

Woche 02 Daten: Was sind Daten?

Skript

Erarbeitet von
Ann-Kathrin Selker

Lernziele	1
Inhalt	1
Einstieg.....	2
Metrische und kategoriale Daten	2
Strukturierte und unstrukturierte Daten	3
Abschluss	6
Quellen	6
Weiterführendes Material.....	6
Disclaimer	6

Lernziele

- Daten anhand von Beispielen in metrische und kategoriale Daten einsortieren
- Daten anhand von Beispielen in strukturierte und unstrukturierte Daten einsortieren
- erklären, wie Bilder im Computer dargestellt werden
- ASCII-kodierten Text in Zeichen umwandeln

Inhalt

Einstieg

Daten begegnen uns überall im Leben. Jeden Tag werden Wetterdaten erhoben, unsere Smartphones verwenden mobile Daten, und laut Datenschutzgrundverordnung sollen unsere persönlichen Daten geschützt werden. Doch was sind Daten überhaupt?

Metrische und kategoriale Daten

In der Informatik handelt es sich bei Daten um Informationseinheiten. Daten können dabei in metrische und kategoriale Daten unterschieden werden. Als metrische oder auch numerische Daten bezeichnen wir vereinfacht gesagt Zahlen, ob natürliche Zahlen (z. B. 1, 2, 3, 4), ganze Zahlen (-2, -1, 0, 1) oder Kommazahlen (1,2 ; 3,4567).

Im Gegensatz zu metrischen Daten kann man mit kategorialen Daten nicht rechnen.

Kategoriale Daten sind entweder eine Kategorie, also z. B. die Fakultäten der HHU, oder eine Skala, zum Beispiel eine Bewertung von 1,0 bis 5,0 einer Klausur. Bei dieser Skala ist nur die Reihenfolge der Werte wichtig, nicht die Werte selber. Es ist also vollkommen egal, dass es in der Notenskala einen Sprung von 4,0 zu 5,0 gibt.

Auch die Zeichen eines Textes sind kategoriale Daten. Im Computer wird ein Zeichen in der Grundkodierung ASCII als Zahl zwischen 0 und 127 gespeichert. ASCII steht für American Standard Code for Information Interchange. Eine Kodierung ist hier eine Tabelle, die jedem Zeichen einen Zahlenwert zuordnet. Der zugeteilte Zahlenwert ist im Grunde bedeutungslos und muss nur konsistent für dasselbe Zeichen verwendet werden.

Dez	Hex	Okt		Dez	Hex	Okt		Dez	Hex	Okt		Dez	Hex	Okt	
0	0x00	000	NUL	32	0x20	040	SP	64	0x40	100	@	96	0x60	140	`
1	0x01	001	SOH	33	0x21	041	!	65	0x41	101	A	97	0x61	141	a
2	0x02	002	STX	34	0x22	042	"	66	0x42	102	B	98	0x62	142	b
3	0x03	003	ETX	35	0x23	043	#	67	0x43	103	C	99	0x63	143	c
4	0x04	004	EOT	36	0x24	044	\$	68	0x44	104	D	100	0x64	144	d
5	0x05	005	ENQ	37	0x25	045	%	69	0x45	105	E	101	0x65	145	e
6	0x06	006	ACK	38	0x26	046	&	70	0x46	106	F	102	0x66	146	f
7	0x07	007	BEL	39	0x27	047	'	71	0x47	107	G	103	0x67	147	g
8	0x08	010	BS	40	0x28	050	(72	0x48	110	H	104	0x68	150	h
9	0x09	011	TAB	41	0x29	051)	73	0x49	111	I	105	0x69	151	i
10	0x0A	012	LF	42	0x2A	052	*	74	0x4A	112	J	106	0x6A	152	j
11	0x0B	013	VT	43	0x2B	053	+	75	0x4B	113	K	107	0x6B	153	k
12	0x0C	014	FF	44	0x2C	054	,	76	0x4C	114	L	108	0x6C	154	l
13	0x0D	015	CR	45	0x2D	055	-	77	0x4D	115	M	109	0x6D	155	m
14	0x0E	016	SO	46	0x2E	056	.	78	0x4E	116	N	110	0x6E	156	n
15	0x0F	017	SI	47	0x2F	057	/	79	0x4F	117	O	111	0x6F	157	o
16	0x10	020	DLE	48	0x30	060	0	80	0x50	120	P	112	0x70	160	p
17	0x11	021	DC1	49	0x31	061	1	81	0x51	121	Q	113	0x71	161	q
18	0x12	022	DC2	50	0x32	062	2	82	0x52	122	R	114	0x72	162	r
19	0x13	023	DC3	51	0x33	063	3	83	0x53	123	S	115	0x73	163	s
20	0x14	024	DC4	52	0x34	064	4	84	0x54	124	T	116	0x74	164	t
21	0x15	025	NAK	53	0x35	065	5	85	0x55	125	U	117	0x75	165	u
22	0x16	026	SYN	54	0x36	066	6	86	0x56	126	V	118	0x76	166	v
23	0x17	027	ETB	55	0x37	067	7	87	0x57	127	W	119	0x77	167	w
24	0x18	030	CAN	56	0x38	070	8	88	0x58	130	X	120	0x78	170	x
25	0x19	031	EM	57	0x39	071	9	89	0x59	131	Y	121	0x79	171	y
26	0x1A	032	SUB	58	0x3A	072	:	90	0x5A	132	Z	122	0x7A	172	z
27	0x1B	033	ESC	59	0x3B	073	;	91	0x5B	133	[123	0x7B	173	{
28	0x1C	034	FS	60	0x3C	074	<	92	0x5C	134	\	124	0x7C	174	
29	0x1D	035	GS	61	0x3D	075	=	93	0x5D	135]	125	0x7D	175	}
30	0x1E	036	RS	62	0x3E	076	>	94	0x5E	136	^	126	0x7E	176	~
31	0x1F	037	US	63	0x3F	077	?	95	0x5F	137	_	127	0x7F	177	DEL

