

MOTOR VEHICLE DEATHS IN THE U.S.

AN ANALYSIS OF SEASONALITY, DAYS OF
THE WEEK, AND AGE FROM 2005-2015

GROUP #3

AMANDA ENSTAD

BARBARA MACGREGOR

MATT RUSSELL

CHI TRAN

WHY THIS TOPIC?

Detailed information,
Well-Structured

Significantly Large Dataset

Easily Accessible (Kaggle)

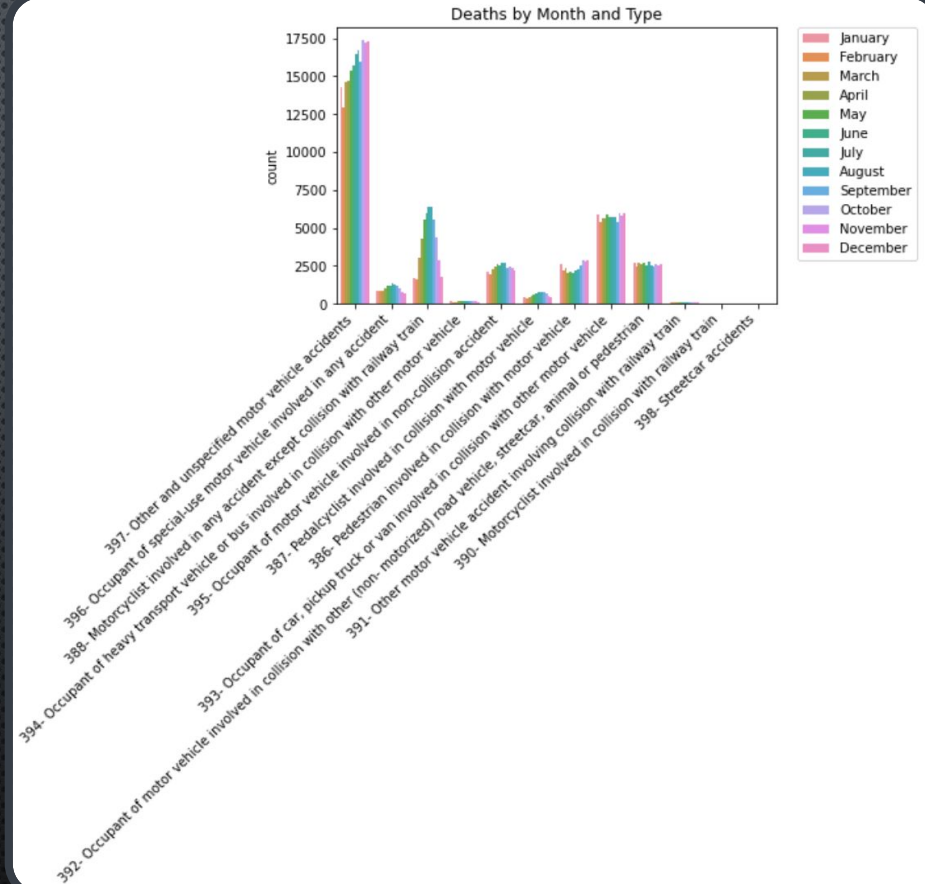
Trustworthy Source

SOURCE

- DATASET: DEATHS IN THE UNITED STATES (CDC, 2017)
 - CAUSES OF DEATH FROM 2005-2015
 - OVER 600 SOURCE CODES (CATEGORIES) IN DATA SET
 - FOCUSED ON MOTOR VEHICLE DEATHS: 12 SOURCE CODES
- CSVs FOR EACH YEAR
- SECONDARY DATA SOURCE: US CENSUS DATA

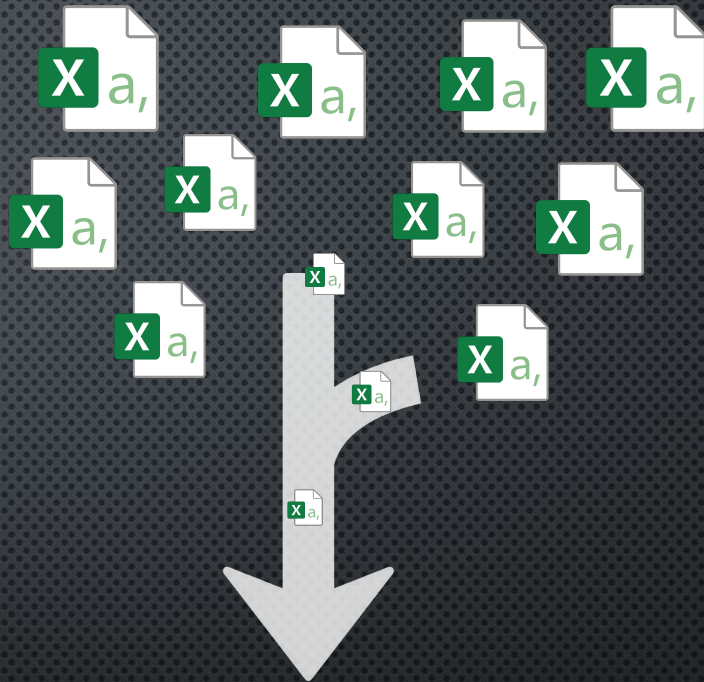
LIMITATIONS

- LACK OF DRIVING & LOCATION DATA
- A LARGE PERCENTAGE OF MOTOR VEHICLES DEATHS WERE “UNSPECIFIED”
 - 45.1% OF MOTOR VEHICLES DEATHS REPORTED WERE CONSIDERED UNSPECIFIED
- SOME STATISTICAL TESTS/HYPOTHESES WE WANTED TO TEST WERE NOT APPLICABLE
 - DATA DID NOT HAVE EQUAL VARIANCES/DISTRIBUTION IN SOME CASES



