# **Exploratory Data Analysis**

Capstone Step 5

# Choices Regarding Data Cleaning/Transformation

During prototyping I observed a few high level improvements possible via treating the raw data. Specific choices are annoted in screenshots on the subsequent slides.

- 1. Location file is entirely redundant and will be descoped.
- 2. A number of columns are redundant and will be removed after verifying the other locations offer equivalent data
- 3. A number of columns can be joined/consolidated
- 4. A number of columns can be converted from varchar to numerical type

# Choices Regarding Pipeline Automation

The prototyping phase helped form the following approach to automation in order to achieve the optimal reliability:

- 1. To handle embedded commas within csv columns, use Linux csvkit's csvcut tool from Python (cut and awk both have tremendous difficulty).
- 2 . Use sqlalchemy create\_engine, pandas read\_csv and dataframe.to\_sql instead of looping through each row to execute mysql\_connector\_python insert statement. The latter fails in cases of NULL values.
- 3. Dot not declare any columns NOT NULL, and err on the side of caution with strings (TEXT instead of VARCHAR).
- 4. Delete files after processing to keep to manageable level.

## Details

## Convert 6 columns to 2 of the part of t	A						G	Н	1	J	К	L	М	N	0				Р			Q	К	S	T	U	V	
DRECT INJURES   INDURECT DEATHS   INDURECT DEATHS   INDURECT DEATHS   INDURED   INDURES   INDURECT DEATHS   INDURED   INDURES   INDURE	4 2	23 207 4 950 6 528 6 1055 25 131 25 140 26 180 27 180 28 183 30 28 183 30 28 183 30 140 9 140 9 140 15 165 173 184 185 185 185 185 185 185 185 185	7 0 20 20 20 20 20 20 20 20 20 20 20 20 2	2004 2008 2008 2002 2002 2007 2007 2007 2008 2004 2003 2004 2009 2009 2009 2009 2009 2009 2009	23 4 6 6 6 6 25 27 25 27 28 13 28 14 9 9 9 6 6 13 31 6	228 951 529 1100 2000 11800 600 2000 330 1930 1530 545 1615 1606 1452 1330 816 809	49216 45823 45826 50353 51466 52434 47310 46523 45618 45617 46730 52853 52853 52853 54986 46987 44988 47487 47487 47487 47255	904094 875896 871938 913730 912708 918106 885808 880416 874455 883021 920637 920638 883678 883679 877686 887356 887356	PENISYLVANIA GEORGIA NORTH CAROLINA TEXAS TEXAS TEXAS TEXAS WEST VIRGINIA VIRGINIA OHIO OHIO PUERTO RICO SOUTH CAROLINA FLORIDA FLORIDA ATLANTIC SOUTH ATLANTIC SOUTH ATLANTIC SOUTH ATLANTIC SOUTH ATLANTIC SOUTH ATLANTIC SOUTH	1.3 3.7 4.8 4.8 4.8 5.4 5.1 5.9 5.9 5.9 5.9 5.9 5.9 5.9 5.9	2020 J 2020 J 2020 J 2020 A 2020 A 2020 J 2020 J 20	ebruary uly uly ugus ey  Albeh Arbh Arbh Arbh Arbh Arbh Arrh Arach	Tornado Tornado Tornado Tornado