

Wissenschaftliche Softwareentwicklung

27. Datenformate: beständig und plattformunabhängig

Till Biskup

Physikalisch-Technische Bundesanstalt

19.02.2024





- 🔑 Formate betreffen nicht nur Roh- und verarbeitete Daten, sondern auch Metadaten, Dokumentation, Abbildungen.
- 🔑 Datenformate sollten über Jahrzehnte lesbar, plattformunabhängig, quelloffen und dokumentiert sein.
- 🔑 Daten über Jahrzehnte lesbar zu archivieren, ist nicht nur eine Frage der Formate, sondern auch der Organisation.
- 🔑 Rohdaten sollten immer (im Originalformat) archiviert und vor (ungewollter) Veränderung geschützt werden.
- 🔑 Das konkrete Datenformat oder die Art der Datenlagerung ist für ein System zur Datenverarbeitung irrelevant.

Kriterien für Datenformate in der Wissenschaft

Beispiele plattform- und sprachunabhängiger Formate

Zum Umgang mit Daten und Metadaten

Bedeutung im Gesamtkontext einer Auswertungssoftware

Warum sich mit Datenformaten befassen?

Ein paar Gründe, warum sie jeden Wissenschaftler etwas angehen



- Wir sind alle Nutzer von Datenformaten.
 - Datenformate haben verschiedene Formen/Funktionen.
 - Der *bewusste* Umgang ist auch hier entscheidend.
- Wissenschaftler tragen Verantwortung für ihre Daten.
 - Daten sind die Grundlage aller empirischen Wissenschaften.
 - Nachvollziehbarkeit erfordert Zugriff auf die Daten.
 - Zugriff erfordert offene und plattformunabhängige Formate.
- Datenformate betreffen nicht nur (Roh-)Daten.
 - Metadaten, Dokumentation der Auswertung, Abbildungen
 - Beständigkeit und offener Zugang sind auch hier wichtig.
- Wir alle erzeugen (oft unbewusst) Datenformate.
 - Zwischenergebnisse werden in Dateien abgelegt.
 - Auch ein handgeschriebenes Laborbuch zählt dazu.

Warum sich mit Datenformaten befassen?

Aufbauten aus Einzelkomponenten mit handgeschriebener Steuerung



- Nicht alle Gerätschaften sind aus einem Guss . . .
 - Steuersoftware ist oft selbst geschrieben.
 - gibt freie Hand bei der Wahl der Datenformate
- ☛ Beispiel: Die Steuersoftware für ein EPR-Spektrometer wurde komplett von einem Doktoranden entwickelt.

- zukunftssicher
 - Das Format sollte Jahrzehnte bestehen können.
 - weitgehend ungelöstes Problem (auch für Archivare)
- plattform- und sprachunabhängig
 - mindestens Windows, Linux/Unix, macOS unterstützen
 - unabhängig von der eingesetzten Programmiersprache
- quelloffen, vollständig dokumentiert, standardisiert
 - ermöglicht unabhängigen Zugriff
- geeignet zur Ablage von Metadaten
 - Daten ohne Metadaten sind wertlos.
 - Metadaten so nah wie möglich bei den Daten ablegen
- versioniert
 - Formate entwickeln sich – wenn auch meist eher langsam.
 - Versionen sollten automatisiert erkennbar/auslesbar sein.

proprietäres Datenformat

auf herstellerspezifischen, nicht veröffentlichten Standards basierend
und damit in der Regel nicht von Dritten lesbar

- Probleme
 - nicht zukunftssicher
 - selten plattform- und sprachunabhängig
 - Exportformate meist ohne Metadaten
- Umgang
 - Rohdaten in diesem Format trotzdem archivieren
 - Metadaten soweit verfügbar separat (manuell) ablegen
 - immer Daten exportieren und ebenfalls archivieren

Kriterien für Datenformate in der Wissenschaft

Beispiele plattform- und sprachunabhängiger Formate

Zum Umgang mit Daten und Metadaten

Bedeutung im Gesamtkontext einer Auswertungssoftware

- reine Textformate
 - ✓ universell lesbar (zumindest ASCII 7-bit)
 - ✓ auch vom Menschen lesbar
 - ✓ ohne zusätzliche Software unmittelbar zugänglich
 - ✗ ggf. langsam im Zugriff
 - ✗ mitunter mit erheblichem Speicherbedarf