DATA WRANGLING

STEP 1: MERGE ALL OF THE CSVs INTO A SINGLE DATAFRAME, USING ONLY THE COLUMNS WE NEED



	month_of_death	sex	detail_age	day_of_week_of_death	current_data_year	manner_of_death	358_cause_recode
207	1	M	32	7	2005	1.0	396
208	1	M	75	5	2005	1.0	387
220	1	M	68	7	2005	1.0	396
234	2	M	21	1	2005	1.0	396
235	2	M	24	1	2005	1.0	396

STEP 2: USED JSON FILES TO TRANSFORM OUR CODED DATAFRAME INTO SOMETHING READABLE

```
▼ "day_of_week_of_death" : { 8 items
  "1" : string "Sunday"
  "2" : string "Monday"
  "3" : string "Tuesday"
  "4" : string "Wednesday"
  "5" : string "Thursday"
  "6" : string "Friday"
  "7" : string "Saturday"
  "9" : string "Unknown"
}
▼ "current_data_year" : { 1 item
  "2008" : string "2008"
}
▼ "injury_at_work" : { 3 items
  "Y" : string "Yes"
  "N" : string "No"
  "U" : string "Unknown"
}
▼ "manner_of_death" : { 8 items
  "1" : string "Accident"
```



```
26 manner_of_death_dict = {
27     1:"Accident",
28     2:"Suicide",
29     3:"Homicide",
30     4:"Pending investigation",
31     5:"Could not determine",
32     6:"Self-Inflicted",
33     7:"Natural"}
34 # "Blank": "Not specified"}
35
36 cause_recode_dict = {
37     385:" 385- Motor vehicle accidents",
38     386:" 386- Pedestrian involved in collision with motor vehicle",
39     387:" 387- Pedalcyclist involved in collision with motor vehicle",
40     388:" 388- Motorcyclist involved in any accident except collision with railway train",
41     389:" 389- Motor vehicle accident involving collision with railway train",
42     390:" 390- Motorcyclist involved in collision with railway train",
43     391:" 391- Other motor vehicle accident involving collision with railway train",
44     392:" 392- Occupant of motor vehicle involved in collision with other (non- motorized) road vehicle, streetcar, animal",
45     393:" 393- Occupant of car, pickup truck or van involved in collision with other motor vehicle",
46     394:" 394- Occupant of heavy transport vehicle or bus involved in collision with other motor vehicle",
47     395:" 395- Occupant of motor vehicle involved in non-collision accident",
48     396:" 396- Occupant of special-use motor vehicle involved in any accident",
49     397:" 397- Other and unspecified motor vehicle accidents",
50     398:" 398- Streetcar accidents"}
51
52
53 clean_df = car_death_data.replace({"month_of_death": month_dict,
54                                   "day_of_week_of_death": day_of_week_dict,
55                                   "manner_of_death": manner_of_death_dict,
56                                   "358_cause_recode": cause_recode_dict})
57
58 clean_df
```

CDC definitions (JSON format)

Replacing coded values in dataframe with new values from dictionaries

STEP 2: USED JSON FILES TO TRANSFORM OUR CODED DATAFRAME INTO SOMETHING READABLE

THE RESULTING DATAFRAME IS CLEAN, EASY TO UNDERSTAND & READY FOR ANALYSIS!

	month_of_death	sex	detail_age	day_of_week_of_death	current_data_year	manner_of_death	358_cause_recode
207	January	M	32	Saturday	2005	Accident	396- Occupant of special-use motor vehicle in...
208	January	M	75	Thursday	2005	Accident	387- Pedalcyclist involved in collision with ...
220	January	M	68	Saturday	2005	Accident	396- Occupant of special-use motor vehicle in...
234	February	M	21	Sunday	2005	Accident	396- Occupant of special-use motor vehicle in...
235	February	M	24	Sunday	2005	Accident	396- Occupant of special-use motor vehicle in...
...
2717184	December	M	67	Friday	2015	Accident	397- Other and unspecified motor vehicle acci...
2717278	December	F	77	Thursday	2015	Accident	397- Other and unspecified motor vehicle acci...
2717674	December	M	73	Wednesday	2015	Accident	397- Other and unspecified motor vehicle acci...
2717998	December	M	70	Thursday	2015	Accident	395- Occupant of motor vehicle involved in no...
2718169	December	F	63	Thursday	2015	Accident	386- Pedestrian involved in collision with mo...

