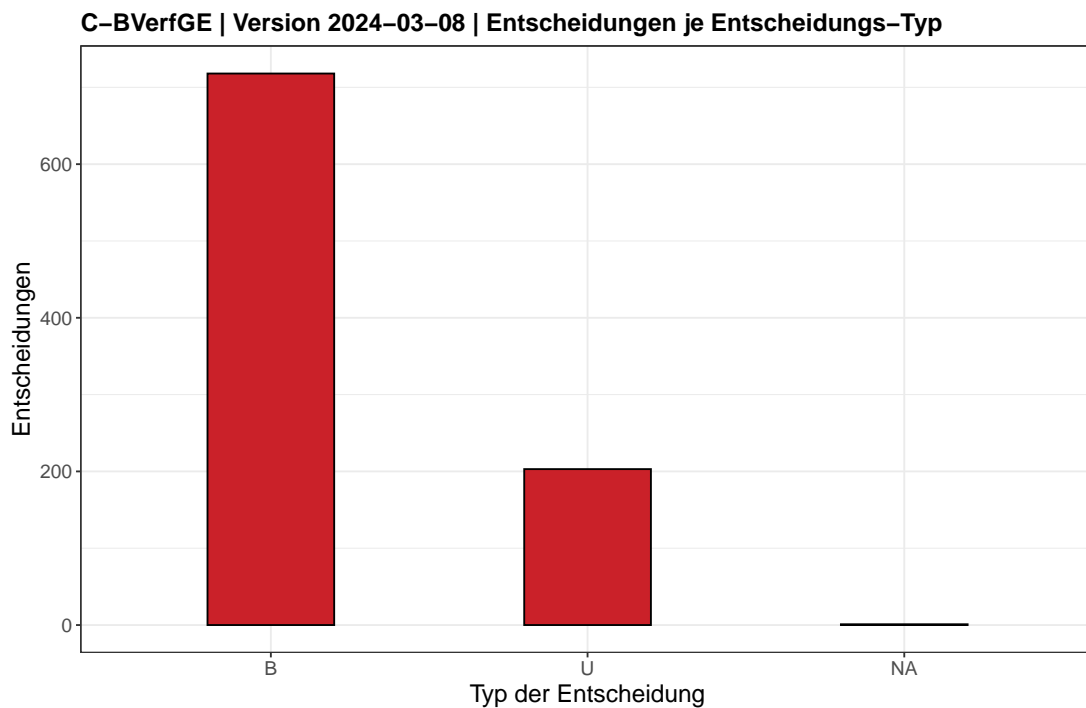
Fobbe | DOI: 10.5281/zenodo.10783177

11 Inhalt des Korpus

11.1 Zusammenfassung

Variable	Anzahl	Min	Quart1	Median	Mittel	Quart3	Max
entscheidungsjahr	40	1951	2001	2007	2007.41	2015.00	2022
eingangsjahr_iso	47	1951	1997	2005	2004.27	2012.00	2022
band	88	1	103	119	121.08	139.75	164
eingangsnummer	461	1	4	207	711.57	1383.00	3588