- Binärformate
 - ✓ oft sehr performant
 - ✓ meist geringer Speicherbedarf
 - ✗ nie ohne spezielle Software lesbar
 - ✗ nicht vom Menschen lesbar

- ☛ Beispiel für die Universalität reiner Textformate:
Unix-Kommandozeile

Format	Struktur
unstrukturierter Text	keine
DSV/CSV	zeilenweise
Windows-INI	blockweise
XML/JSON/YAML	hierarchisch

- Unterschied bzgl. (wiederkehrenden) Strukturen
 - (Mess-)Daten sind meist strukturell hochrepetitiv
 - Metadaten sind i.d.R. hierarchisch
 - unstrukturierter Text ggf. in Einzelfeldern erlaubt (Kommentare)
- Jedes Format hat seine Berechtigung und Anwendungen.
- Die Liste erhebt keinen Anspruch auf Vollständigkeit.

Listing 1: Beispiel für unstrukturierten Text

```
Lorem ipsum dolor sit amet, consectetur adipiscing elit.  
Vestibulum varius gravida arcu, rutrum finibus  
ligula mollis sed.  
Fusce blandit suscipit ultricies. Curabitur lacus libero,  
accumsan sed velit id, semper ultrices ante.
```

```
Vestibulum hendrerit finibus ex sit amet tincidunt. Morbi  
nec quam id arcu convallis porta.
```

- keinerlei Struktur
- Zeilenumbrüche und Trennungen von Absätzen willkürlich
- gut geeignet für Text ohne strikt wiederkehrende Struktur (z.B. Kommentare, Berichte mit viel Fließtext)
- logische Textauszeichnung durch Steuerelemente möglich

Listing 2: Beispiel für zeichengetrennte Werte (DSV)

```
root:x:0:0:root:/root:/bin/zsh
daemon:x:1:1:daemon:/usr/sbin:/bin/sh
bin:x:2:2:bin:/bin:/bin/sh
sys:x:3:3:sys:/dev:/bin/sh
sync:x:4:65534:sync:/bin:/bin/sync
games:x:5:60:games:/usr/games:/bin/sh
man:x:6:12:man:/var/cache/man:/bin/sh
lp:x:7:7:lp:/var/spool/lpd:/bin/sh
mail:x:8:8:mail:/var/mail:/bin/sh
news:x:9:9:news:/var/spool/news:/bin/sh
uucp:x:10:10:uucp:/var/spool/uucp:/bin/sh
```

- zeilenweise Struktur
- eindeutiges Trennzeichen (hier: Doppelpunkt)
- Soll das Trennzeichen in einem Feld verwendet werden, wird ein Backslash (\) vorangestellt.

Listing 3: Beispiel für „kommagetrennte“ Werte (CSV)

```
Year,Make,Model,Description,Price
1997,Ford,E350,"ac, abs, moon",3000.00
1999,Chevy,"Venture "Extended Edition","",,4900.00
1999,Chevy,"Venture "Extended Edition, Very Large","",,5000.00
1996,Jeep,Grand Cherokee,"MUST SELL!
air, moon roof, loaded",4799.00
```

- zeilenweise Struktur
- oft Spaltenüberschriften in erster Zeile
- diverse Trennzeichen möglich, meist Komma
- jeder Eintrag, der (Leer- oder) Sonderzeichen enthält, von Anführungszeichen eingeschlossen
- Format sehr uneinheitlich implementiert (vgl. RFC 4180)

“ *While there are various specifications and implementations for the CSV format [...], there is no formal specification in existence, which allows for a wide variety of interpretations of CSV files. This section documents the format that seems to be followed by most implementations ...*

– Y. Shafranovich, RFC 4180

- ☛ Beispiel dafür, wie ein Format *nicht* aussehen sollte
- ☛ Aufgrund der Uneinheitlichkeit keine Kompatibilität untereinander zu gewährleisten
- ☛ Sollte *nicht* verwendet werden – Alternative: DSV

