Data processing using UNIX

Chris Reynolds | Info 287: FA.24

How many unique countries does the data contain? 24

cut -d ',' -f1 covid19_excess_deaths_large.csv | sort | uniq | wc -l

There are 4 parts to this code delineate via a pipe (|):

- via the "cut" command, we tell the system to look at only the first column of data in the Covid data set csv, which contains the information regarding which country the data is coming from;
- the system sorts the country data alphabetically;
- the system looks for unique items within the set, because it is sorted alphabetically it only needs the "unqi" command to compare each instance to the instances above and blow the one it is looking at, at a time;
- we are asking for a line count;
- we manually subtract 1 from the number presented as the header of "country" is in this count.

	M2 — -zsh — 95×63
christinereynolds@Ezra M2 % cut -d l 25	',' -f1 covid19_excess_deaths_large.csv sort uniq wc -
christinereynolds@Ezra M2 % 🗌	

How many unique regions in Italy does the data contain? 22

cut -d "," -f1,2 covid19_excess_deaths_large.csv | grep "Italy" | sort -t "," -k2 | uniq | wc -l

- Here we are telling the system to "cut" and only look at the first 2 columns of data in the file;
- "grep" tells the system to look for instances of the word Italy anywhere within the data (because of the nature of this data set, I know it will only be seeing/counting in the first column so that doesn't need to be specified);
- this time the system sorts the 2nd column alphabetically;
- looking for unique instances within the regions column;
- again we're asking for a line count;
- additionally we need to subtract 1 from the number presented as the header "region" isn't apart of what we want to count.



Find all records which pertain to the "Lombardia" in Italy.

grep "Lombardia" covid19_excess_deaths_large.csv

Written this way, the command retrieves all rows of data that contain Lombardia in any cell and writes them in the terminal.

	🚞 M2 — -zsh — 95×63	
[christinereynol	ds@Ezra M2 % grep "Lombardia" covid19_excess_deaths_large.csv	
Italy,Lombardia	region,week,2015,1,2,2015-01-15,2399,2399,0,0	
Italy,Lombardia	region,week,2015,1,3,2015-01-22,2348,2348,0,0	
Italy,Lombardia	region,week,2015,1,4,2015-01-29,2316,2338,-22,-0.9409751924721945	
Italy,Lombardia	region,week,2015,2,5,2015-02-05,2380,2361,19,0.8047437526471839	
Italy,Lombardia	region,week,2015,2,6,2015-02-12,2461,2286,175,7.655293088363948	
Italy,Lombardia	region,week,2015,2,7,2015-02-19,2384,2142,242,11.297852474323065	
Italy,Lombardia	region,week,2015,2,8,2015-02-26,2199,2135,64,2.997658079625282	
Italy,Lombardia	region,week,2015,3,9,2015-03-05,2145,2145,0,0	
Italy,Lombardia	region,week,2015,3,10,2015-03-12,2034,2034,0,0	
Italy,Lombardia	region,week,2015,3,11,2015-03-19,1984,1984,0,0	
Italy,Lombardia	region,week,2015,3,12,2015-03-26,2012,1921,91,4.737116085372193	
Italy,Lombardia	region,week,2015,4,13,2015-04-02,1963,1890,73,3.8624338624338606	
Italy,Lombardia	region,week,2015,4,14,2015-04-09,1956,1956,0,0	
Italy,Lombardia	region,week,2015,4,15,2015-04-16,1915,1845,70,3.794037940379397	
Italy,Lombardia	region,week,2015,4,16,2015-04-23,1895,1854,41,2.2114347357065753	
Italy,Lombardia	region,week,2015,4,17,2015-04-30,1729,1738,-9,-0.5178365937859581	
Italy,Lombardia	region,week,2016,1,2,2016-01-15,2005,2399,-394,-16.42350979574823	
Italy,Lombardia	region,week,2016,1,3,2016-01-22,2112,2348,-236,-10.051107325383299	
Italy,Lombardia	region,week,2016,1,4,2016-01-29,2013,2338,-325,-13.900769888793846	
Italy,Lombardia	region,week,2016,2,5,2016-02-05,1945,2361,-416,-17.619652689538327	
Italy,Lombardia	region,week,2016,2,6,2016-02-12,2018,2286,-268,-11.723534558180233	
Italy Lombardia	region_week.2016.2.7.2016-02-19.1942.21422009.337068160597568	

Copy the records obtained (in step 3) to a new file where the records are organized by the number of deaths in descending order.

head -1 covid19_excess_deaths_large.csv > Lombardia_deaths.csv

grep "Lombardia" covid19_excess_deaths_large.csv | sort -t "," -k8 -nr >> Lombardia_deaths.csv

This is done in two separate code chunks.