Einblendung ASCII-Tabelle (Quelle [1])

Die ASCII-Kodierung enthält die lateinischen Groß- und Kleinbuchstaben, Ziffern und Satzzeichen. Dabei wird ein kleines a durch die Zahl 97 und ein großes A durch die Zahl 65 repräsentiert. Die Ziffer 7 wird mit einer 55 dargestellt und ein Punkt hat die Zahl 46. Die ASCII-Tabelle enthält aber auch nicht darstellbare Zeichen, genannt Steuerungszeichen, wie das Tabulatorzeichen mit der 9. Die Sonderzeichen einiger Sprachen, wie zum Beispiel die deutschen Umlaute, sind nicht enthalten. Heutzutage wird die erweiterte Kodierung Unicode benutzt, die zwar mehr Speicherplatz pro Zeichen verbraucht, aber dafür auch sämtliche Sonderzeichen aller Sprachen und Emojis darstellen kann. Inzwischen enthält Unicode etwa 145.000 Zeichen. Die ASCII-Tabelle ist dabei integriert, sodass sich der Zahlenwert für z. B. ein kleines a nicht geändert hat.

Durch Aneinanderreihung einzelner Zeichen entstehen sogenannte Strings. Da die Zeichen eines Strings auch Leer- und Satzzeichen beinhalten, sind Strings daher nicht nur einzelne Wörter wie "KI", sondern auch ganze Texte wie „Herzlich willkommen in der Vorlesung ,KI für alle‘! Schön, dass ihr da seid.“

Strukturierte und unstrukturierte Daten

In den späteren Anwendungen wollen wir aber auch komplexere Daten abbilden können. Wir unterscheiden zwischen strukturierten und unstrukturierten Daten. Strukturierte Daten sind in einem festgelegten Format erstellt. Ein Beispiel dafür sind Daten in Tabellen. Tabellen bestehen aus Zeilen und Spalten, bei denen jede Zeile für ein Datenobjekt steht und jede Spalte für Eigenschaften dieses Objekts, sogenannte Features. Features können dabei verschiedene Datentypen haben.

Matrikelnummer	Name	Vorname	Studienfach
1234567	Fischer	Anna	Philosophie
1289365	Ostrowska	Waleria	Mathematik
1314532	Hietanen	Toni	Germanistik
1349823	Osterhagen	Maximilian	Computerlinguistik
1453211	Glöckner	Marina	Pharmazie
1444478	Wei	Fen	Rechtswissenschaften
1521223	Rana	Chandra	Pharmazie
1523898	Yuryeva	Nora	Kunstgeschichte
1554599	Page	Sean	Informatik
1577266	Kuster	Martin	Medizin
1612923	Schulze	Anja	Betriebswirtschaftslehre
1682789	Yamauchi	Hisako	Wirtschaftschemie
1725521	Chia	Qiang	Medizin
1921456	Preciado Mena	Pirro	Informatik

Einblendung Tabelle Studierendendaten

In dieser Tabelle beinhaltet jede Zeile die Daten eines Studierenden, in der ersten Zeile z. B. von Anna Fischer mit Matrikelnummer 1234567 und Studienfach Philosophie. Hier werden Matrikelnummern als natürliche Zahlen dargestellt, Namen als Strings, und das Feature „Studienfach“ beinhaltet Kategorien. Dieses Schema ist für alle Studierenden gleich, d. h. der erste Wert ist immer eine Matrikelnummer etc.