Listing 4: Beispiel für das Windows-INI-Format

```
; last modified 1 April 2001 by John Doe
[owner]
name=John Doe
organization=Acme Widgets Inc.

[database]
server=192.0.2.62
port=143
file="payroll.dat"
```

- blockweise Struktur
- zwei Hierarchieebenen: Blöcke und Schlüssel-Wert-Paare
- ursprünglich für Konfigurationen entwickelt
- Grundidee lässt sich auch anderweitig einsetzen
 - Beispiel folgt in der nächsten Vorlesung

Listing 5: Beispiel für XML

```
<?xml version="1.0"?>
<verzeichnis>
  <titel>Wikipedia Ortsverzeichnis</titel>
  <eintrag>
    <stichwort>Genf</stichwort>
    <eintragstext>Genf ist der Sitz von ...</eintragstext>
  </eintrag>
  <eintrag>
    <stichwort>Mainz</stichwort>
    <eintragstext>Mainz ist eine Stadt, die ...</eintragstext>
  </eintrag>
</verzeichnis>
```

- beliebig verschachtelte, hierarchische Struktur
- syntaktisch validierbar gegen ein definiertes Schema (ohne Notwendigkeit, die Semantik zu verstehen)
- mühsam für Menschen les- und insbesondere schreibbar

Listing 6: Beispiel für JSON

```
{
  "id": 1,
  "name": "Foo",
  "price": 123,
  "tags": [
    "Bar",
    "Eek"
  ],
  "stock": {
    "warehouse": 300,
    "retail": 20
  }
}
```

- beliebig verschachtelte, hierarchische Struktur
- gut geeignet für die Persistenz von Datenstrukturen
- (deutlich) sparsamer als XML

Listing 7: Beispiel für YAML

```
id: 1
name: Foo
price: 123
tags:
  - Bar
  - Eek
stock:
  warehouse: 300
  retail: 20
```

- beliebig verschachtelte, hierarchische Struktur
- einfacher menschenles- und schreibbar als JSON
 - keine Klammern oder Anführungszeichen notwendig
- Übermenge zu JSON (seit YAML 1.2)
- Unterstützung durch viele Programmiersprachen

Hinweise zum Umgang mit Reintextformaten

Ggf. Beschränkung auf druckbare Zeichen aus dem ASCII-Zeichensatz



Hex	0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
00	NUL ^@ 0	SOH ^A 1	STX ^B 2	ETX ^C 3	EOT ^D 4	ENQ ^E 5	ACK ^F 6	BEL ^G 7	BS ^H 8	TAB ^I 9	LF ^J 10	VT ^K 11	FF ^L 12	CR ^M 13	SO ^N 14	SI ^O 15
10	DLE ^P 16	DC1 ^Q 17	DC2 ^R 18	DC3 ^S 19	DC4 ^T 20	NAK ^U 21	SYN ^V 22	ETB ^W 23	CAN ^X 24	EM ^Y 25	SUB ^Z 26	ESC ^[27	FS ^\ 28	GS ^] 29	RS ^^ 30	US ^? 31
20		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
30	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
40	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
50	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
60	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
70	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

☛ Zeichen 20_{Hex} bis 7e_{Hex} (32 bis 126) sind „druckbar“.

ASCII: American Standard Code for Information Interchange

- Grund für Steuer- und Sonderzeichen
 - Auch ASCII enthält nicht druckbare Zeichen.
 - Notwendigkeit zur Darstellung weiterer Zeichen
- Strategie der geringsten Überraschung (*least surprise*)
 - Unterstützung der „Backslash“-Konvention

Zeichen	Bedeutung
<code>\n</code>	neue Zeile (new line)
<code>\t</code>	Tabulator
<code>\r</code>	Wagenrücklauf (carriage return)
<code>\xnn</code>	Zeichen mit Hexadezimalwert <code>nn</code>
<code>\unnnn</code>	Unicodezeichen mit Hexadezimalwert <code>nnnn</code>
<code>\\</code>	<code>\</code>

Sinnvolle Strategien der Komprimierung:

- Vermeiden redundanter Information
 - Beispiel: äquidistante Achse
Start-, Endwert und Schrittweite reichen aus
- Kompression der gesamten Textdatei
 - Kompression unabhängig vom Dateiformat
 - Platzersparnis mitunter erheblich
 - Beispiel: Open Document Format (OpenOffice etc.)