- The system sets the new file "Lombardia_deaths.csv" with only the first row of information from the Covid file in it, essentially setting up a table header;
- here we're using the same "grep" code to find the instances of "Lombardia";
- this time we're asking the system to sort the 8th column (the number of deaths) in reverse numerical order (from largest to smallest);
- ">>" appends that data to the new csv.

```
M2 — -zsh — 95×63
Ichristinereynolds@Ezra M2 % head -1 covid19_excess_deaths_large.csv
gountry, region, period, year, month, week, date, deaths, expected_deaths, excess_deaths, excess_deaths_p ct
Ichristinereynolds@Ezra M2 % head -1 covid19_excess_deaths_large.csv > Lombardia_deaths.csv
christinereynolds@Ezra M2 % grep "Lombardia" covid19_excess_deaths_large.csv | sort -t "," -k8
I-nr >> Lombardia_deaths.csv
christinereynolds@Ezra M2 %
```

Lombardia deaths

country	region	period	year	month	week	date	deaths	expected_deaths	excess_deaths	excess_deaths_pct
Italv	Lombardia region	week	2020	3	12	2020-03-25	7830	1921	5909	307.600208224883
Italy	Lombardia region	week	2020	4	13	2020-04-01	6949	1890	5059	267 6719576719580
	Lombardia region	week	2020		11	0000 00 10	6105	1000	4101	007 7116005 49397
	Lombardia region	WEEK	2020			2020-03-18	6105	1964	4121	207.7110933463870
taly	Lombardia region	week	2020	4	14	2020-04-08	5269	1956	3313	169.3762781186090
taly	Lombardia region	week	2020	4	15	2020-04-15	4207	1845	2362	128.0216802168020
taly	Lombardia region	week	2020	3	10	2020-03-11	3875	2034	1841	90.5113077679449
taly	Lombardia region	week	2020	4	16	2020-04-22	3310	1854	1456	78.5329018338727
taly	Lombardia region	week	2017	1	2	2017-01-15	2914	2399	515	21.46727803251350
talv	Lombardia region	week	2017	1	3	2017-01-22	2722	2348	374	15 92844974446340
	Lombardia region	week	2017		47	2017-01-22	2722	2340	000	54 4597 479904040
Laiy	Lombardia region	week	2020	4	17	2020-04-29	2627	1738	889	51.150/4/9861910
taly	Lombardia region	week	2017	1	4	2017-01-29	2585	2338	247	10.56458511548330
taly	Lombardia region	week	2018	1	2	2018-01-15	2553	2399	154	6.41934139224678
taly	Lombardia region	week	2015	2	6	2015-02-12	2461	2286	175	7.6552930883639
taly	Lombardia region	week	2020	3	9	2020-03-04	2419	2145	274	12.77389277389280
talv	Lombardia region	week	2019	1	4	2019-01-29	2416	2338	78	3.336184773310520
	Lemberdie meien	week	0010			2010 02 05	0410	0.001		0.16010165194044
aly	Lombardia region	week	2019	2	5	2019-02-05	2412	2361	51	2.16010165184244
taly	Lombardia region	week	2018	1	3	2018-01-22	2407	2348	59	2.5127768313458
taly	Lombardia region	week	2015	1	2	2015-01-15	2399	2399	0	
taly	Lombardia region	week	2019	2	7	2019-02-19	2389	2142	247	11.531279178338
taly	Lombardia region	week	2019	2	6	2019-02-12	2386	2286	100	4.3744531933508
talv	Lombardia region	week	2015	2	7	2015-02-19	2384	2142	242	11.2978524743231
,	Lemberdie meien	week	0015	-		2015 02 05	0200	0.961	10	0.00474075064710
taly	Lombardia region	week	2015	2	5	2015-02-05	2380	2361	19	0.804743752647184
taly	Lombardia region	week	2017	2	5	2017-02-05	2361	2361	0	
taly	Lombardia region	week	2015	1	3	2015-01-22	2348	2348	0	
taly	Lombardia region	week	2018	1	4	2018-01-29	2338	2338	0	
taly	Lombardia region	week	2019	2	8	2019-02-26	2326	2135	191	8.9461358313817
taly	Lombardia region	Week	2015		4	2015-01-29	231F	2.30		-0.94097519247219
y	Lombardia		2013	-	-	2019 00 07	0000	2030	-22	.0 54100000000
ually	Lombardia region	week	2018	2	5	2018-02-05	2301	2361	-60	-2.5412960609911
aly	Lombardia region	week	2019	1	2	2019-01-15	2300	2399	-99	-4.12671946644436
aly	Lombardia region	week	2017	2	6	2017-02-12	2286	2286	0	
aly	Lombardia region	week	2019	3	9	2019-03-05	2269	2145	124	5.7808857808857
taly	Lombardia region	week	2019	1	3	2019-01-22	2267	2348	-81	-3.44974446337308
taly	Lombardia maion	Week	2015	-		2015-02-22	2100	0105	01	2 9976580706250
als:	Lombardia region	wee'	2010		d	2020-04-25	2199	2135		7 300 4067 4000
any	Lombardia region	week	2020	1	4	2020-01-29	2165	2338	-1/3	-1.399486/408041
aly	Lombardia region	week	2018	3	9	2018-03-05	2161	2145	16	0.7459207459207
aly	Lombardia region	week	2015	3	9	2015-03-05	2145	2145	0	
aly	Lombardia region	week	2018	2	7	2018-02-19	2142	2142	0	
taly	Lombardia region	week	2018	2	8	2018-02-26	2135	2135	0	
-	Lombardia region	week	2020	1	3	2020-01-22	2129	2348	-219	-9 3270868824531
uny .	Lombardia region	week	2010		10	0010 00 10	0110	2040	210	4 1007025102044
Laiy	Lombardia region	Week	2018	3	10	2010-03-12	2110	2034	04	4.1297935103244
taly	Lombardia region	week	2016	1	3	2016-01-22	2112	2348	-236	-10.0511073253833
taly	Lombardia region	week	2020	2	6	2020-02-12	2108	2286	-178	-7.7865266841644
