

# Project: Transcribing Data

Daniel T. Kaplan

2023-01-12

In order to use data wrangling and graphics techniques, the data frames themselves must be properly formed. Unfortunately, most people are unfamiliar with the basic principles behind data-frame organization. Consequently, even machine-readable documents such as spreadsheets look organized but often contain flaws that prevent their full use as data.

This project is meant to help you internalize principles of proper data organization so that they become second nature to you. Another goal of the project is to develop your skills in recognizing spreadsheet pitfalls that interfere with wrangling. Finally, the project will introduce you to the concept of a *relational database*. We will not be making much use of relational databases in Math 300, but anyone literate with data needs to be aware of this critical and ubiquitous tool and why it is so important to working with data.

## The project context: the US Census

As you may know, Article 1 Section 2 of the US Constitution requires that an “actual Enumeration” of all residents of the US be made every ten years.

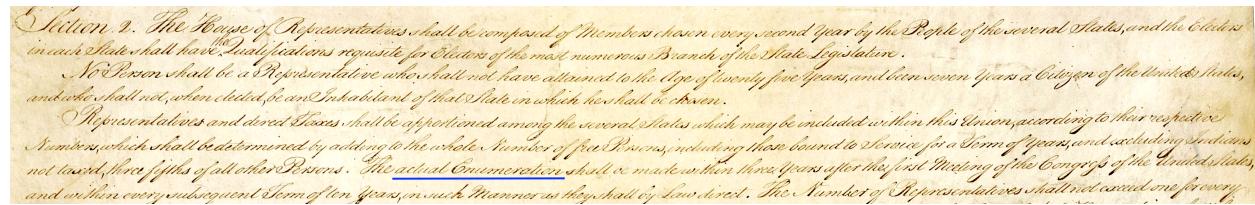


Figure 1: ARTICLE 1, Section 2 of the US Constitution, requiring a census every 10 years. The section also contains the infamous "three-fifths" provision for counting enslaved people.

A basic input to the census enumeration is the “population schedule,” which lists individual persons as rows in a spreadsheet. By law, the Census Bureau must keep private the information on individuals, and is not even able to share it with other government agencies. Only the summary tabulations made via data wrangling can be published.

However, 72 years after each census the population schedules can be released. In 2022, the population schedules for the 1950 census were released. We are going to use population schedules from the 1940 census, which are more accessible than the newly released 1950s sheets.

@fig-pop-schedule-small shows a population schedule from North Nevada Street in Colorado Springs, CO. This is a spreadsheet in the original sense of the word: a broad sheet of paper used for accounting.

## Spreadsheet structure

We call a document like that in @fig-pop-schedule-small a “spreadsheet” because it is not yet organized appropriately for a data frame. (This is understandable, because the concept of relational databases originates in the 1970s. Today’s Census Bureau, like every data organization, uses relational spreadsheets.)

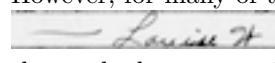
@fig-pop-schedule-portion zooms in on the upper-left side of the sheet to make the structure more apparent.

State Colorado Incorporated place Broadmoor City of city.....  
 Township or other division of county Block No. 1  
 County Boulder Institution (Check if institution and list in which section)

DEPARTMENT OF COMMERCE—BUREAU OF THE CENSUS U.D. No. 2 E.D. No. 507  
 SIXTEENTH CENSUS OF THE UNITED STATES: 1940  
 POPULATION SCHEDULE

LOCATION	REGISTERED DATE	NAME	RELATION	FEDERAL CENSUS NUMBER	EDUCATION	PLACE OF BIRTH	RESIDENT APRIL 1, 1940										PEOPLES IN YEARS OLD AND OVER—EMPLOYMENT STATUS																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																					
							1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580	581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700	701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724	725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748	749	750	751	752	753	754	755	756	757	758	759	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	760	761	762	763	764	765	766	767	768	769	770	771	772	773	774	775	776	777	778	779	770	771	772	773	774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797	798	799	790	791	792	793	794	795	796	797	798	799	800	801	802	803	804	805	806	807	808	809	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	810	811	812	813	814	815	816	817	818	819	820	821	822	823	824	825	826	827	828	829	820	821	822	823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846	847	848	849	840	841	842	843	844	845	846	847	848	849	850	851	852	853	854	855	856	857	858	859	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	860	861	862	863	864	865	866	867	868	869	870	871	872	873	874	875	876	877	878	879	870	871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894	895	896	897	898	899	890	891	892	893	894	895	896	897	898	899	900	901	902	903	904	905	906	907	908	909	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919	910	911	912	913	914	915	916	917	918	919	920	921	922	923	924	925	926	927	928	929	920	921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944	945	946	947	948	949	940	941	942	943	944	945	946	947	948	949	950	951	952	953	954	955	956	957	958	959	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969	960	961	962	963	964	965	966	967	968	969	970	971	972	973	974	975	976	977	978	979	970	971	972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995	996	997	998	999	990	991	992	993	994	995	996	997	998	999	1000	1001	1002	1003	1004	1005	1006	1007	1008	1009	1000	1001	1002	1003	1004	1005	1006	1007	1008	1009	1010	1011	1012	1013	1014	1015	1016	1017	1018	1019	1010	1011	1012	1013	1014	1015	1016	1017	1018	1019	1020	1021	1022	1023	1024	1025	1026	1027	1028	1029	1020	1021	1022	1023	1024	1025	1026	1027	1028	1029	1030	1031	1032	1033	1034	1035	1036	1037	1038	1039	1030	1031	1032	1033	1034	1035	1036	1037	1038	1039	1040	1041	1042	1043	1044	1045	1046	1047	1048	1049	1040	1041	1042	1043	1044	1045	1046	1047	1048	1049	1050	1051	1052	1053	1054	1055	1056	1057	1058	1059	1050	1051	1052	1053	1054	1055	1056