Was man vermeiden sollte:

- Kompression von Teilen (Feldern) einer Textdatei
- Binärcodierung von Teilen (Feldern) einer Textdatei

- ANSI/IEEE 754-1985
 - Standard für die Repräsentation von Gleitkommazahlen
 - Möglichkeit, *reine* Zahlenkolonnen abzulegen
 - Achtung: mehr als ein Standard, nur für 64 Bit eindeutig
 - benötigt zusätzliche Informationen (Metadaten) zum konkreten Format und der Dimension der Daten

 - HDF5
 - hierarchisches Binärformat
 - selbstbeschreibend mit Unterstützung für Metadaten
 - extrem performant
 - Erlaubt die Arbeit mit (Teilen von) Datensätzen, die viel zu groß für den Arbeitsspeicher wären.
- ☞ Es gibt auch hier noch (viele) weitere Formate . . .

Kriterien für Datenformate in der Wissenschaft

Beispiele plattform- und sprachunabhängiger Formate

Zum Umgang mit Daten und Metadaten

Bedeutung im Gesamtkontext einer Auswertungssoftware

- Rohdaten sollten immer archiviert werden.
 - Was als Rohdaten gilt, ist nicht immer eindeutig...
 - Unversehrtheit sollte sichergestellt und überprüfbar sein
- Daten(sätze) sollten eineindeutig identifizierbar sein.
 - konsistentes Schema für Benennung und Ablage
 - Zugriff weitgehend unabhängig von der Art der Speicherung
- Metadaten sollten mit den Daten verbunden sein.
 - Daten ohne Metadaten sind wertlos.
 - in der gleichen Datei oder direkt daneben ablegen
- Funktionierende Backups sind essentiell.
 - Daten sind wertvoll und von grundlegender Bedeutung.
 - Backups regelmäßig auf korrekte Funktion überprüfen

- Was sind Rohdaten?
 - kontextabhängig – meist relativ klare Vorstellungen
 - Widerstandswerte eines Thermometers wohl eher nicht . . .
 - Rohdaten eines Bildsensors dagegen wohl eher schon . . .
 - Interne Verarbeitungsschritte sollten ggf. bekannt sein.
- Schutz vor ungewollter Veränderung
 - Insbesondere für Rohdaten von essentieller Bedeutung.
 - Bei der Langzeitarchivierung können sich Datenfehler einschleichen bzw. sind oft unvermeidbar.
- Sicherstellung der Unversehrtheit
 - kryptographische Hashes: Überprüfung auf Veränderungen
 - regelmäßig überprüfte Backups zur Datensicherung

kryptographische Hash-Funktion

Funktion, die eine Zeichenfolge beliebiger Länge auf eine solche fester Länge abbildet, kollisionsresistent sein sollte und immer eine Einwegfunktion (unumkehrbar) ist.

- Prüfsummen (*hashes*) viel kleiner als eigentliche Daten
 - lassen sich viel einfacher und schneller vergleichen
 - lassen sich schnell eindeutig aus den Daten erzeugen
- Tipps aus der Praxis
 - Prüfsummen über Daten und Metadaten ggf. trennen
 - ggf. Prüfsumme über Prüfsummen von Einzeldateien bilden
- 👉 Standard-Hash-Algorithmen weit verbreitet und verfügbar

- Zielstellung
 - einfacher Zugriff auf einen Datensatz
 - eindeutiges Etikett, das die Zuordnung ermöglicht
 - Nachvollziehbarkeit von Auswertungen
 - Rohdaten und verarbeitete Daten ansprechbar
 - mögliche Lösung
 - eindeutiger Bezeichner für jeden Datensatz
 - Zugriff erfolgt über Bezeichner und Zuordnungstabelle, die den Bezeichner mit dem Speicherort verknüpft
 - Vorteil: unabhängig vom eingesetzten Ablagesystem (Verzeichnishierarchie, Datenbank, Netzwerkspeicher, ...)
- ☛ Die Zuordnungstabelle sollte gut gesichert werden.
Alternative: Zuordnungstabelle automatisch generierbar