STEP 3: CHECK THE DATA QUALITY & REMOVE ANY INACCURATE OR AMBIGUOUS DATA

```
In [5]: 1 # check quality of data
        2 # list all unique values in each columns
        3
        4 colNames = list(clean_df.columns)
        5 for col in colNames:
        6     print(col)
        7     print(f"{clean_df[col].unique()}")
        8     print("-----")

month_of_death
['January' 'February' 'March' 'April' 'June' 'May' 'July' 'August'
 'September' 'October' 'November' 'December']
-----
sex
['M' 'F']
-----
detail_age
[ 32  75  68  21  24  25  44  49  40  11  14  64  57  18   9  54  31  42
 27  43  58  46  33  22  13  61  56  16  53  30  23  26  20  83  34  38
 85  15   1  47  65  72  52  17  28  12  19  48  45  50  66  81  41  62
 59  69  39  37  63  78  80  35  73  36  90  999  82  55  51   2  60  70
 29  91  79   7  84  76   8   5  77  89  74   4  10  88  86   6   3  71
 67  92 101  94  87  93  96 102 100  99  95  98 106 104 103 107]
-----
day_of_week_of_death
['Saturday' 'Thursday' 'Sunday' 'Tuesday' 'Wednesday' 'Monday' 'Friday'
 'Unknown']
-----
current_data_year
[2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015]
-----
manner_of_death
['Accident' 'Natural' 'Could not determine' 'Homicide' nan
 'Pending investigation' 'Suicide']
-----
358_cause_recoded
[' 396- Occupant of special-use motor vehicle involved in any accident'
 ' 387- Pedalcyclist involved in collision with motor vehicle']
```

```
In [6]: 1 # Clean up
        2 # Remove not Logical data
        3 # ie: age of 999, day of week : Unknown, manner_of_death nan, need filter by Accident
        4 finalDf = clean_df[clean_df["detail_age"] != 999]
        5 finalDf = finalDf[finalDf["manner_of_death"] == "Accident"]
        6 finalDf = finalDf[finalDf["day_of_week_of_death"] != "Unknown"]
        7 finalDf.head()
```


NOW, WE CAN USE OUR FINAL, CLEAN DATAFRAME TO ANSWER
OUR RESEARCH QUESTIONS!



DOES THE AVERAGE NUMBER OF VEHICLE DEATHS VARY BY MONTH?

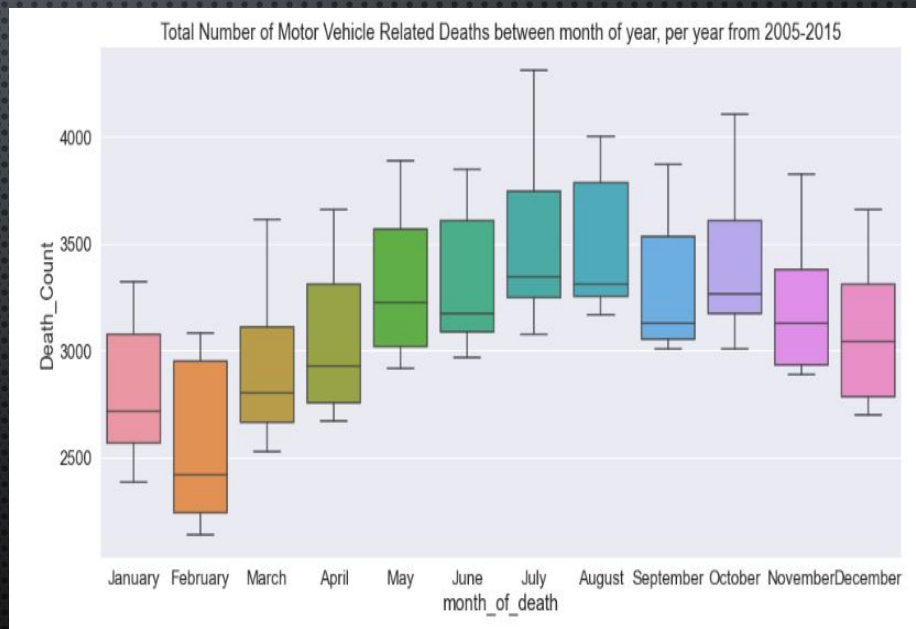
- 1.) GROUP THE DATAFRAME BY MONTH, THEN SUM HOW MANY DEATHS OCCURRED IN EACH MONTH PER YEAR
- 2.) LEVENE'S TEST TO TEST FOR EQUAL VARIANCE
- 3.) ANOVA TO DETERMINE IF THERE IS A DIFFERENCE BETWEEN GROUPS IN OUR POPULATION
- 4.) F-ONEWAY TEST TO COMPARE SPECIFIC GROUPS

```
=====
Levene Summary
=====
      W      pval  equal_var
-----
0.072   1.000    True

=====
ANOVA SUMMARY
=====

Source      ddof1    ddof2      F      p-unc      np2
-----
month_of_death      11      120    7.664    0.000    0.413

=====
Feb vs. Aug
=====
F_onewayResult(statistic=42.425595897501985, pvalue=2.385867742925349e-06)
=====
Jul vs. Aug
=====
F_onewayResult(statistic=0.00010788018826291518, pvalue=0.9918157990446276)
```



DOES THE AVERAGE NUMBER OF VEHICLE DEATHS VARY BY SEASON?

