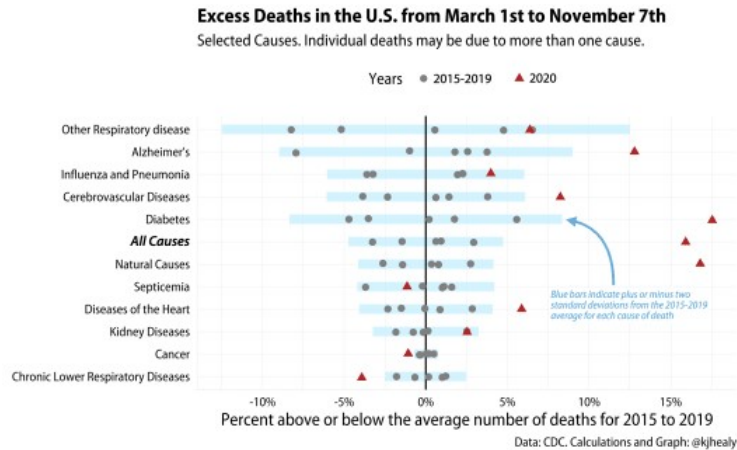


As I was [saying the other day](#), calculating excess deaths can be a tricky business, especially if your focus is on understanding counterfactuals like how many people died of some cause who would not have died due to some other competing risk over the period of interest. Moreover, even setting the counterfactuals aside, the whole business of accurately counting and classifying deaths on the scale of a country as large and variegated as the United States is an enormous challenge in itself. The CDC has been putting out its own [estimates of deaths due to COVID-19](#), and they make various efforts (such as weighting the estimates and so on) to account for delayed reporting and other issues.

I'll do something a little simpler here, but I think still useful. Using the [weekly counts for 2019-2020](#) and the [final counts for 2015-2019](#) we can examine selected causes for evidence of excess mortality beyond the baseline of expectations set by the past five years.

Here's the overall figure.



Excess deaths by Cause.

The idea here is to look at selected non-COVID causes of death (and also All-Cause mortality, i.e. everything) between March 1st and September 1st of this year, in comparison to the same causes between 2015 and 2019. We set the baseline as the mean number of deaths for each cause between 2015 and 2019. Then we calculate how far off each year is from that mean, for each cause. Obviously, it's not going to be the case that exactly the same number of people die of a given cause each year. But the U.S. is a large country, so there's a lot of stability from year to year, too. Most causes bounce around their average, but some are more variable than others. Cancer deaths, for instance, do not move around much from year to year. Others, such as Alzheimer's, and infectious diseases like the 'flu, are more variable. In the figure here, each gray dot is one of the years from 2015 to 2019, bouncing around that "No different from average" zero line. I've banded them with a blue bar showing twice the standard deviation from the mean. While not a super-formal test, anything outside two standard deviations from average is probably worth paying attention to here. We restrict ourselves to deaths that take place from March 1st to September 1st each year, as COVID wasn't causing fatalities in the US before March. The September 1st cutoff is mostly because data after this point (right now) gets quite noisy, with some states, such as Connecticut and North Carolina, not providing timely provisional counts by cause to the CDC.

I think the patterns here are interesting, and pretty clear. Most of the variation is just bouncing around within five percentage points of the mean, and all of the 2015-2019 years are within two standard deviations of their means. But 2020 is clearly different for many causes. All Cause mortality is everything, and is way up, as are deaths attributed to Diabetes, Alzheimer's, and Natural Causes. COVID isn't on the graph as a cause because there's no 2015-2019 baseline to compare it to, as it didn't exist. But it is a cause in the 2020 data.