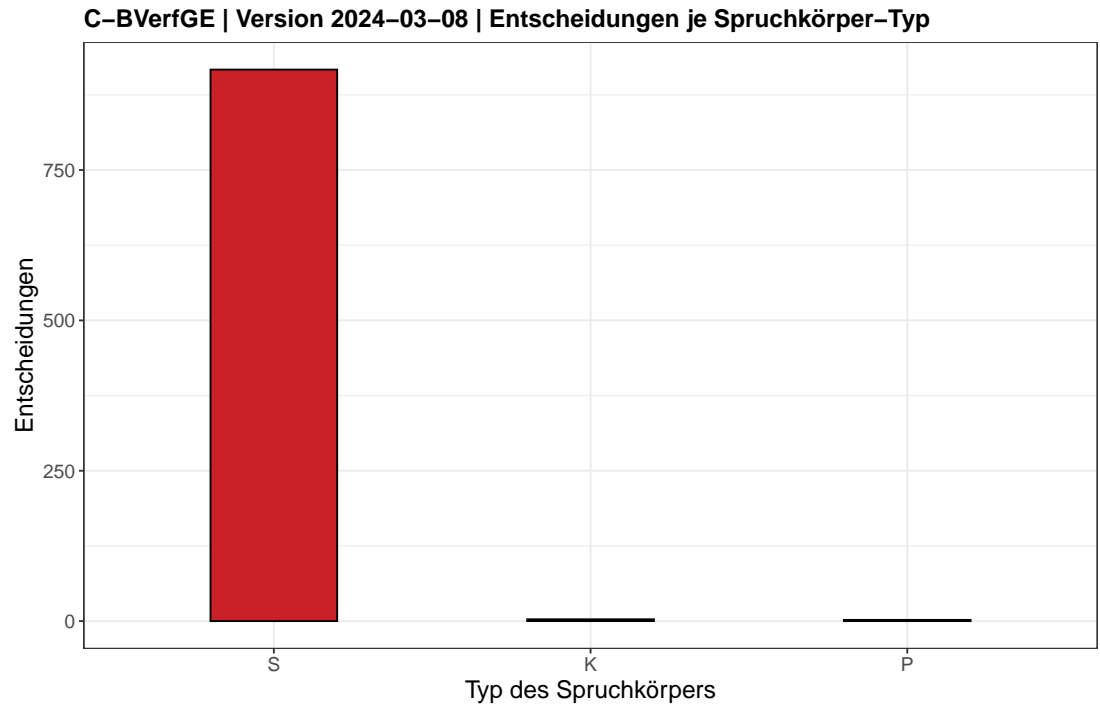
11.2 Nach Typ der Entscheidung



Typ	Entscheidungen	% Gesamt	% Kumulativ
NA	1	0.11	0.11
B	718	77.87	77.98

U	203	22.02	100.00
Total	922	100.00	100.00

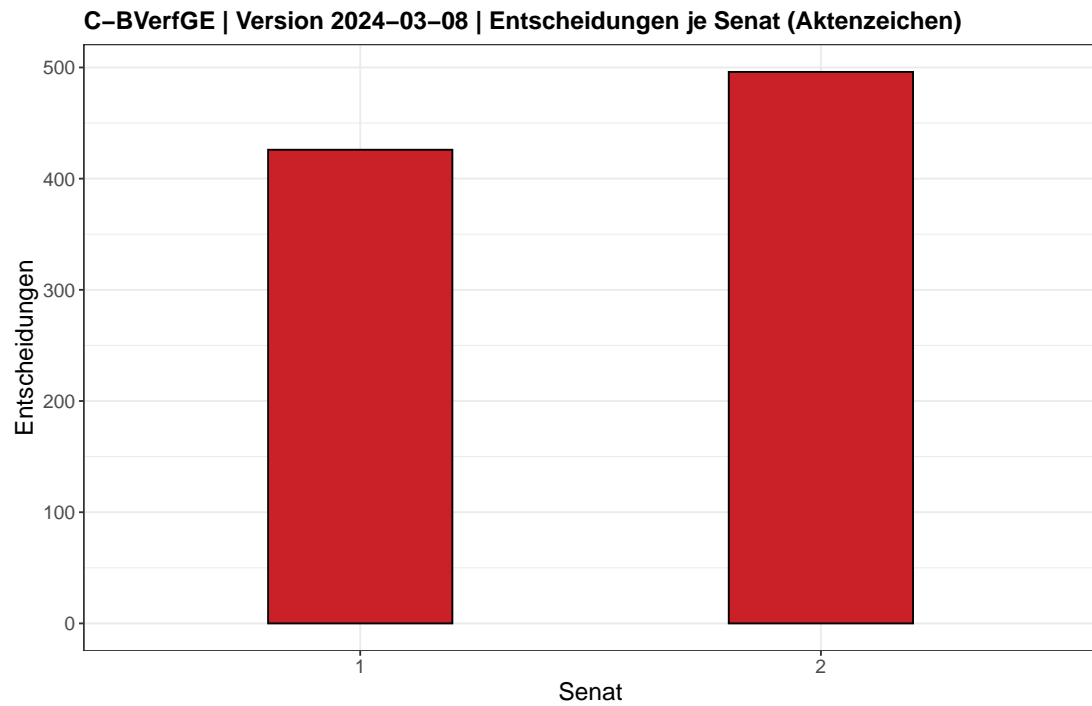
11.3 Nach Typ des Spruchkörpers



Fobbe | DOI: 10.5281/zenodo.10783177

Typ	Entscheidungen	% Gesamt	% Kumulativ
K	3	0.33	0.33
P	2	0.22	0.54
S	917	99.46	100.00
Total	922	100.00	100.00

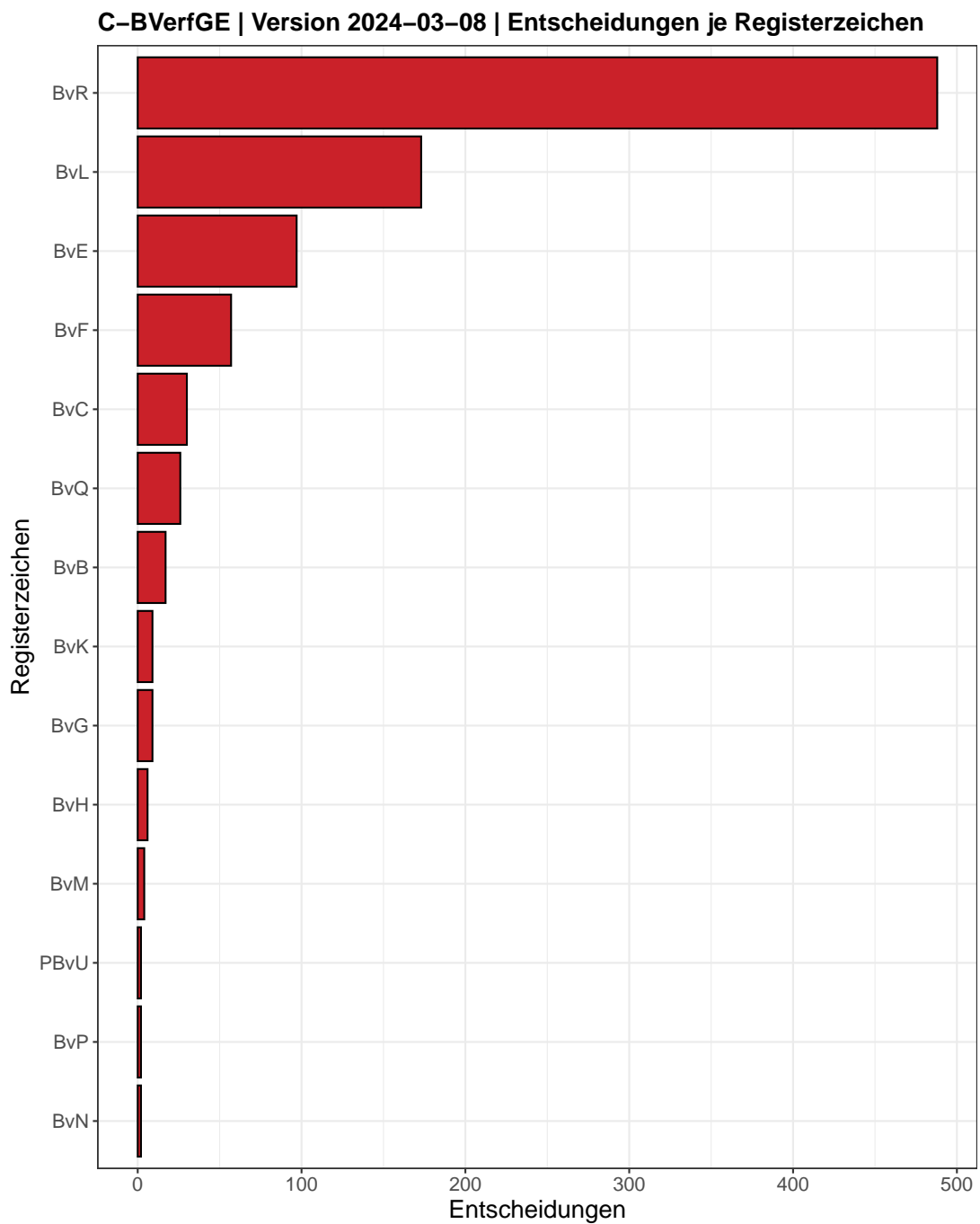
11.4 Nach Spruchkörper (Aktenzeichen)