taly	Lombardia region	week	2020	1	2	2020-01-15	2108	2399	-291	-12.1300541892455
taly	Lombardia region	week	2020	2	8	2020-02-26	2105	2135	-30	-1.40515222482436
halu	Lombardia majon	wook	2020	2	5	2020.02.05	2000	1961	- 262	11 0060027006612
	Lombardia region	week	2020			2020-02-03	2035	2301	-202	-11.0808827880012
aly	Lombardia region	week	2017	2		2017-02-19	2079	2142	-63	-2.9411/64/05882
taly	Lombardia region	week	2020	2	7	2020-02-19	2076	2142	-66	-3.0812324929972
taly	Lombardia region	week	2019	3	11	2019-03-19	2057	1984	73	3.6794354838709
aly	Lombardia region	week	2017	3	9	2017-03-05	2056	2145	-89	-4.1491841491841
talv	Lombardia region	week	2019	3	10	2019-03-12	2051	2034	17	0.83579154375614
-	Lombardia region	week	2018	2	6	2018-02-12	2051	2286	-235	-10 279965004374
	Londardia region	week	0045	-	40	2010 02 12	2001	2200	200	10.27000004074
aly	Lombardia region	week	2015	3	10	2015-03-12	2034	2034	0	
taly	Lombardia region	week	2016	2	6	2016-02-12	2018	2286	-268	-11.7235345581802
taly	Lombardia region	week	2016	1	4	2016-01-29	2013	2338	-325	-13.9007698887938
taly	Lombardia region	week	2015	3	12	2015-03-26	2012	1921	91	4.7371160853721
taly	Lombardia region	week	2016	1	2	2016-01-15	2005	2399	-394	-16.423509795748
alv	Lombardia region	week	2018	3	11	2018-03-19	2004	1984	20	1.0080645161290
	Lombardia region	week	2010			2010-03-18	2004	1304	20	0.4.47000000000000
Laiy	Lombardia region	Week	2019	4	14	2019-04-09	1990	1930	42	2.14/239203003/
taly	Lombardia region	week	2017	2	8	2017-02-26	1994	2135	-141	-6.6042154566744
taly	Lombardia region	week	2016	2	8	2016-02-26	1989	2135	-146	-6.8384074941452
taly	Lombardia region	week	2015	3	11	2015-03-19	1984	1984	0	
taly	Lombardia region	week	2018	4	14	2018-04-09	1969	1956	13	0.66462167689161
taly	Lombardia region	week	2015		19	2015-04-02	1962	1800	79	3.86243386243396
	Lombardia	weet.	2013			2015 01 07	1055	1080	-	
udiy	Lombardia region	week	2015	4	14	2010-04-09	1956	1956	0	
aly	Lombardia region	week	2019	3	12	2019-03-26	1953	1921	32	1.66579906298803
taly	Lombardia region	week	2016	3	10	2016-03-11	1950	2034	-84	-4.129793510324
taly	Lombardia region	week	2017	3	11	2017-03-19	1945	1984	-39	-1.96572580645162
taly	Lombardia region	week	2016	2	5	2016-02-05	1945	2361	-416	-17.6196526895383
taly	Lombardia region	week	2016	2	7	2016-02-19	1942	2142	-200	-9.3370681605975
taly	Lombardia region	Week	2019		16	2018-04-22	1930	1854	76	4,099244875942
y	Lombardia	weet.	2010	-	10	2019 02 07	1000	1034		
any .	Lunch	week	2018	3	12	2010-03-26	1921	1921	0	
aly	Lombardia region	week	2016	3	12	2016-03-25	1921	1921	0	
taly	Lombardia region	week	2018	4	13	2018-04-02	1917	1890	27	1.42857142857142
aly	Lombardia region	week	2016	3	9	2016-03-04	1916	2145	-229	-10.6759906759907
aly	Lombardia region	week	2015	4	15	2015-04-16	1915	1845	70	3.7940379403794
aly	Lombardia region	week	2015		16	2015-04-22	1895	1854		2.21143473570859
talv	Lombardia mai	work	2010	-	10	2010.04.02	1000	1034	41	
	Lombardia region	wedt	2019	4	13	2010-04-02	1000	1890	-	0.476100470407
ау	Lompardia region	week	2016	4	13	2016-04-01	1881	1890	-9	-0.47619047619048
aly	Lombardia region	week	2016	3	11	2016-03-18	1881	1984	-103	-5.19153225806451
aly	Lombardia region	week	2017	3	10	2017-03-12	1877	2034	-157	-7.7187807276302
aly	Lombardia region	week	2016	4	14	2016-04-08	1874	1956	-82	-4.1922290388548
aly	Lombardia region	week	2017	3	12	2017-03-26	1862	1921	-59	-3.07131702238418
		Weer	2010		10	2019-04-22	1954	1054	-	
alv	1 ombardia main	week	2019	4	10	2019-04-23	1004	1854	0	0.400/
aly	Lombardia region				15	2018-04-16	1847	1845	2	0.108401084010850
aly aly	Lombardia region	week	2018	4	10				0	
taly taly taly	Lombardia region Lombardia region	week week	2018 2017	4	15	2017-04-16	1845	1845	-	
aly aly aly aly	Lombardia region Lombardia region Lombardia region	week week week	2018 2017 2017	4 4 4	15	2017-04-16 2017-04-02	1845 1836	1845	-54	-2.8571428571428
taly taly taly taly	Lombardia region Lombardia region Lombardia region	week week week	2018 2017 2017 2017	4	15 13 17	2017-04-16 2017-04-02 2017-04-22	1845 1836	1845	-54	-2.8571428571428
taly taly taly taly taly	Lombardia region Lombardia region Lombardia region Lombardia region	week week week week	2018 2017 2017 2017	4 4 4 4	15 13 17	2017-04-16 2017-04-02 2017-04-30	1845 1836 1824	1845 1890 1738	-54	-2.8571428571428
taly taly taly taly taly taly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week	2018 2017 2017 2017 2016	4 4 4 4 4 4	15 13 17 16	2017-04-16 2017-04-02 2017-04-30 2016-04-22	1845 1836 1824 1792	1845 1890 1738 1854	-54 86 -62	-2.8571428571428 4.9482163406214 -3.34412081984897