Here's a closer look at the code and the tables the graph was produced from. As usual, I used R and the [covdata](#) package.

```
1
2 df_yr <- nchs_wdc %>%
3   filter(jurisdiction == "United States") %>%
4   filter(year > 2014,
5         week >= 9 &
6         week <= 34) %>%
7   group_by(cause, year) %>%
8   summarize(period_deaths = sum(n, na.rm = TRUE))
9
10
11 baseline_deaths <- nchs_wdc %>%
12   filter(jurisdiction == "United States") %>%
13   filter(year %in% c(2015:2019),
14         week >= 9 &
15         week <= 34) %>%
16   group_by(year, cause) %>%
17   summarize(total_n = sum(n, na.rm = TRUE)) %>%
18   group_by(cause) %>%
19   summarize(baseline = mean(total_n, na.rm = TRUE),
20             baseline_sd = sd(total_n, na.rm = TRUE))
21
22 df_excess <- left_join(df_yr, baseline_deaths) %>%
23   mutate(excess = period_deaths - baseline,
```

```

24     pct_excess = (excess / period_deaths)*100,
25     pct_sd = (baseline_sd/baseline)*100) %>%
26     rename(deaths = period_deaths)
27
28

```

Our excess deaths table (covering March 1st to September 1st) looks like this:

```

1
2
3 > df_excess
4 # A tibble: 82 x 8
5 # Groups:   cause [15]
6   cause      year deaths baseline baseline_sd excess pct_excess pct_sd
7
8 1 All Cause    2015 1318237 1359816      32421. -41579      -3.15    2.38
9 2 All Cause    2016 1338615 1359816      32421. -21201      -1.58    2.38
10 3 All Cause    2017 1367439 1359816      32421.   7623       0.557    2.38
11 4 All Cause    2018 1372444 1359816      32421.  12628       0.920    2.38
12 5 All Cause    2019 1402345 1359816      32421.  42529       3.03    2.38
13 6 All Cause    2020 1641133 1359816      32421. 281317      17.1    2.38
14 7 Alzheimer's  2015   51412   55788.    2684.  -4376.     -8.51    4.81
15 8 Alzheimer's  2016   55137   55788.    2684.   -651.     -1.18    4.81
16 9 Alzheimer's  2017   57448   55788.    2684.  1660.      2.89    4.81
17 10 Alzheimer's 2018   56828   55788.    2684.  1040.      1.83    4.81
18 # ... with 72 more rows
19

```

We also make a little tibble of the medians and standard deviations to make drawing the bars more convenient.

```

1
2
3 df_meds <- df_excess %>%
4   summarize(med = median(pct_excess))
5
6 df_sd <- df_excess %>%
7   filter(cause %nin% c("COVID-19 Underlying", "COVID-19 Multiple cause", "Other")) %>%
8   group_by(cause) %>%
9   slice(1) %>%
10  select(cause, pct_sd) %>%
11  mutate(lwr = -2*pct_sd,
12         upr = 2*pct_sd) %>%
13  left_join(df_meds)
14

```

The core of the plot is produced like this:

```

1
2 df_excess %>%
3   filter(cause %nin% c("COVID-19 Underlying", "COVID-19 Multiple cause", "Other")) %>%
4   mutate(yr_ind = ifelse(year == 2020, TRUE, FALSE)) %>%
5   ggplot(aes(x = pct_excess/100, y = reorder(cause, pct_excess, median), color = yr_ind, group =
6     year)) +
7     geom_linerange(data = df_sd, mapping = aes(xmin = lwr/100, xmax = upr/100, y = reorder(cause, med)),
8       color = "deepskyblue1", alpha = 0.4, inherit.aes = FALSE, size = 3) +
9     geom_vline(xintercept = 0, color = "black") +
10    geom_jitter(size = 2, position = position_jitter(height = 0.05)) +
11    scale_color_manual(values = c("gray50", "firebrick"),
12      labels = c("2015-2019", "2020")) +
13    scale_x_continuous(breaks = c(-10, -5, 0, 5, 10, 15, 20)/100, labels = scales::percent_format(
14      accuracy = 1)) +
15    labs(x = "Percent above or below the average number of deaths for 2015 to 2019",
16      y = NULL,
17      color = "Years",
18      title = "Excess Deaths in the U.S. from March 1st to September 1st",
19      subtitle = "Selected Causes, arranged by median excess deaths.",
20      caption = "Data: CDC. Calculations and Graph: @kjhealy")