Fobbe | DOI: 10.5281/zenodo.10783177

Senat	Entscheidungen	% Gesamt	% Kumulativ
1	426	46.2	46.2
2	496	53.8	100.0
Total	922	100.0	100.0

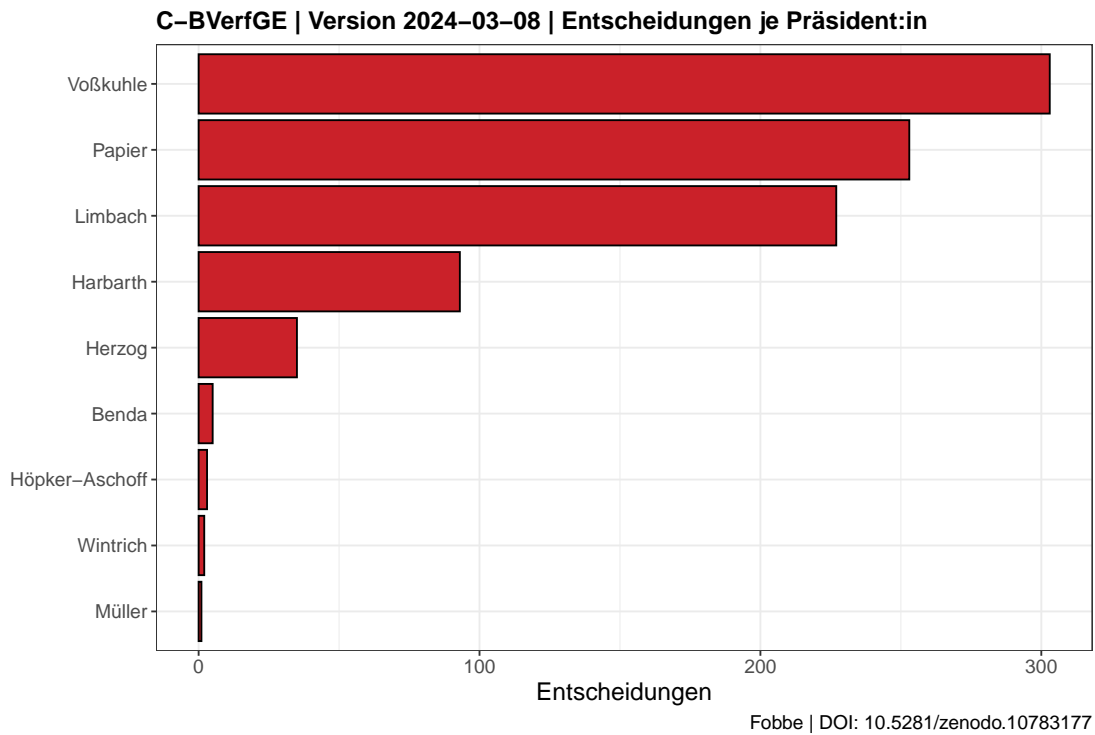
11.5 Nach Registerzeichen



Fobbe | DOI: 10.5281/zenodo.10783177

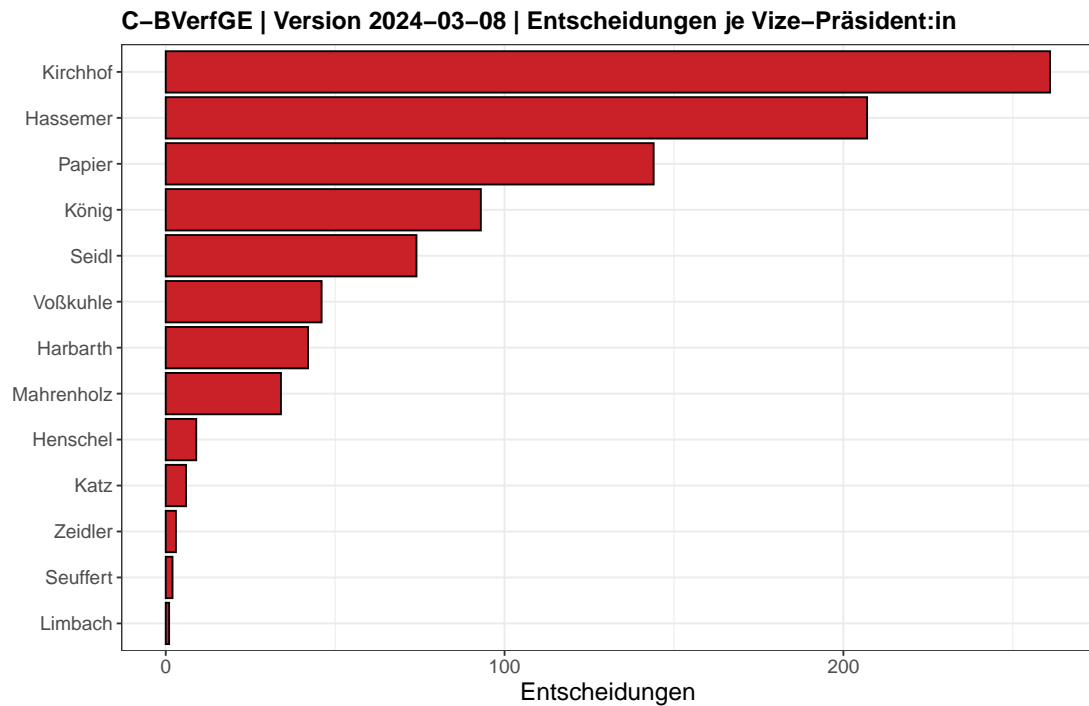
Registerzeichen	Entscheidungen	% Gesamt	% Kumulativ
BvB	17	1.84	1.84
BvC	30	3.25	5.10
BvE	97	10.52	15.62
BvF	57	6.18	21.80
BvG	9	0.98	22.78
BvH	6	0.65	23.43
BvK	9	0.98	24.40
BvL	173	18.76	43.17
BvM	4	0.43	43.60
BvN	2	0.22	43.82
BvP	2	0.22	44.03
BvQ	26	2.82	46.85
BvR	488	52.93	99.78
PBvU	2	0.22	100.00
Total	922	100.00	100.00

11.6 Nach Präsident:in



Präsident:in	Entscheidungen	% Gesamt	% Kumulativ
Benda	5	0.54	0.54