taly taly taly taly taly taly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week	2018 2017 2017 2017 2016 2016	4 4 4 4 4 4 4 4 4	15 13 17 16 15	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15	1845 1836 1824 1792 1792	1845 1890 1738 1854 1845	-54 86 -62 -53	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.87262872628726
taly taly taly taly taly taly taly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week	2018 2017 2017 2017 2016 2016 2017	4 4 4 4 4 4 4 4	15 15 13 17 16 15 14	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09	1845 1836 1824 1792 1792 1780	1845 1890 1738 1854 1854 1845 1956	-54 -62 -53 -176	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.87262872628726 -8.9979550102249
taly taly taly taly taly taly taly taly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week week	2018 2017 2017 2017 2016 2016 2016 2017 2017	4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-23	1845 1836 1824 1792 1792 1780 1776	1845 1890 1738 1854 1854 1845 1956 1854	-54 -54 -62 -53 -176 -78	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.87262872628726 -8.9979550102249 -4.2071197411003
aly aly aly aly aly aly aly aly aly aly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week week week	2018 2017 2017 2017 2016 2016 2017 2017 2017	4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16 16 15	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-23 2019-04-16	1845 1836 1824 1792 1792 1780 1776 1775	1845 1890 1738 1854 1855 1956 1854 1854	-54 -54 -62 -53 -176 -78 -70	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.872628726287262 -8.9979550102249 -4.2071197411003 -3.7940379403704
aly aly aly aly aly aly aly aly aly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week week week	2018 2017 2017 2016 2016 2016 2017 2017 2017	4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16 15 15	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-23 2019-04-16	1845 1836 1824 1792 1792 1780 1776 1775	1845 1890 1738 1854 1855 1956 1854 1855	-54 -54 -62 -53 -176 -78 -70	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.872628726287268 -8.9979550102249 -4.2071197411003 -3.7940379403794
taly taly taly taly taly taly aly aly aly aly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week week week	2018 2017 2017 2016 2016 2017 2017 2017 2019 2019	4 4 4 4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16 15 14 16 15 17	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-23 2019-04-16 2019-04-16	1845 1836 1824 1792 1792 1780 1776 1775 1762	1845 1890 1738 1854 1855 1956 1854 1855 1738	-54 -54 -62 -53 -176 -78 -70 24	-2.857142857142857142857142857142857142857142857142857452162872628726287262872628726287262872628
aly aly aly aly aly aly aly aly aly aly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week week week	2018 2017 2017 2017 2016 2016 2017 2017 2019 2019 2018	4 4 4 4 4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16 15 14 16 15 17 17	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-23 2019-04-16 2019-04-30 2018-04-30	1845 1836 1824 1792 1792 1780 1776 1775 1762 1738	1845 1890 1738 1854 1854 1855 1956 1854 1854 1738 1738	-54 86 -62 -53 -176 -78 -70 24 0	-2.8571428571428 4.9482163406214 -3.34412081984887 -2.87262872628726 8.9979550102249 -4.2071197411003 -3.79403794037944 1.38089758342923
taly taly taly taly taly taly taly taly	Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region Lombardia region	week week week week week week week week	2018 2017 2017 2016 2016 2016 2017 2017 2019 2019 2018 2015	4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16 15 17 17 17	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-23 2019-04-16 2019-04-30 2018-04-30 2015-04-30	1845 1836 1824 1792 1792 1780 1776 1775 1762 1738 1729	1845 1890 1738 1854 1845 1956 1854 1855 1738 1738 1738	-54 -54 86 -62 -53 -176 -78 -70 -24 -70 -24 -9	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.87262872628726 -8.9979550102249 -4.2071197411003 -3.7940379403794 1.38089758342923 -0.51783659378595
taly taly taly taly taly taly taly taly	Lombardia region Lombardia region	week week week week week week week week	2018 2017 2017 2016 2016 2016 2017 2019 2019 2019 2018 2015 2016	4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4	15 15 13 17 16 15 14 16 15 14 16 15 17 17 17 17	2017-04-16 2017-04-02 2017-04-30 2016-04-22 2016-04-15 2017-04-09 2017-04-09 2017-04-23 2019-04-16 2019-04-30 2018-04-30 2015-04-30	1845 1836 1824 1792 1792 1780 1776 1775 1762 1738 1729 1728	1845 1890 1738 1854 1855 1956 1854 1855 1738 1738 1738	-54 -54 -62 -53 -176 -78 -70 -24 -0 -9 -9 -10	-2.8571428571428 4.9482163406214 -3.34412081984897 -2.87262872628726 -8.9979550102249 -4.2071197411003 -3.7940374403794 1.38089758342923 -0.51783659378595 -0.57537399309551