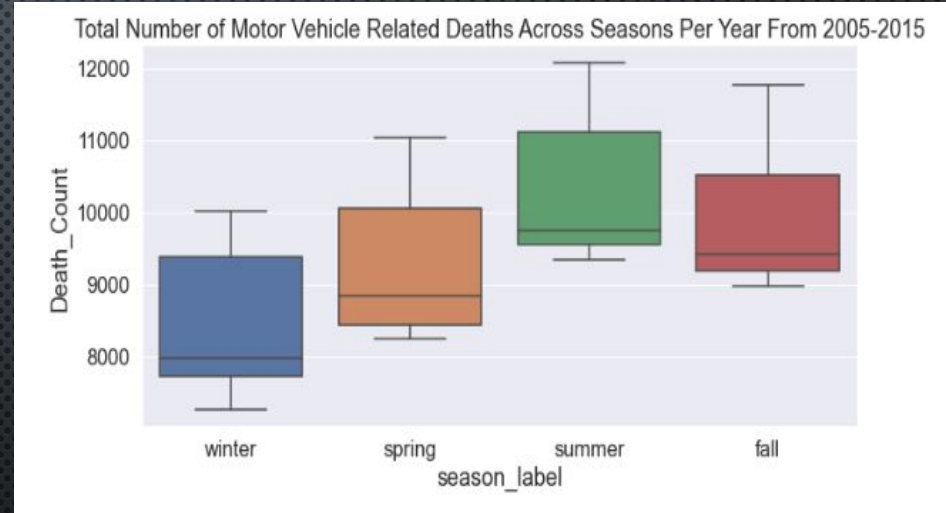
1.) SEASON DEFINITIONS

- A. Winter: Dec-Feb
- B. Spring: March – May
- C. Summer: June-Aug
- D. Fall: Sept-Nov

2.) GROUP THE DATAFRAME BY SEASON, THEN SUM HOW MANY DEATHS OCCURRED IN EACH SEASON PER YEAR

3.) LEVENE'S TEST TO TEST FOR EQUAL VARIANCE

4.) ANOVA TO DETERMINE IF THERE IS A DIFFERENCE BETWEEN GROUPS IN OUR POPULATION



DOES THE AVERAGE NUMBER OF VEHICLE DEATHS VARY BY SEASON?

1.) P-VALUE < 0.05 FROM ANOVA REJECT NULL HYPOTHESIS

2.) PAIRWISE TUKEY TEST TO COMPARE DIFFERENT SEASONS

3.) COMPARING SUMMER AND WINTER, PAIRWISE TUKEY AND F-ONEWAY HAVE DIFFERENT RESULTS

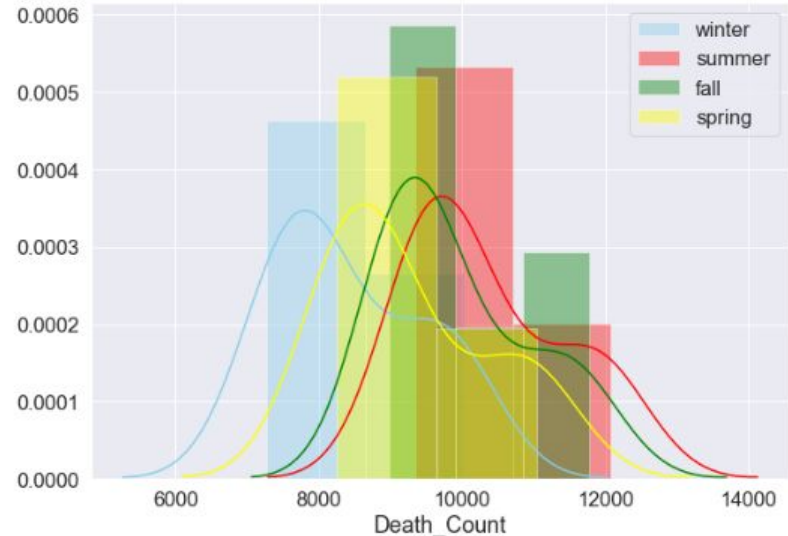
=====
POST HOC TESTS
=====

A	B	mean(A)	mean(B)	diff	se	tail	T	p-tukey	hedges
fall	spring	8441.000	9267.182	-826.182	445.349	two-sided	-1.855	0.252	-0.761
fall	summer	8441.000	10350.000	-1909.000	445.349	two-sided	-4.287	0.001	-1.758
fall	winter	8441.000	9940.455	-1499.455	445.349	two-sided	-3.367	0.005	-1.381
spring	summer	9267.182	10350.000	-1082.818	445.349	two-sided	-2.431	0.076	-0.997
spring	winter	9267.182	9940.455	-673.273	445.349	two-sided	-1.512	0.433	-0.620
summer	winter	10350.000	9940.455	409.545	445.349	two-sided	0.920	0.768	0.377

=====
FOneWay : winter vs. Summer
=====

F_onewayResult(statistic=18.781165691792996, pvalue=0.00032243751192492423)

Total Number of Motor Vehicle Related Deaths Across Seasons Per Year From 2005-2015



IS THERE AN INTERACTION BETWEEN AGE AND DAY OF WEEK FOR MOTOR VEHICLE DEATHS?