Harbarth	93	10.09	10.63
Herzog	35	3.80	14.43
Höpker-Aschoff	3	0.33	14.75
Limbach	227	24.62	39.37
Müller	1	0.11	39.48
Papier	253	27.44	66.92
Voßkuhle	303	32.86	99.78
Wintrich	2	0.22	100.00
Total	922	100.00	100.00

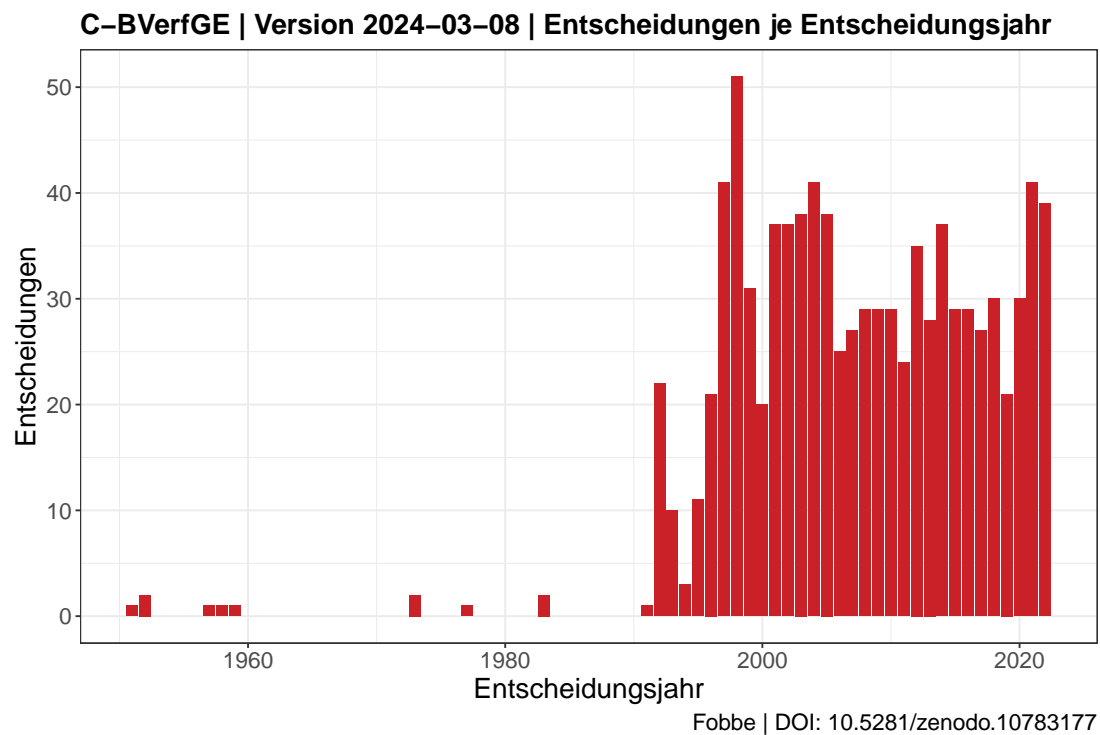
11.7 Nach Vize-Präsident:in



Fobbe | DOI: 10.5281/zenodo.10783177

Vize-Präsident:in	Entscheidungen	% Gesamt	% Kumulativ
Harbarth	42	4.56	4.56
Hassemer	207	22.45	27.01
Henschel	9	0.98	27.98
Katz	6	0.65	28.63
Kirchhof	261	28.31	56.94
König	93	10.09	67.03
Limbach	1	0.11	67.14
Mahrenholz	34	3.69	70.82
Papier	144	15.62	86.44
Seidl	74	8.03	94.47
Seuffert	2	0.22	94.69
Voßkuhle	46	4.99	99.67
Zeidler	3	0.33	100.00
Total	922	100.00	100.00