Hurricane Hurricane Hurricane Hurricane Hurricane Hash Flood Flo	CZ, TYP  C C C C C Z Z Z Z C C C C C C C C C C	5 17 15 45 256 442 216 101 51 37 89 167 127 47 5 136 350 350	AMITE BUCKS BARTOW CLEVELAND CLEVELAND CLEVELAND COASTAL WILL KLEBERG COL ORANGE WEBSTER DICKENSON DARKE LICKING WASHINGTON SAN JUAN JASPER ALLENDALE EASTERN ALA S SANTEER T S SANTEER T S SANTEER T S SANTEER T	CHUA O EDISTO BE/ O EDISTO BB/ BBAY FROM L'	ACH SC OUT 20NM ITLE CREEK, VA, 1	TO CAPE HENRY, VA, IN	CLUDING THE CHESAPAGIAN BORDA	EAKE BAY BRIDGE TUNNER BRYOND SIM OF SHO	BRC 25-JUL CRP 25-JUL LCH 26-AUG RLX 25-MAY ILN 28-MAR ILN 28-MAR ILN 28-MAR SJU 09-SEP SJU 09	20 05:2 00 20 10:5 00 20 10:5 00 20 10:5 00 20 10:5 00 20 20 20 20 20 20 20 20 20 20 20 20	51-5 GFE 51-6 S-JUL- 51-6 S-JUL- 51-6 7-AUG 51-6 7-AUG 51-6 7-AUG 51-5 3-APR 51-5 8-MAR 51-5 3-APR	20 05 9:00 20 11 00:00 0 2 00:00 0 1 00:00 0 2 00:00 0 3 00:00 0 4 20:00:00 4 20:00:00 4 20:00:00 4 30:00:00 4 10:00:00 4 10:00:00 4 10:00:00 4 10:00:00 4 10:00:00 4 10:00:00 4 10:00	No.   No.	JRIES_INC	
V	dt Convert 6	6 column	ns to 2			n/end	i datetin	i	V	1	(de aft ve	elete er rifying		V	i		V					V		Redundant	r verifyin	i ng equivale	i ence)	
O   O   O   O   O   O   O   O   O   O	U ES.DIRECT INJUR	V		RECT DEATHS	X INDIRECT D	Y AMAGE PR				A MAGN	AB WITUDE N	AC	AD	AE	V - Vi dt - d i - ini	archar latetime		AI H TOR, OTHER, W			AL PSTOR OTHER CZ NAM	AM E BEGIN RANGE	AN E BEGIN AZIMU	JTH BEGIN LOCATIO	AO	EN		
	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0		0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	10 77 77 12 13 55 1 1 3 3 5 5 0 0 0 0 1 1 1 5 5 5 5 5 5 5 5 5 5 5	00.00K  0.00K  0.00K	0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.0	DK D	Emergency Mana WWS Storm Surv WWS Storm Surv WWS Storm Surv WWS Storm Surv Budy Mesonet Emergency Mana Department of He 311 Call Center state Official Law Enforcement Herforcement Law Enforcement 111 Call Center Law Enforcement Wesonet Wesonet Wesonet Wesonet Wesonet	ger ger 35 50 41 49 42 36 51	E N N N	EG AG AG AG AG	Heavy Rain Heavy Rain Heavy Rain Heavy Rain Heavy Rain Heavy Rain		EF2 EF2	18.3 0.06 0.33	1760 500 200	PHI	MS PA GA	113 101	PIKE PHILADELPHIA GORDON	2 2 1 1 1 1 1 1 2 1 4	W WNW ESE NNW E E E N NW W WNW	EAST FORK ORNMELLS HOT FOLSOM ARCHDALE  ERBACON OSBORNS GAP GREENVILLE PATASKALA ZEU MACKSBURG SAN JUAN ALLENDALE OSW FOLLYANS FOLLYANS ISLA OMESAPEACH SULLIVANS ISLA OMESAPEACH	INE ARPT  WALD ARP  IND  AY BRIDGE TUN	111 2 2 4 4 4 2 2 2 2 2 1 1 3 4 4 2 2 1 1 1 1 1		1
													i - int	rchar														

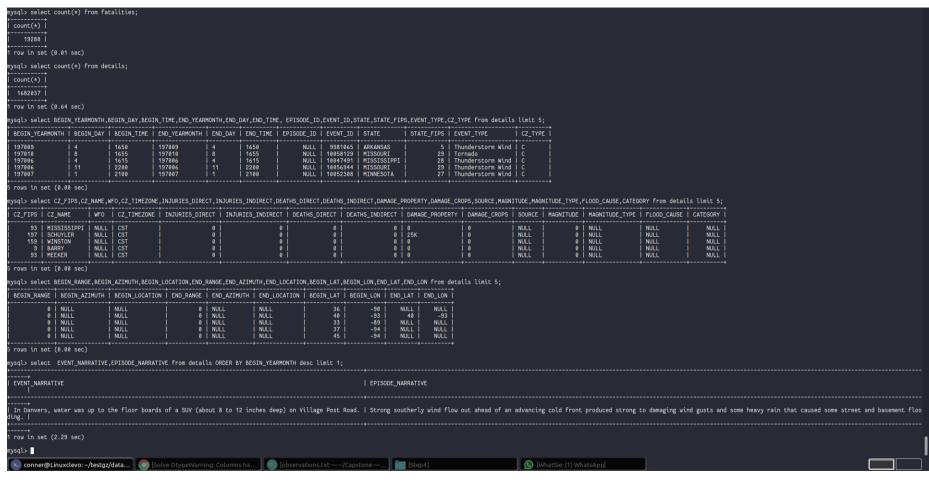
Details (cont)