Analyzing COVID-19 Data using R

Chris Reynolds | Info 287: FA.24

Load the CSV file into an R dataframe; view the dataframe and display the category of each variable (column)

1	#set directory
2	<pre>setwd("/Users/christinereynolds/Desktop/Problem Solving/M5")</pre>
3	#read data from CSV
4	<pre>coviddata <- read.csv("covid19_excess_deaths_large_1.csv", header=TRUE, sep=","</pre>
5	#show table in RStudio
6	View(coviddata)
7	#write column titles and data types in console
8	str (coviddata)
9	#clean the data
10	coviddata2 <- na.omit(coviddata)

Explain:

- Start by setting up the directory & pointing the code to the file/data we're starting with
- Asking the system to show us all of the column titles and types of data represented in the columns to better understand the data we're working with
- Asking the system to clean the data by telling it to omit any rows with blank/missing data Visual:

> #write column ti	tles and	data types in console
> str (coviddata)		
'data.frame': 53	947 obs.	of 11 variables:
\$ country	: chr	"Italy" "Italy" "Italy" "Italy"
\$ region	: chr	"Bergamo city" "Bergamo city" "Bergamo city" "Bergamo city"
\$ period	: chr	"week" "week" "week"
\$ year	: int	2015 2015 2015 2015 2015 2015 2015 2015
\$ month	: int	1 1 1 2 2 2 2 3 3 3
\$ week	: int	2 3 4 5 6 7 8 9 10 11
\$ date	: chr	"2015-01-15" "2015-01-22" "2015-01-29" "2015-02-05"
\$ deaths	: num	39 31 41 31 38 35 24 46 30 25
<pre>\$ expected_deaths</pre>	: num	36 31 39 31 35 32 24 33 30 27
<pre>\$ excess_deaths</pre>	: num	3 0 2 0 3 3 0 13 0 -2
<pre>\$ excess_deaths_p</pre>	ct: num	8.33 0 5.13 0 8.57

Pick one country and use R to filter the data using that country

12	#filtering for Netherlands (data frame) Opt 1
13	df_Netherands <- coviddata2[coviddata2\$country == 'Netherlands',]
14	#filtering for Netherlands (view w/tidyverse) Opt 2
15	library(tidyverse)
16	coviddata2 %>%
17	filter(country == 'Netherlands')

Explain:

I have two options here as I was having trouble with what was shown in lecture, so the second option of using tidyverse was actually the one I was first successful with. Opt 1 basically reads:

- create data frame: df_Netherands f
- rom the coviddata2 (cleaned data set)
- where the column of country is Netherlands
- include all other columns [,]

Visual:

🗈 Covid_with.R 🗴 🚺 df_NL2020weeks 🗴 🚺 df_2020weeks 🗶 🚺 df_Netherands 🗙						herands 🗙	df_covid 🔉		
🗲 🔿 🛛 🛪 🖌 Filter								Q	
^	country	region	period	year	month	week	date	deaths	expected
2205	Netherlands	Netherlands	week	2010	1	2	2010-01-15	2830	
2206	Netherlands	Netherlands	week	2010	1	3	2010-01-22	2846	
2207	Netherlands	Netherlands	week	2010	1	4	2010-01-29	2865	
2208	Netherlands	Netherlands	week	2010	2	5	2010-02-05	2858	
2209	Netherlands	Netherlands	week	2010	2	6	2010-02-12	2831	
2210	Netherlands	Netherlands	week	2010	2	7	2010-02-19	2900	

Filter the dataframe for the year 2020

```
19 #filtering for the year 2020 (data frame) Opt 1
20 df_covid2020 <- coviddata2[coviddata2$year == 2020, ]
21 #filtering for the year 2020 Opt 2
22 coviddata2 %>%
23 filter(year == 2020)
```

Explain:

Again I first figured out the tidyverse before understanding the first option shown in class. Opt one reads:

- create data frame df_covid2020
- from the cleaned data set
- where the column year is 2020
- include all other columns

Visual:

vith.R	× 🚺 df_N	L2020weeks 🛛 📔	df_2020we	eks 🛛 📔	df_Nethera	ands ×	df_covid202	0×	
•	▶ / 🛋 🛉 🔻	Filter					Q		
^	country	region	period	year	month	week	date	deaths	expe
81	Italy	Bergamo city	week	2020	1	2	2020-01-15	33	
82	Italy	Bergamo city	week	2020	1	3	2020-01-22	29	
83	Italy	Bergamo city	week	2020	1	4	2020-01-29	24	
84	Italy	Bergamo city	week	2020	2	5	2020-02-05	26	
85	Italy	Bergamo city	week	2020	2	6	2020-02-12	27	

For the new dataframe (generated in the last question), sum the number of deaths for each week



Explain:

I realized here that these were supposed to be stacking on each other at this point, but figured I'd get a little extra practice in so I did both the larger set from 2020 and then do the data frames for the Netherlands. That second chunk of code reads:

- Create a data frame from the previous Netherlands dataset and filter it by year
- Then create another new set where the deaths are added together for the whole country by week

oviddata 🗙	Covid	l_with.R* ×	df_NL2020weeks ×	df_2020weeks ×	📕 df_Netherands × » 👝 🗖
(л 📔 🍸 Filt	er			Q
^	Weeks	SumDeaths			
1	2	3364			
2	3	3153			
3	4	3043			
4	5	3160			
5	6	3192			
6	7	3198			

For the country you have picked, in the year 2020, is there a relationship between the week number and number of deaths? Report correlation values for these and comment on them

Answer: According to linear regression, no. -0.01161864

Visual:



Explain:

There is a super weak negative correlation: as the weeks go on there is the slightest decrease in deaths. I tried out a Kendall and Pearson correlation matrix, there is a "stronger" correlation shown with Kendall (which is typically used for ordered data); I thought Peasrson might work as the scatter plot showed a bit of a curve, but isn't normalized. Visual:



Generate visualization for the correlation matrix obtained in the last step





Explain: Using the library GGally you can create a visualization of the correlation matrix. These are more interesting when there are

- more values to look at, this visualization
- -0.5 basically says there isn't a correlation between
- -1.0 the deaths and the weeks.

Perform regression using Week as the predictor and Number of deaths as the response/outcome variable. Show the regression graph with the regression line



Explain:

I found this interesting as this is essentially the same visual that I got from the scatter plot only this time it came with the correlation line. I tried out both codes given to us (linear and curved). I can clearly see why the curved version could be overfitting.



Write the line equation

Answer: SumDeath = -1.251*Weeks+3425.732

```
55 #line equation
56 lm(SumDeaths~Weeks, df_NL2020weeks)
```

Explain:

Lm extracts the coefficients from the data frame in $y \sim x$ form because we're using x (the week) to predict y (the number of deaths in the country)

Visual

```
Call:
lm(formula = SumDeaths ~ Weeks, data = df_NL2020weeks)
Coefficients:
(Intercept) Weeks
3425.732 -1.251
```

Analyzing COVID-19 Data using Python

Chris Reynolds | Info 287: FA.24

Load the CSV file into a Python dataframe. Write appropriate commands to view the top and bottom few rows of the dataframe. Print the category of each variable.

Code:

```
df = pd.read_csv("/content/drive/MyDrive/Classes/PS/Colab
Notebooks/covid10_excess_deaths_large.csv")
print (df.head, df.tail)
print(df.dtypes)
```

Explain:

- First line of code tells the Google Colab environment where to find the csv and to read it in as a dataframe named "df".
- The second line asks for the first 5 and last 5 rows of the dataframe to be printed, so we can double check that the read is complete.
- The third line gives us more information about what variables are in the dataset and what type of data they are.

Visual

s [10] #load the data df = pd.read_csv("/content/drive/MyDr	rive/Classes/Prob	lem Solving/Colab	Notebooks/covid19_exces	s_deaths_large.csv")
s C	<pre>#print variable t print (df.dtypes)</pre>	ypes				
Ð	r country region period year month week date deaths expected_deaths excess_deaths excess_deaths_pct dtype: object	object object int64 float64 object float64 float64 float64 float64 float64 float64				
s [12] print (df.head, c	lf.tail)				
- FI	sbound method NDF 0 Ital 1 Ital 2 Ital 3 Ital 4 Ital 3 Ital 4 Ital 53942 Switzerlan 53943 Switzerlan 53945 Switzerlan 53946 Switzerlan 53946 Switzerlan 6 Sy 0 39.0	rame.head of y Bergamo city we y Bergamo city we y Bergamo city we y Bergamo city we y Bergamo city we is Switzerland we id Switzerland we	country eek 2015 1 eek 2015 1 eek 2015 2 eek 2015 2 eek 2015 2 eek 2020 5 eek 2020 5 eek 2020 5 eek 2020 6 eek 2020 6 eek 2020 6 eek 2020 6	region period 2.0 2015-01-15 3.0 2015-01-22 4.0 2015-01-29 5.0 2015-02-05 6.0 2015-02-12 20.0 2020-05-17 21.0 2020-05-31 23.0 2020-06-07 24.0 2020-06-14 4.0 2020-06-14 5.0 2020-06-14	year month week	date \

Pick one of the two regions: "New York City", "Lombardia region"

Use Python to filter the data using that region.

Answer: New York City

Code:

df_NYC = df[df["region"]].isin(["New York City"])]

Explain:

• This codes creates a new dataframe from the larger set by telling the computer to find the "region" column and return with all instances of "New York City"

Visual:

/ [13] df_NYC = df[df["region"].isin(["New York City"])]

```
Check the type and shape of df_NYC using:
```

Code:

```
print(type(df_NYC))
```

print(df_NYC.shape)

Report your findings: 'pandas.core.frame.DataFrame' - (378, 11)

Explain:

- We asked what type of data is our dataframe, and what dimensions is it.
- Python told us that it is a dataframe that was created with pandas, with 378 rows and 11 columns.