11.8 Nach Entscheidungsjahr

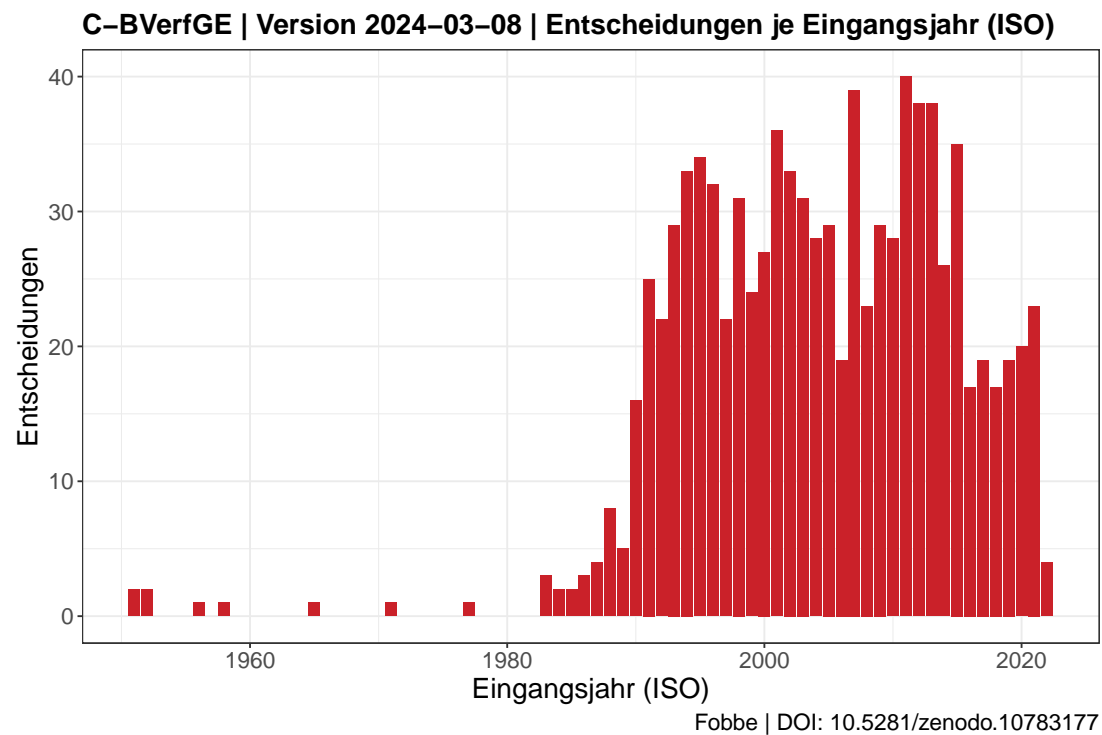


Jahr	Entscheidungen	% Gesamt	% Kumulativ
1951	1	0.11	0.11
1952	2	0.22	0.33
1957	1	0.11	0.43
1958	1	0.11	0.54
1959	1	0.11	0.65
1973	2	0.22	0.87
1977	1	0.11	0.98
1983	2	0.22	1.19
1991	1	0.11	1.30
1992	22	2.39	3.69
1993	10	1.08	4.77
1994	3	0.33	5.10
1995	11	1.19	6.29

(continued)

Jahr	Entscheidungen	% Gesamt	% Kumulativ
1996	21	2.28	8.57
1997	41	4.45	13.02
1998	51	5.53	18.55
1999	31	3.36	21.91
2000	20	2.17	24.08
2001	37	4.01	28.09
2002	37	4.01	32.10
2003	38	4.12	36.23
2004	41	4.45	40.67
2005	38	4.12	44.79
2006	25	2.71	47.51
2007	27	2.93	50.43
2008	29	3.15	53.58
2009	29	3.15	56.72
2010	29	3.15	59.87
2011	24	2.60	62.47
2012	35	3.80	66.27
2013	28	3.04	69.31
2014	37	4.01	73.32
2015	29	3.15	76.46
2016	29	3.15	79.61
2017	27	2.93	82.54
2018	30	3.25	85.79
2019	21	2.28	88.07
2020	30	3.25	91.32
2021	41	4.45	95.77
2022	39	4.23	100.00
Total	922	100.00	100.00

11.9 Nach Eingangsjahr (ISO)



Jahr	Entscheidungen	% Gesamt	% Kumulativ
1951	2	0.22	0.22
1952	2	0.22	0.43
1956	1	0.11	0.54
1958	1	0.11	0.65
1965	1	0.11	0.76
1971	1	0.11	0.87
1977	1	0.11	0.98
1983	3	0.33	1.30
1984	2	0.22	1.52
1985	2	0.22	1.74
1986	3	0.33	2.06
1987	4	0.43	2.49
1988	8	0.87	3.36

(continued)