SPECIFICALLY:

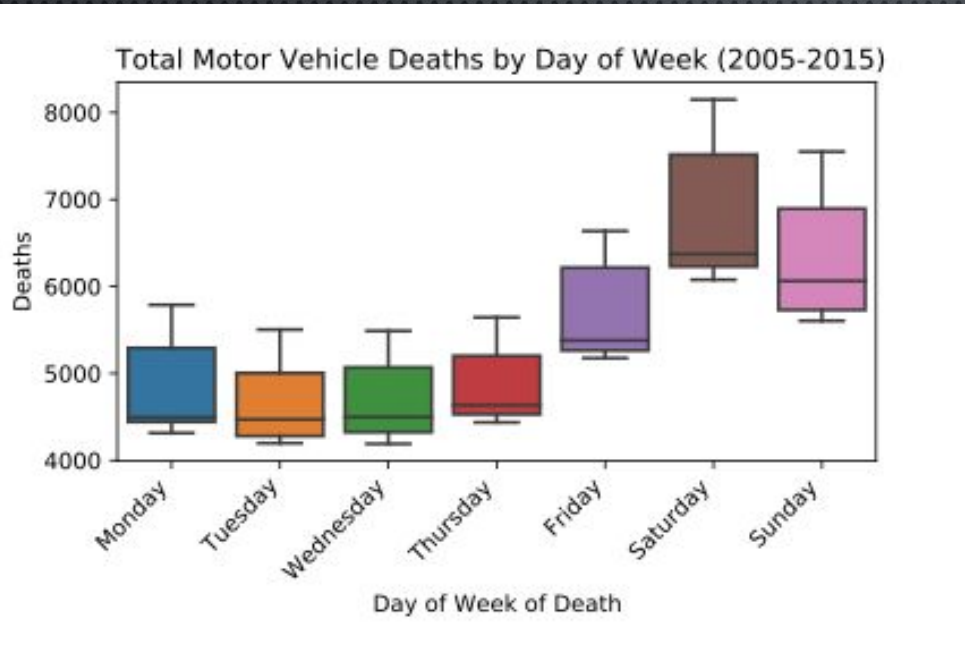
- DO YOUNGER PEOPLE SHOW A STRONGER TENDENCY TOWARD WEEKEND FATALITIES?
- H_0 IN ALL CASES IS THAT THERE IS NO EFFECT (WHETHER AGE, DAY OF WEEK, OR BOTH)
- MOTIVATION: POPULAR CULTURE AND HIGH SCHOOL EXPERIENCE

There's a Wikipedia page...



https://en.wikipedia.org/wiki/List_of_car_crash_songs

MOTOR VEHICLES DEATHS BY DAY OF WEEK



LEvene's Test

P-VALUE = 8.27 E-01 (SAME FOR STATS PACKAGE AND PINGOUIN)

$P \geq 0.05$

➤ SAMPLE VARIANCES SIMILAR ENOUGH FOR ANOVA

ANOVA F-ONEWAY

STATISTIC = 7.66

$P = 6.62 \text{ E-}10$ (STATS PACKAGE)

$P(\text{UNCORR.}) = 1.11 \text{ E-}14$ (PINGOUIN)

➤ REJECT THE NULL HYPOTHESIS THAT THE SAMPLE MEANS ARE EQUAL

MOTOR VEHICLE DEATHS BY AGE GROUP

LEVENE'S TEST

P-VALUE = 5.79 E-09 (PINGOIN PACKAGE)