A key aspect of data-frame organization is the idea of a “unit of observation”: the kind of thing represented by a row. Looking carefully at @fig-pop-schedule-portion, it is evident that each row corresponds to an individual person. For example, line 42 is about Louise W. Service, a 38-year-old married woman.

However, for many of the individuals, there are blank cells. Many of the names start with a long dash, as in . The point of such elisions is to avoid duplicating entries. For instance, Louise W. shares the last name of her husband, William C., as well as the same address, 1415 N. Nevada Street.

Avoiding unnecessary duplication is a fundamental principle in the organization of data.<sup>1</sup> In 1940, avoiding duplication reduced the amount of writing needed for data entry. In the 2020s, avoiding unnecessary duplication is still important as a way of avoiding possible inconsistencies and making it clear that two or more cells must be identical.

The use of blanks or dashes in the spreadsheet is a common-sense way of indicating that information from one row is shared by an adjacent row. But it is not consistent with the proper organization of data frames and databases.

A fundamental, though unexpected and non-intuitive principle of data-frame organization is that the *order of the rows* does not matter. Instead, all information about relationships is represented by the data itself, regardless of order. This principle has important consequences for the organization of data and a major part of the motivation for relational databases.

A modern way to interpret what’s going on in the population schedule spreadsheet is that there are actually *two different units of observation* involved: i) the individual person and ii) the household. Since a data frame can have only a single unit of observation, representing the spreadsheet in a modern format will require two different data frames: one for the individual persons and another for the individual households.

The requirement that every data frame have its unit of observation provides many advantages. For instance, suppose it was realized after the data entry that the rent on the household’s home is incorrect. Fixing it can be accomplished by changing a single number in a single data frame. Or, consider what happens when it is discovered that a person was missed in the enumeration of a household. Adding this new person can be accomplished by appending a new row to the data frame containing individuals.

## Tasks for the project

Note that you **do not need** to do any work in R for this project.

You will need two URLs and the `pop_schedule_ID` from this roster.

1. Population schedule URL: links to an image of a population schedule. Download this to your laptop.
2. Data entry URL: links to a Google sheet. You will edit the Google sheet from a browser on your laptop (or any other machine.) The Google sheet contains two tabs, for the *household* and the *persons* data frames.
  - Fill in the *household* data frame first. You will need it for the *persons* data frame.
  - Some of the population schedules are continued from a previous sheet. Start your transcription with the first *complete* household.

## Data-entry tips

1. When entering data for a categorical variable, decide on a set of levels before you start entering the data.
2. There is a `household_ID` variable in both the *households* and *persons* data frames. Similarly, there is a `pop_schedule_ID`. These must be absolutely consistent across the two data frames.
3. Each population schedule contains a column (#3) labeled “Number of household in order of visitation.” Use this as the `household_ID` in your data frames. It is unique within the population schedule. Together with the `pop_schedule_ID` it is unique universally, since every population schedule has a unique `pop_schedule_ID`.

---

<sup>1</sup>An example of *necessary* duplication appears in lines 52 and 53, where two people happen to have the same name.

4. In the *persons* spreadsheet, we have collapsed columns 21-25 into one multi-level categorical variable named **workplace**. Read the column 21-25 headers to determine a set of levels that will properly encode the information in those columns.

#### Notes in draft for instructors

1. Google sheet for distributing population schedules and data-entry spreadsheets: <https://docs.google.com/spreadsheets/d/1ZxbiNCKrGeCAqYTd1R9QLSRBLia9HBk4k8DxHxKqiok/edit?usp=sharing>
2. An example data-entry spreadsheet: <https://docs.google.com/spreadsheets/d/1nUMP1VQfTmhQy1qlsKGSlJ825GpnyXto8M8tkUhtOUU/edit?usp=sharing>