- im Speicher
 - flüchtig, nur für Zwischenschritte sinnvoll
 - Dateien im Dateisystem
 - Standardsituation in den meisten Fällen in der Wissenschaft
 - Die Wahl eines geeigneten Formates ist entscheidend.
 - Datenbank
 - große Vorteile beim gezielten Zugriff
 - in zentralem Repository über Netzwerk
 - hauptsächlich bei sehr großen Datenmengen
oder bei der Notwendigkeit eines verteilten Zugriffs
- ☛ Jede Art hat eigene Ansprüche an das gewählte Format.
- ☛ sollte ein austauschbares Implementationsdetail sein (DIP)

Kriterien für Datenformate in der Wissenschaft

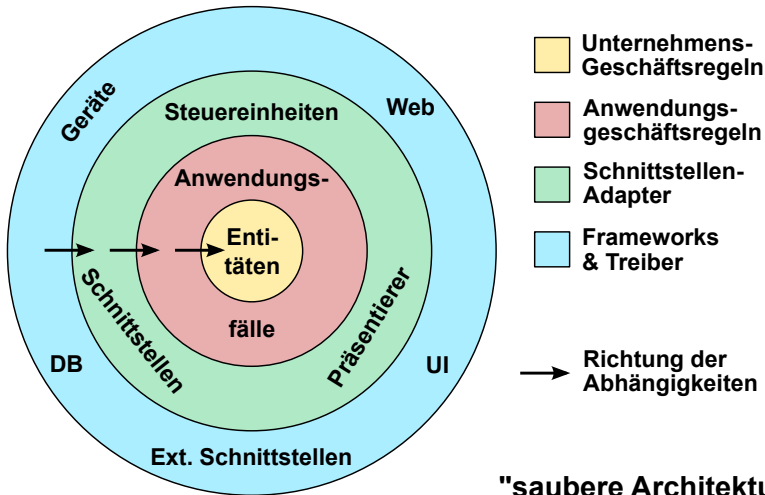
Beispiele plattform- und sprachunabhängiger Formate

Zum Umgang mit Daten und Metadaten

Bedeutung im Gesamtkontext einer Auswertungssoftware

Datenformate sind peripher

Auswertungsroutinen hängen nicht von Datenformaten ab.



verändert nach Robert C. Martin: Clean Architecture. Prentice Hall, Boston 2018, S. 203

Dependency-Inversion-Prinzip

Abstraktionen sollten nicht von Details abhängen.

- Speicherung von (Roh-)Daten ist ein Detail.
 - Das Datenformat ist für die Auswertungsroutine egal.
 - Die Organisation der Datenablage ist ebenfalls egal.
- Datenformate sind austauschbar.
 - Alles, was benötigt wird, ist eine Importroutine.
- Die Organisation der Datenablage ist austauschbar.
 - Zugriff auf Datensätze über einen abstrakten Schlüssel (ID) und eine (beliebig implementierbare) Zuordnungstabelle

- Jede Abstraktionsebene hat eigene Ansprüche.
 - Die Repräsentation der Daten in Auswertungsroutinen kann komplett anders sein als bei der Archivierung.
 - Jede Abstraktionsebene sollte die ihr entsprechende und bestmögliche Repräsentation wählen.
- Abhängigkeiten zeigen immer nach innen.
 - Datenformate sollten die durch sie implizierten Strukturen *nicht* nach innen durchreichen.
 - Auswertungsroutinen wissen nichts von der Datenablage.
 - Der Austausch einer Verzeichnishierarchie durch eine Datenbank ist problemlos möglich.

☛ Flexibilität, Modularität, Wiederverwendbarkeit



- 🔑 Formate betreffen nicht nur Roh- und verarbeitete Daten, sondern auch Metadaten, Dokumentation, Abbildungen.
- 🔑 Datenformate sollten über Jahrzehnte lesbar, plattformunabhängig, quelloffen und dokumentiert sein.
- 🔑 Daten über Jahrzehnte lesbar zu archivieren, ist nicht nur eine Frage der Formate, sondern auch der Organisation.
- 🔑 Rohdaten sollten immer (im Originalformat) archiviert und vor (ungewollter) Veränderung geschützt werden.
- 🔑 Das konkrete Datenformat oder die Art der Datenlagerung ist für ein System zur Datenverarbeitung irrelevant.