P >= 0.05

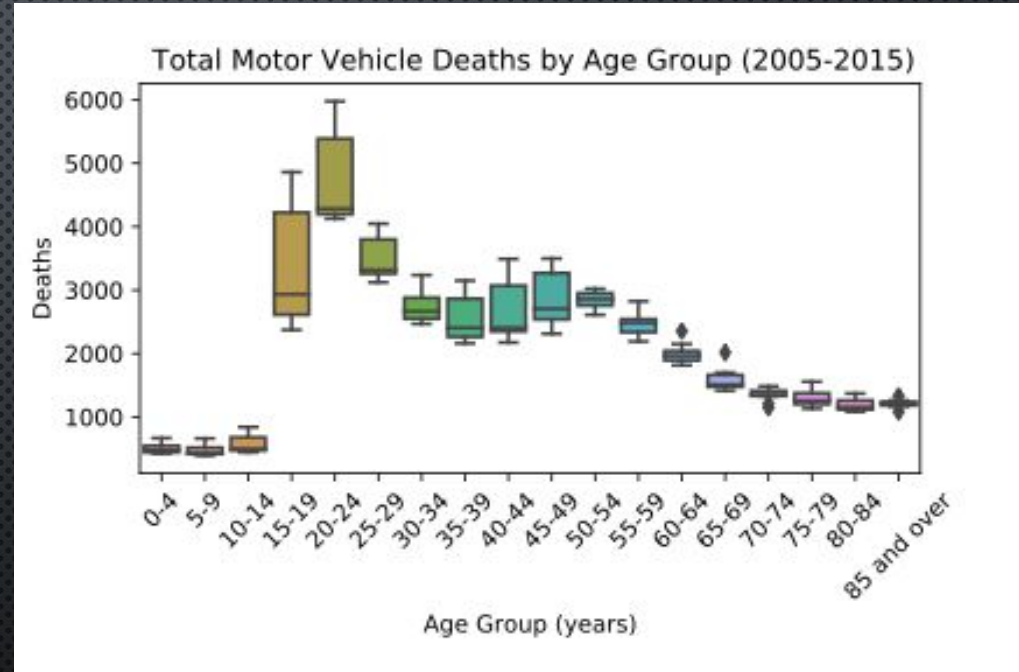
➤ SAMPLE VARIANCES TOO DISSIMILAR FOR ANOVA

KRUSKAL-WALLIS H-TEST FOR INDEPENDENT SAMPLES

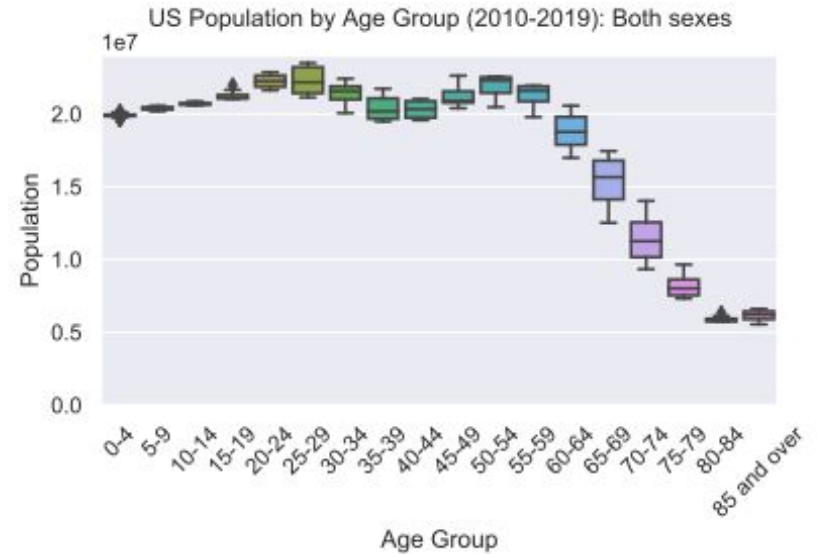
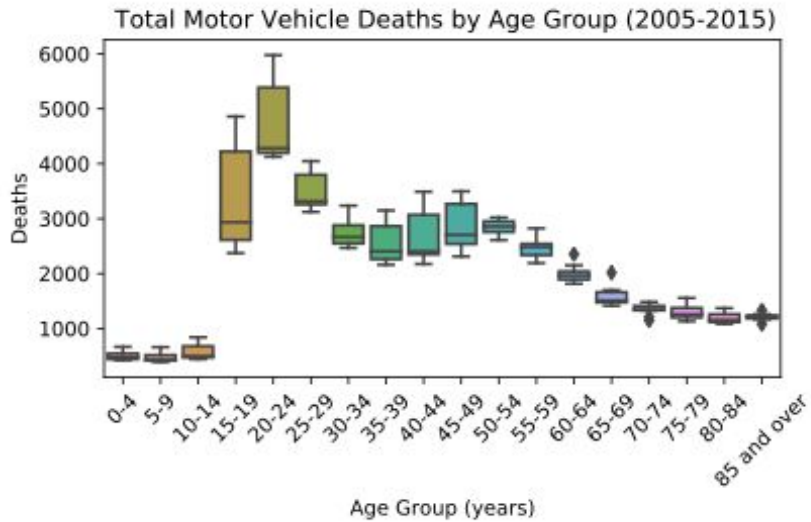
STATISTIC = 183.13

P = 6.84 E-30 (PINGOIN PACKAGE)