```

COVID and All-Cause mortality

We can also take a look at the `df_excess` table to see what's happening with All-Cause mortality and COVID-19 specifically. We need to wrangle the table a little to get the estimates side by side.

```

1 start_week <- 9
2 end_week <- 34
3
4 df_yr <- nchs_wdc %>%
5   filter(jurisdiction == "United States") %>%
6   filter(year > 2014,
7         week >= start_week &
8         week <= end_week) %>%
9   group_by(jurisdiction, cause, year) %>%
10  summarize(period_deaths = sum(n, na.rm = TRUE))
11
12
13 baseline_deaths <- nchs_wdc %>%
14   filter(jurisdiction == "United States") %>%
15   filter(year %in% c(2015:2019),
16         week >= start_week &
17         week <= end_week) %>%
18   group_by(jurisdiction, year, cause) %>%
19   summarize(total_n = sum(n, na.rm = TRUE)) %>%
20   group_by(jurisdiction, cause) %>%
21   summarize(baseline = mean(total_n, na.rm = TRUE),
22         baseline_sd = sd(total_n, na.rm = TRUE))
23
24 df_excess <- left_join(df_yr, baseline_deaths) %>%
25   mutate(excess = period_deaths - baseline,
26         pct_excess = (excess / period_deaths)*100) %>%
27   rename(deaths = period_deaths)
28
29 excess_count <- df_excess %>%
30   filter(year == 2020 &
31         cause %in% c("All Cause", "COVID-19 Multiple cause"))
32
33 excess_table <- excess_count %>%
34   mutate(col_cause = janitor::make_clean_names(cause)) %>%
35   select(jurisdiction, col_cause, deaths:pct_excess) %>%
36   group_by(jurisdiction) %>%
37   select(-cause) %>%
38   pivot_wider(names_from = col_cause, values_from = deaths:pct_excess) %>%
39   select(-pct_excess_covid_19_multiple_cause, -excess_covid_19_multiple_cause,
40 -baseline_covid_19_multiple_cause,
41 -baseline_sd_covid_19_multiple_cause)
42
43 colnames(excess_table) <- c("jurisdiction", "all_cause", "covid", "baseline", "baseline_sd", "excess",
44 "pct_excess")
45
46
47 excess_table <- excess_table %>%
48   mutate(deficit = excess - covid,
49         pct_covid = (covid / all_cause) * 100,
50         pct_deficit = (deficit / all_cause) * 100) %>%
51   select(jurisdiction, all_cause, baseline, baseline_sd, excess, covid, deficit, everything())
52

```

Which (finally) gives us this:

```

1
2 excess_table
3
4 # A tibble: 1 x 10
5 # Groups:   jurisdiction [1]
6 jurisdiction all_cause baseline baseline_sd excess covid deficit pct_excess pct_covid pct_deficit
7
8 1 United States 1641133 1359816 32421. 281317 179303 102014 17.1 10.9 6.22
9
10

```

So in these data (remember, the numbers are updated regularly, we're looking at March 1 to September 1 only, and this is a rough-and-ready calculation), we have 1,641,133 All-Cause deaths in comparison to a baseline 2015-2019 average of 1,359,816. In this period the raw excess is

281,317 deaths. COVID-19 was listed as a cause of 179,303 of these, leaving a deficit—a remaining excess—of 102,014. Overall excess mortality from March 1st to September 1st is 17.1% above the baseline, with COVID-19 accounting for 10.9 of those percentage points, with a 6.22 percentage point excess distributed across other causes.

Some proportion of the COVID-19 deaths would have succumbed to other causes of death this year. Some proportion of the non-COVID excess deaths are directly or indirectly attributable to COVID. Directly, for example, by someone dying of COVID in a care home, but having the cause recorded as Alzheimer's or Natural Causes. Indirectly, for instance, by someone suffering a stroke or a heart attack but being reluctant or unable to seek treatment until it was too late. And COVID has also, weirdly, probably resulted in some lives saved as a result of, say, fewer car accidents as a consequence of lockdown. Parceling out these effects, or trying to, will be a job for demographers and public health people for some time to come. But the sheer size of the direct and indirect mortality shock due to COVID just seems undeniable, and my feeling is that it won't be made to disappear even if its indirect and counterfactual effects get chiseled out or shifted a little at the margins as better data comes in and better estimates become possible.