Visual:



For the new dataframe, df_NYC (generated in the last question), find out the correlation between the following pairs:

Code format:

dataframe_name["indep_verible"].corr(dataframe_name["depent_veriable"])

- year and excess_deaths
- year and expected_deaths
- month and excess_deaths
- month and expected_deaths
- week and excess_deaths
- week and expected_deaths
- expected_deaths and excess_deaths

Show your code & report correlation values for these pairs and comment on them:

```
df_NYC["year"].corr(df_NYC["excess_deaths"])
```

```
0.23396172917730804
```

```
df_NYC["year"].corr(df_NYC["expected_deaths"])
```

```
0.06610680808556683
```

```
df_NYC["month"].corr(df_NYC["excess_deaths"])
        -0.0869475912348743
df_NYC["month"].corr(df_NYC["expected_deaths"])
        -0.4350711961277345
df_NYC["week"].corr(df_NYC["excess_deaths"])
        -0.08899497249801491
df_NYC["week"].corr(df_NYC["expected_deaths"])
        -0.44655513477878506
df_NYC["expected_deaths"].corr(df_NYC["excess_deaths"])
        -0.011660100774849353
```

Comment:

• As the numbers get closer to 1 or -1 the correlation is stronger. The positive and negative indicates the direction of the correlation.

For the variable excess_deaths (and NOT expected_deaths), pick the two most strongly related variables (obtained in the previous step) and explain (in a few lines) why there might be a strong relationship between them

Year: 0.23396172917730804 Week: -0.08899497249801491

These have the strongest correlation with excess deaths. Year has the strongest relationship which is due to COVID 19. Both week and month seemed pretty weak to me as they were both less than .01; if I had to guess I'd say this could be due to how many data points there are at the week and month level over this many years.

Should excess_deaths be your predictor or outcome? Why?

Outcome as it is the dependent variable.

Based on your results obtained in the previous steps (Q5), pick one outcome and two predictor variables (Important: For regression, the variables that show the strongest correlation with the outcome are often picked as predictors - if all the predictors show weak correlations, pick those with higher correlation than others). State which variables you have picked as predictors (two variables) and the outcome (one variable).

Year: Predictor Variable Week: Predictor Variable Excess_deaths: Outcome Variable

Perform regression using one predictor (one of the two in the previous step) and one outcome variable. Show your code.

Code:

```
import numpy as np
from sklearn.linear_model import LinearRegression
x1 = np.array(df_NYC['year']).reshape((-1,1))
y1 = np.array(df_NYC['excess_deaths'])
```

```
mode4 = LinearRegression()
mode4.fit(x1, y1)
r2 = mode4.score(x1,y1)
intercept = mode4.intercept_
slope = mode4.coef_
```

pred4 = mode4.predict(x1)

Visual - see Appendix: Code

Show the regression graph with the regression line. Write the line equation.

Code:

```
import matplotlib.pyplot as plt
plt.scatter(df_NYC['year'], df_NYC['excess_deaths'])
plt.plot(x1, pred4, color='#EE5397', linewidth=3)
plt.grid(color = '#00C4D4', linestyle = '--', linewidth = 0.5)
plt.xlabel('Year')
plt.ylabel('Excess Deaths')
plt.title('# Excess Deaths by Year in New York City')
plt.show()
```

Explain:



• These lines use matplot Library and call for it via plt

• .scatter creates plots the data points from the dataframe as a scatter plot where x is years (independent variable) and y is excess deaths (dependent variable)

• The .plot created the the regression line

• I chose to add a dotted line grid to help with understanding the numbers being represented, and added x and y axis labels as well as a graph title.

Line Equation

Start coding or generate with AI.



from google.colab import drive
drive.mount('/content/drive')

→ Drive already mounted at /content/drive; to attempt to forcibly remount, call

import pandas as pd

#load the data
df = pd.read_csv("/content/drive/MyDrive/Classes/Problem Solving/Colab Notebooks/cc

print (df.head, df.tail)

$\overline{\mathbf{F}}$	<bound< th=""><th colspan="4">d method NDFrame.head of</th><th colspan="4">country</th><th colspan="3">region period year</th><th>I</th></bound<>	d method NDFrame.head of				country				region period year			I
	0	I	taly	Bergamo	city	week	2015		1	2.0	2015-01-15		
	1	I	taly	Bergamo	city	week	2015		1	3.0	2015-01-22		
	2	I	taly	Bergamo	city	week	2015		1	4.0	2015-01-29		
	3	I	taly	Bergamo	city	week	2015		2	5.0	2015-02-05		
	4	I	taly	Bergamo	city	week	2015		2	6.0	2015-02-12		
	•••		• • •		• • •	• • •	•••	• •	•	• • •			
	53942	Switzer	land	Switzer	land	week	2020		5	20.0	2020-05-17		
	53943	Switzer	land	Switzer	land	week	2020		5	21.0	2020-05-24		
	53944	Switzer	land	Switzer	land	week	2020		5	22.0	2020-05-31		
	53945	Switzer	land	Switzer	land	week	2020		6	23.0	2020-06-07		
	53946	Switzer	land	Switzer	land	week	2020		6	24.0	2020-06-14		
		deaths	avna	teab bata	hc ov	cess de	aathc			death	s nct		
	0		слреч	36	. 113 CZ	(cess_u		CALC	33_	_עכמנוו פיז	22222		
	1	31 0		31	0		0 0			0.0	00000		
	2	11 0		30	0		2 0			5 1	28205		
	2	31 0		31	0		0 0			0 0	20205		
	1	30 0		35			3 0			0.0 Q 5	71/20		
	4	20.0			.0		5.0			0.0	/1429		
	52042	1127 2		11/1	••		12 7			1 2	00701		
	52042	1127.J		1141		-	10 0			-1.2	00701		
	52011	1072 2		1190	0.0	-	00 0			-1.J	2204		
	52045	10/2.2		11/1		-	-90.0 6 7			-0.4	26707		
	53945	1140.0		1100	0.0		-0.2			-0.5	20/9/		
	53940	1088./		1111	0	-	-22.3			-2.0	07201		
	[53947	rows x	11 co [.]	lumns]> <	bound	method	NDFra	me.ta	il	of	COL	untry	
	0	I	taly	Bergamo	city	week	2015		1	2.0	2015-01-15		
	1	I	taly	Bergamo	city	week	2015		1	3.0	2015-01-22		
	2	I	taly	Bergamo	city	week	2015		1	4.0	2015-01-29		
	3	I	talv	Bergamo	citv	week	2015		2	5.0	2015-02-05		