Unter unstrukturierten Daten verstehen wir Daten, die keine normalisierte, feste Struktur aufweisen. Darunter fallen Bilder, Texte und ähnliche Daten. Damit diese verarbeitet werden können, muss ihnen erst noch eine Art Struktur gegeben werden. Bei Bildern ist dies zum Beispiel durch Raster möglich. Dabei wird ein Bild in Zeilen und Spalten aufgeteilt. Jedes dadurch entstandene Kästchen wird als Pixel bezeichnet und je nach Anwendung durch einen oder mehrere Zahlenwerte repräsentiert. Bei sogenannten Grayscale-Bildern (also Graustufenbildern) wird ein Pixel durch eine Zahl zwischen 0 und 255 dargestellt, die die Intensität des Lichts angibt. Dabei steht 0 für kein Licht, also schwarz, und 255 für nur Licht, also weiß. Die Zahlen dazwischen sind Abstufungen von grau. Da es insgesamt 256 mögliche Pixelwerte gibt, können auch nur 256 verschiedene Graustufen angezeigt werden.



Einblendung Schmetterlingsbild (Quelle [2])

Das ursprüngliche Schmetterlingsbild besitzt über 300.000 Pixel. Wenn wir jetzt in der Kopfgegend des Schmetterlings in das Bild herein zoomen, erhalten wir diesen Teilabschnitt mit 5 Zeilen und 5 Spalten und damit insgesamt 25 Pixeln,



Einblendung gezoomtes Bild + Tabelle mit Pixelwerten

bei dem man klar das Raster erkennen kann. Zum Beispiel hat das Pixel in der linken oberen Ecke den Wert 250, das daneben den Wert 176 usw.

Um Farbbilder darzustellen, reichen Lichtintensitäten nicht aus. Bei RGB-Bildern wird jedes Pixel durch drei Zahlenwerte zwischen 0 und 255 dargestellt: Dabei steht der erste Wert für rot (R), der zweite für grün (G), und der dritte für blau (B). Diese drei Farben bezeichnen wir auch als Grundfarben. Je niedriger der Zahlenwert für eine Grundfarbe, desto weniger Anteile hat diese an der resultierenden Farbe.

R	G	B	Farbe
0	0	255	Blue
255	128	0	Orange
0	0	0	Black
255	255	255	White

Einblendung Tabelle mit Farbwerten

Zum Beispiel hat reines Blau den Wert (0, 0, 255), reines Orange den Wert (255, 128, 0), schwarz den Wert (0, 0, 0) und weiß den Wert (255, 255, 255). Da jede Kombination von Rot-, Grün- und Blauwerten eine neue Farbe angibt, können insgesamt $256 * 256 * 256$ verschiedene Farben angezeigt werden, also mehr als 16 Millionen.

Abschluss

In diesem Video habt ihr gelernt, was wir unter metrischen und kategorialen Daten verstehen, was wir unter unstrukturierten und strukturierten Daten verstehen, wie Textkodierung funktioniert und wie Bilder im Computer dargestellt werden können.

Quellen

Quelle [1] Bild „ASCII“ von Hador, Hervorhebungen hinzugefügt.
<https://de.wikiversity.org/wiki/Datei:ASCII.pdf>
Lizenz CC BY-SA 3.0

Quelle [2] pixabay.com

Weiterführendes Material

<http://christianherta.de/lehre/infoDarstellung.pdf>

Disclaimer

Transkript zu dem Video „Woche 02 Daten: Was sind Daten?“, Ann-Kathrin Selker. Dieses Transkript wurde im Rahmen des Projekts ai4all des Heine Center for Artificial Intelligence and Data Science (HeiCAD) an der Heinrich-Heine-Universität Düsseldorf unter der Creative Commons Lizenz [CC-BY](https://creativecommons.org/licenses/by/4.0/) 4.0 veröffentlicht. Ausgenommen von der Lizenz sind die verwendeten Logos, alle in den Quellen ausgewiesenen Fremdmaterialien sowie alle als Quellen gekennzeichneten Elemente.