Jahr	Entscheidungen	% Gesamt	% Kumulativ
1989	5	0.54	3.90
1990	16	1.74	5.64
1991	25	2.71	8.35
1992	22	2.39	10.74
1993	29	3.15	13.88
1994	33	3.58	17.46
1995	34	3.69	21.15
1996	32	3.47	24.62
1997	22	2.39	27.01
1998	31	3.36	30.37
1999	24	2.60	32.97
2000	27	2.93	35.90
2001	36	3.90	39.80
2002	33	3.58	43.38
2003	31	3.36	46.75
2004	28	3.04	49.78
2005	29	3.15	52.93
2006	19	2.06	54.99
2007	39	4.23	59.22
2008	23	2.49	61.71
2009	29	3.15	64.86
2010	28	3.04	67.90
2011	40	4.34	72.23
2012	38	4.12	76.36
2013	38	4.12	80.48
2014	26	2.82	83.30
2015	35	3.80	87.09
2016	17	1.84	88.94
2017	19	2.06	91.00
2018	17	1.84	92.84
2019	19	2.06	94.90

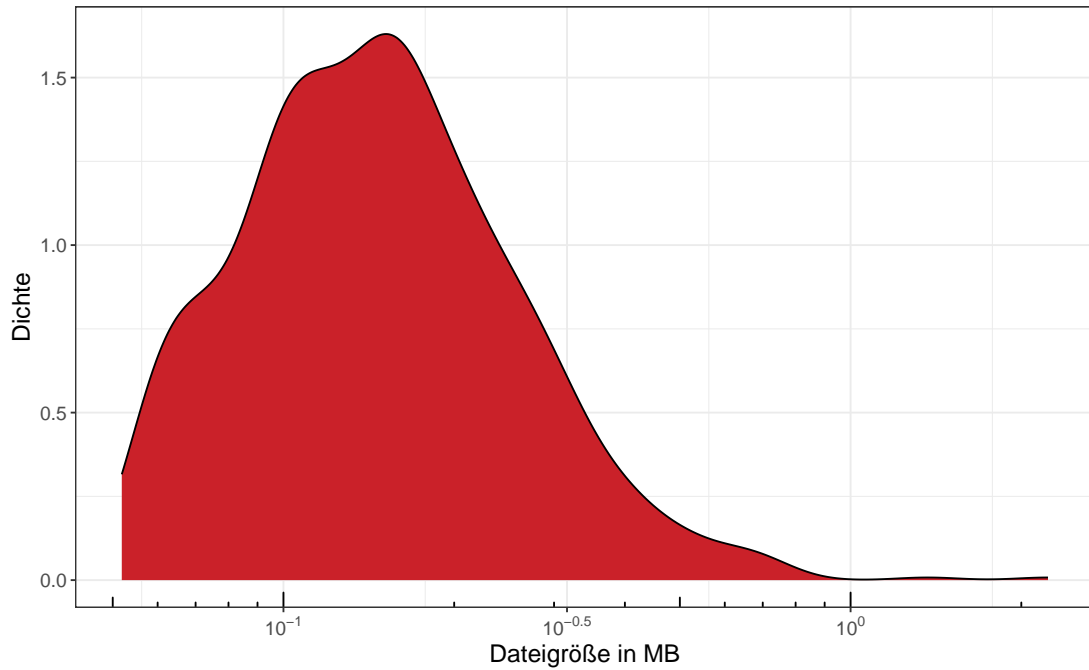
(continued)

Jahr	Entscheidungen	% Gesamt	% Kumulativ
2020	20	2.17	97.07
2021	23	2.49	99.57
2022	4	0.43	100.00
Total	922	100.00	100.00

12 Dateigrößen

12.1 Verteilung PDF-Dateigrößen

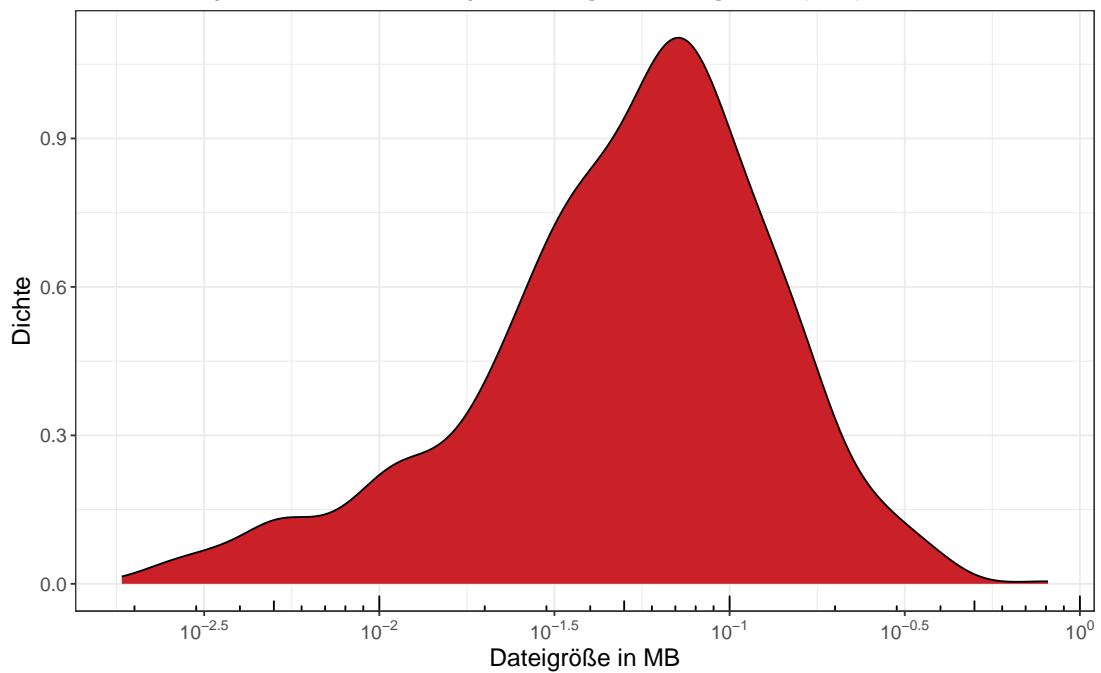
C-BVerfGE | Version 2024-03-08 | Verteilung der Dateigrößen (PDF)