	4	I	talý	Bergamo cit	ý week	2015	2	6.0	2015-02-12
	53942	Switzer	land	Switzerlan	d week	2020	5	20.0	2020-05-17
	530/3	Switzer	land	Switzerlan	d week	2020	5	2010	2020 05 1/
	52044	Switzer	land	Switzerlan	u week	2020	J F		
	53944	Switzer	tand	Switzertan	а week	2020	2	22.0	2020-05-31
	53945	Switzer	land	Switzerlan	d week	2020	6	23.0	2020-06-07
	53946	Switzer	land	Switzerlan	d week	2020	6	24.0	2020-06-14
		doathc	ovno	ctod doothc	020000	doathc	02000	doath	c nct
	0		exhe		excess_		EXCESS		5_µCL
	0	39.0		30.0		3.0		8.3	33333
	1	31.0		31.0		0.0		0.0	00000
	2	41.0		39.0		2.0		5.1	28205
	3	31.0		31.0		0.0		0.0	00000
	4	38.0		35.0		3.0		8.5	71429
		1127 2		1141 0		10 7		1 7	
	53942	112/.3		1141.0		-13./		-1.2	00/01
	53943	1180.0		1198.0		-18.0		-1.5	02504
	53944	1072.2		1171.0		-98.8		-8.4	37233
	53945	1148.8		1155.0		-6.2		-0.5	36797
	53946	1088.7		1111.0		-22.3		-2.0	07201
#ori	[53947	rows x	11 co	lumns]>					
#prin	t (df.d	able lyp tynes)	es						
ргтп		cypes,							
7	countr	V		obiect					
	region	,		object					
	neriod			object					
	periou			int64					
	year								
	montn			10164					
	week			float64					
	date			object					
	deaths			float64					
	expect	ed death	S	float64					
	excess			float64					
	evress	_deaths	nct	float64					
	dtvpe:	_ucutins_ object	ρει	1 cou co+					
	~- 7-------------	,							
			·	· · / [1151	V I C'I				
at_N	r c = dT	lat [reg	1TOU]	. TSIN(["New	TOPK LIT	A1)]			
prin	t(type(df_NYC))							
	+/ 4 E MV	(Cabana)							

pr print(df_NYC.shape)

<class 'pandas.core.frame.DataFrame'> (378, 11)

df_NYC["year"].corr(df_NYC["excess_deaths"])

0.23396172917730804

df_NYC["year"].corr(df_NYC["expected_deaths"])

0.06610680808556683

df_NYC["month"].corr(df_NYC["excess_deaths"])

-0.0869475912348743

df_NYC["month"].corr(df_NYC["expected_deaths"])

-0.4350711961277345

```
df_NYC["week"].corr(df_NYC["excess_deaths"])
```

-0.08899497249801491

df_NYC["week"].corr(df_NYC["expected_deaths"])

-0.44655513477878506

df_NYC["expected_deaths"].corr(df_NYC["excess_deaths"])

-0.011660100774849353

```
import numpy as np
from sklearn.linear_model import LinearRegression
```

import matplotlib.pyplot as plt

np.corrcoef(df_NYC['excess_deaths'], df_NYC['deaths'])

array([[1. , 0.99174825], [0.99174825, 1.]])

```
#Regression using scikitlearn - y = a +bx
#y is outcoming x is predictor a is y-intercept b is the slope
```

```
x1 = np.array(df_NYC['year']).reshape((-1,1))
y1 = np.array(df_NYC['excess_deaths'])
mode4 = LinearRegression()
mode4.fit(x1, y1)
r2 = mode4.score(x1,y1)
intercept = mode4.intercept_
```

```
slope = mode4.coef_
print("Intercept:", intercept)
print("Slop:", slope)
print("R2:", r2)
    Intercept: -119013.32539623718
    Slop: [59.05833713]
    R2: 0.05473809071963642
pred4 = mode4.predict(x1)

plt.scatter(df_NYC['year'], df_NYC['excess_deaths'])
plt.plot(x1, pred4, color='#EE5397', linewidth=3)
plt.grid(color = '#00C4D4', linestyle = '---', linewidth = 0.5)
plt.xlabel('Year')
plt.ylabel('Excess Deaths')
plt.title('# Excess Deaths by Year in New York City')
plt.show()
```