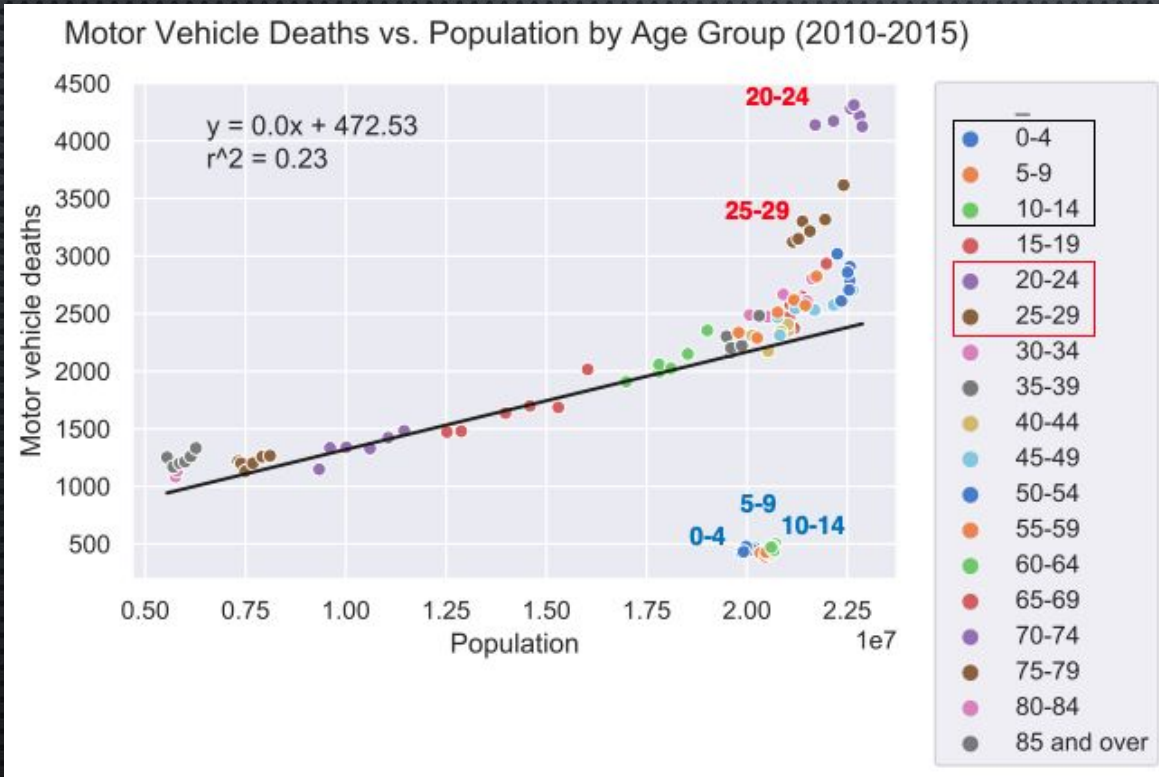
➤ REJECT THE NULL HYPOTHESIS THAT THE SAMPLE MEDIAN ARE EQUAL



CAN SOME AGE-GROUP EFFECTS BE ATTRIBUTED TO DIFFERING COHORT SIZES?

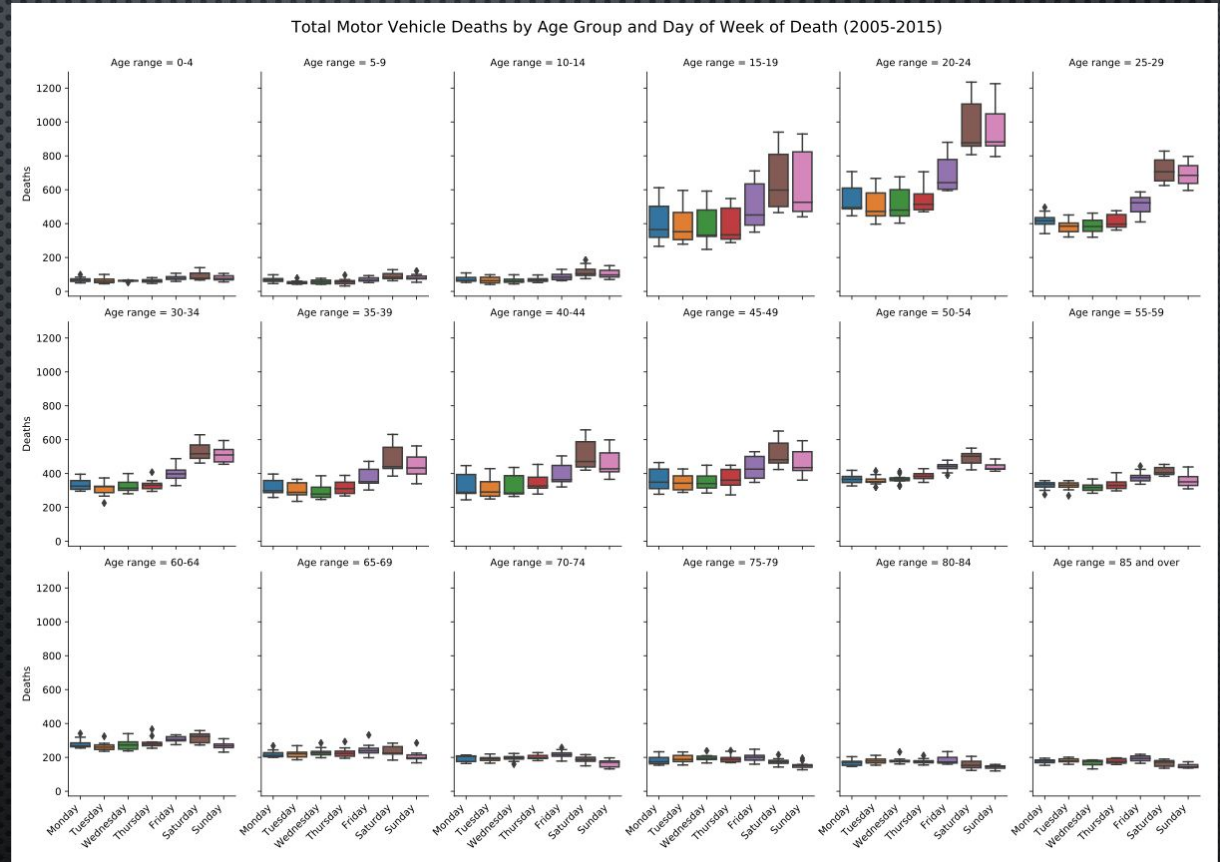


MOTOR VEHICLE DEATHS BY AGE GROUP



LINEAR REGRESSION OF MOTOR VEHICLE DEATHS VS. POPULATION SUGGESTS THERE ARE OTHER FACTORS, PARTICULARLY FOR CHILDREN AND YOUNG ADULTS