Fobbe | DOI: 10.5281/zenodo.10783177

12.2 Verteilung TXT-Dateigrößen

C-BVerfGE | Version 2024-03-08 | Verteilung der Dateigrößen (TXT)



Fobbe | DOI: 10.5281/zenodo.10783177

12.3 Gesamtgröße je ZIP-Archiv

Datei	Größe in MB
C-BVerfGE_2024-03-08_DE_ANALYSE.zip	1.40
C-BVerfGE_2024-03-08_DE_CSV_Annotiert.zip	126.80
C-BVerfGE_2024-03-08_DE_CSV_Datensatz.zip	19.68
C-BVerfGE_2024-03-08_DE_CSV_Metadaten.zip	0.14
C-BVerfGE_2024-03-08_DE_CSV_Segmentiert.zip	20.67
C-BVerfGE_2024-03-08_DE_HTML_Datensatz.zip	27.73
C-BVerfGE_2024-03-08_DE_PDF_Datensatz.zip	152.28
C-BVerfGE_2024-03-08_DE_TXT_Datensatz.zip	21.54
C-BVerfGE_2024-03-08_Source_Files.zip	0.60

13 Signaturprüfung

13.1 Allgemeines

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 Hash-Werte in zwei verschiedenen Verfahren berechnet und in einer CSV-Datei dokumentiert.

In **Phase II** wird diese CSV-Datei 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 ist.

Dieses Codebook ist vollautomatisch erstellt und prüft die kryptographisch sicheren SHA3-512 Signaturen (»hashes«) aller ZIP-Archive, sowie die GPG-Signatur der CSV-Datei, welche die SHA3-512 Signaturen enthält. SHA3-512 Signaturen werden durch einen system call zur OpenSSL library auf Linux-Systemen berechnet. Eine erfolgreiche Prüfung meldet »Signatur verifiziert!«. Eine gescheiterte Prüfung meldet »FEHLER!«

13.2 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

14 Changelog

14.1 Version 2024-03-08

- Vollständige Aktualisierung der Daten (bis einschließlich Band 164)
- Über 40 neue historische Entscheidungen aus dem Zeitraum 1970 bis 1998 (u.a. Mauerschützen)
- Versionskontrolle aller verwendeten Software mit Docker
- Vereinfachung der Repository-Struktur
- Erstellung des Source-Archivs aus dem Git Manifest
- Entfernung des Submodules und Überführung der Funktionen in das Hautprojekt
- Entfernung des GPG Checks im Codebook um einn durchgängig automatisierten Prozess zu gewährleisten. Mit GPG signiert werden in Zukunft Compilation Report, Codebook und CSV-Datei der Signaturen.
- Funktionen werden nicht mehr vollständig im Compilation Report abgedruckt um die Fehlerrate bei der LaTeX-Kompilierung zu senken

14.2 Version 2023-02-20

- Vollständige Aktualisierung der Daten (bis einschließlich Band 160)
- 50 neue historische Entscheidungen aus dem Zeitraum 1951 bis 1998 (u.a. Elfes, Schleyer-Entführung, Kurzarbeitergeld, Nachtarbeiterinnen)
- Aktenzeichen aus dem Eingangszeitraum 2000 bis 2009 nun korrekt mit führender Null formatiert (z.B. 1 BvR 44/02 statt 1 BvR 44/2)
- Überarbeitung der Namen der Entscheidungen, u.a. Einfügung von Bindestrichen um Lesbarkeit zu verbessern und weitere Standardisierung

14.3 Version 2022-06-20

- Vollständige Aktualisierung der Daten (bis einschließlich Band 158)
- Überarbeitung der codierten Entscheidungsamen
- Standardisierung der Befangenheitsanträge als “Selbstablehnung” und “Richterablehnung” in der Variable “name”
- Strenge Versionskontrolle von R packages mit **renv**
- Kompilierung jetzt detailliert konfigurierbar, insbesondere die Parallelisierung
- Parallelisierung nun vollständig mit *future* statt mit *foreach* und *doParallel*
- Codebook-Erstellung stark beschleunigt durch Verwendung vorberechneter Diagramme
- Fehlerhafte Kompilierungen werden vor der nächsten Kompilierung vollautomatisch aufgeräumt
- Alle Ergebnisse werden automatisch fertig verpackt in den Ordner ‘output’ sortiert
- README und CHANGELOG sind jetzt externe Markdown-Dateien, die bei der Kompilierung automatisch eingebunden werden
- Source Code des Changelogs zu Markdown konvertiert
- REGEX-Tests im Detail kommentiert