- AGE-GROUP VARIANCES ARE TOO DISSIMILAR FOR ANOVA
- CAN AGE GROUPS BE TESTED ONE BY ONE?
- VISUAL IMPRESSION IS THAT ANY WEEKEND EFFECT MAY DISSIPATE WITH AGE



- VARIANCES WITHIN AGE GROUPS ARE SUITABLE FOR ANOVA TEST OF DAY-TO-DAY DIFFERENCES

	W	pval	equal_var
levene	1.94635	0.085272	True
	W	pval	equal_var
levene	0.882933	0.512125	True
	W	pval	equal_var
levene	1.078084	0.384021	True
	W	pval	equal_var
levene	0.933086	0.476952	True
	W	pval	equal_var
levene	0.661903	0.680509	True
	W	pval	equal_var
levene	1.668938	0.141543	True
	W	pval	equal_var
levene	0.805239	0.569251	True
	W	pval	equal_var
levene	0.988245	0.440008	True
	W	pval	equal_var
levene	0.35972	0.901849	True
	W	pval	equal_var
levene	0.830179	0.550596	True
	W	pval	equal_var
levene	0.866909	0.523659	True
	W	pval	equal_var
levene	0.565013	0.75671	True
	W	pval	equal_var
levene	0.519195	0.791876	True
	W	pval	equal_var
levene	0.343863	0.911163	True
	W	pval	equal_var
levene	0.409808	0.870158	True
	W	pval	equal_var
levene	0.563515	0.757874	True
	W	pval	equal_var
levene	1.066215	0.391112	True
	W	pval	equal_var
levene	0.840581	0.5429	True

	p (uncorr.)
Overall	0
Age group	
0-4	0.000039
5-9	0.000002
10-14	8.33E-07
15-19	0.000006
20-24	6.80E-19
25-29	5.00E-28
30-34	2.24E-23
35-39	3.37E-12
40-44	1.43E-09
45-49	2.40E-09
50-54	1.22E-18
55-59	7.35E-11
60-64	0.000145
65-69	0.066528
70-74	0.000002
75-79	0.000216
80-84	0.000015
85 and over	4.41E-07

- ONE-WAY ANOVA SHOWED SIGNIFICANT DIFFERENCES AMONG DAYS FOR ALL BUT THE 65-69 AGE GROUP
- AT THIS LEVEL OF ANALYSIS, YOUNGER PEOPLE ARE NOT DIFFERENT

POST-HOC TESTS SHOWED SIGNIFICANT DIFFERENCES IN MOTOR VEHICLE FATALITY COUNTS AMONG ALL PAIRS OF WEEKDAYS AND AGE GROUPS, EXCEPT THOSE SHOWN HERE

1	Post-hoc tests corrected for multiple-comparisons													
2	Contrast	A	B	Paired	Parametric	T	dof	Tail	p-unc	p-corr	p-adjust	BF10	hedgesVNAME?	
3	day_of_week_of_death	Monday	Thursday	TRUE	TRUE	-1.276	197	two-sided	0.204	0.204	fdr_bh	0.177	-0.016\	
4	day_of_week_of_death	Tuesday	Wednesday	TRUE	TRUE	-1.816	197	two-sided	0.071	0.074	fdr_bh	0.398	-0.024\	
5														
6	Age_group	45-49	50-54	TRUE	TRUE	0.229	76	two-sided	0.819	0.819	fdr_bh	0.129	0.023\	
7	Age_group	80-84	85 and over	TRUE	TRUE	-1.039	76	two-sided	0.302	0.304	fdr_bh	0.211	-0.134\	
8	Age_group	35-39	55-59	TRUE	TRUE	1.368	76	two-sided	0.175	0.178	fdr_bh	0.307	0.166\	
9	Age_group	30-34	40-44	TRUE	TRUE	1.493	76	two-sided	0.14	0.142	fdr_bh	0.363	0.085\	
10	Age_group	15-19	25-29	TRUE	TRUE	-1.53	76	two-sided	0.13	0.134	fdr_bh	0.383	-0.112\	
11														

- SOME TENDENCY FOR AGE GROUPS TO CLUSTER
- NO EVIDENCE YET TO CLEANLY SEPARATE WEEKENDS FROM WEEKDAYS
- LIKELY NEXT TEST WOULD BE TO GROUP (SOME) WEEKDAYS VS. WEEKEND
- ANALYSES TO DATE ARE EASIEST WITH EQUAL SAMPLE SIZES, SO WOULD CHOOSE E.G. SAT/SUN VS. TUE/WED

CONCLUSIONS/ FURTHER QUESTIONS

- THERE ARE SIGNIFICANT DIFFERENCES IN MOTOR VEHICLE DEATHS BETWEEN:
 - SUMMER & WINTER
 - AGE
 - DAYS OF WEEK
- QUESTIONS WE HAVE AFTER OUR ANALYSES:
 - WHAT IS THE CORRELATION BETWEEN AUTO INSURANCE RATES AND AGE OF DRIVER?
 - FUTURE ANALYSIS: WILL MOTOR VEHICLE DEATHS DECREASE WITH THE NEWLY ENFORCED TEXTING AND DRIVING LAWS