14.4 Version 2021-09-19

- Vollständige Aktualisierung der Daten

- Neue Variablen: Lizenz, Typ der Entscheidung, Zeichenzahl, Pressemitteilung, Zitier-vorschlag, Aktenzeichen (alle), Verfahrensart, Kurzbeschreibung und Richter
- Neue Variante: Linguistischen Annotationen
- Neue Variante: Segmentiert
- Neue Variante: HTML
- Erweiterung der Codebook-Dokumentation
- Strenge Kontrolle und semantische Sortierung der Variablen-Namen
- Abgleich der selbst berechneten ECLI mit der in der HTML-Fassung dokumentierten ECLI
- Variable für Entscheidungstyp wird nun aus dem Zitier-vorschlag berechnet um eine höhere Genauigkeit zu gewährleisten

14.5 Version 2021-01-03

- Vollständige Aktualisierung der Daten
- Veröffentlichung des vollständigen Source Codes
- Deutliche Erweiterung des inhaltlichen Umfangs des Codebooks
- Einführung der vollautomatischen Erstellung von Datensatz und Codebook
- Einführung von Compilation Reports um den Erstellungsprozess exakt zu dokumentieren
- Einführung von Variablen für Versionsnummer, Concept DOI, Version DOI, ECLI, Entscheidungs-namen, BVerfGE-Band, BVerfGE-Seite, Typ des Spruchkörpers, Präsi-dentIn, Vize-PräsidentIn und linguistische Kennzahlen (Tokens, Typen, Sätze)
- Automatisierung und Erweiterung der Qualitätskontrolle
- Einführung von Diagrammen zur Visualisierung von Prüfergebnissen
- Einführung kryptographischer Signaturen
- Alle Variablen sind nun in Kleinschreibung und Snake Case gehalten
- Variable ‘Suffix’ in ‘kollision’ umbenannt.
- Variable ‘Ordinalzahl’ in ‘eingangsnummer’ umbenannt.
- Umstellung auf Stichtags-Versionierung

14.6 Version 1.1.0

- Vollständige Aktualisierung der Daten
- Angleichung der Variablen-Namen an andere Datensätze der CE-Serie⁷
- Einführung der Variable ‘Suffix’ um weitere Entscheidungen korrekt erfassen zu können; aufgrund der fehlenden Berücksichtigung des Suffix wurden einige wenige Entscheidungen in Version 1.0.0 irrtümlich von der Sammlung ausgeschlossen.
- Stichtag: 2020-08-09

14.7 Version 1.0.0

- Erstveröffentlichung
- Stichtag: 2020-05-16

⁷ Siehe: <https://zenodo.org/communities/sean-fobbe-data/>

15 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] parallel stats graphics grDevices utils datasets methods
## [8] base
##
## other attached packages:
## [1] future.apply_1.10.0 future_1.32.0 spacyr_1.2.1
## [4] quanteda_3.2.4 readtext_0.81 data.table_1.14.8
## [7] scales_1.2.1 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 RcppTOML_0.2.2
## [16] magick_2.7.4 rmarkdown_2.20
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.10 here_1.0.1 svglite_2.1.1 lattice_0.20-45
## [5] listenv_0.9.0 png_0.1-8 rprojroot_2.0.3 digest_0.6.31
## [9] utf8_1.2.3 parallelly_1.34.0 R6_2.5.1 evaluate_0.20
## [13] pillar_1.8.1 rlang_1.0.6 curl_5.0.0 rstudioapi_0.14
## [17] Matrix_1.5-1 reticulate_1.28 qpdf_1.3.0 labeling_0.4.2
## [21] webshot_0.5.4 stringr_1.5.0 selectr_0.4-2 tinytex_0.44
## [25] munsell_0.5.0 compiler_4.2.2 xfun_0.37 pkgconfig_2.0.3
## [29] askpass_1.1 systemfonts_1.0.4 globals_0.16.2 htmltools_0.5.4
## [33] tidyselect_1.2.0 tibble_3.2.0 codetools_0.2-18 fansi_1.0.4
## [37] viridisLite_0.4.1 dplyr_1.1.0 withr_2.5.0 rappdirs_0.3.3
## [41] grid_4.2.2 jsonlite_1.8.4 gtable_0.3.1 lifecycle_1.0.3
## [45] magrittr_2.0.3 RcppParallel_5.1.7 cli_3.6.0 stringi_1.7.12
## [49] farver_2.1.1 xml2_1.3.3 stopwords_2.3 generics_0.1.3
## [53] vctrs_0.5.2 fastmatch_1.1-3 tools_4.2.2 glue_1.6.2
## [57] fastmap_1.1.1 yaml_2.3.7 colorspace_2.1-0
```

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, and Akitaka Matsuo. 2020. *Spacyr: Wrapper to the spaCy 'Nlp' Library*. <https://spacyr.quanteda.io>.
- Benoit, Kenneth, and Adam Obeng. 2021. *Readtext: Import and Handling for Plain and Formatted Text Files*. <https://github.com/quanteda/readtext>.
- 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>.
- 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>.
- Ewing, Mark. 2021. *Mgsub: Safe, Multiple, Simultaneous String Substitution*.
- Ooms, Jeroen. 2023a. *Magick: Advanced Graphics and Image-Processing in R*.
- . 2023b. *Pdftools: Text Extraction, Rendering and Converting of Pdf Documents*.
- R Core Team. 2022. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- . 2022. *Rvest: Easily Harvest (Scrape) Web Pages*.
- . 2023. *Httr: Tools for Working with Urls and Http*.
- Wickham, Hadley, Winston Chang, Lionel Henry, Thomas Lin Pedersen, Kohske Takahashi, Claus Wilke, Kara Woo, Hiroaki Yutani, and Dewey Dunnington. 2023. *Ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics*.
- Wickham, Hadley, and Dana Seidel. 2022. *Scales: Scale Functions for Visualization*.
- 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 Golemund. 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*.