								AQ	AR	AS	AT	AU	AV	AW	AX	AY
1							END_	AZIMUTH	END_LOCATION	BEGIN_L/	AT BEGIN_LO	N END_LAT	END_LON	EPISODE_N EV	ENT_NAF D	TA_SOURCE
2							E		MARS HILL	31.2047	-90.7432	31.2705	-90.446	A strong col A v	ery large C	V I
3							W			40.0815	-74.9592	40.0822	-74.9599	Tropical Stor A t	ornado to C	V
3							NW		FOLSOM	34.3951	-84.8631	34.3968	-84.8576	A line of thu A I	lational V C	V
5							SSE		KINGS MTN	35.177	-81.413	35.196	-81.325	Unusually hir NV	S storm C	7
6														Hurricane HvHu		<b>\</b> /
Ť														Hurricane HrTC	DON site C	
Ŕ														A tropical we The		V.
- 8							E		ERBACON	38.5367	-80.5887	38.5186		Scattered ar Co		
9 10							NE			37.1988	-82.5269			An intense I At		
11							S			40.1	-84.62			Thunderstor Hig		
12							NW			39.97	-82.62			Scattered th Hig		
13							NE			39.6391		39.6392	-81.4539	An intense I On	the West Ct	v A
14							NW			18.4149	-66.0901			An upper-lev Flo		
15							NW		SAN JUAN	18,414				An upper-lev Flo		
16														Strong gradi La		
<del>-19</del>							WNW		ALLENDALE OSWALD ARP	33.01	-81.4	33.01	-81.4	An area of le Alle	ndale Co	
16									THE CONTROL OF THE	00.01	02.4	00.01		Very strong Str		
18 19							WSW		FOLLY BEACH	32.65	-79.94	32.65	-79 94	A severe qu Th	Weather C	<i>)</i> \
20							F			32.76	-79.82			A severe gu Th		
- 51							N N		CHESAPEAKE BAY BRIDGE TUNNEL 3RD ISLAN		-76.08			Scattered si Wi		
-53							N			47.18	-87.23			A potent sto The		
-55									DIVINE NO ROCK ETCH	41.10	01.20	41.20	01.20	rt potent ste m	- Ottarina O	
- 25 - 25 - 26										-	-	-				
- 5 =							_ V		V	- T	- 1	- 1	- 1	V	V	Delete
26												_				
<del>- 57</del>																
26																
28 29 30																
30												_	_			
31					_						_	_	_			
31																
31 32																
31 32 33																
31 32 33 34																
31 32 33 34 35																
31 32 33 34 35 36																
31 32 33 34 35 36 37																
31 32 33 34 35 36 37 38																
31 32 33 34 35 36 37 38																
31 32 33 34 35 36 37 38 39 40																
31 32 33 34 35 36 37 38 39 40 41																
31 32 33 34 35 36 37 38 39 40 41 42																
31 32 33 34 35 36 37 38 39 40 41 42 43																

### Fatalities

Α	В	C	l D	E	F	_ (	G	н	l ı	J	K
FAT YEARMONTH									FATALITY SEY	FATALITY LOCATION	
202106		0	42960	953511	D D	06/09/2021		70	M	Golfing	202106
202100		0	43206	961309	ı	07/20/2021			M	Other	202100
202107		0	43207		1	07/20/2021			M	Other	202107
		0						31			
202107		<b>y</b>			!	07/20/2021			M	Other	202107
202107	27	<u> </u>	44279		I	07/27/2021			•	Other	202107
202103		0			D	03/28/2021			F	In Water	202103
202104		0	42962	954336	D	04/24/2021		2		Permanent Home	202104
202106		0	42963	954408	D	06/10/2021			M	In Water	202106
202106		0			D	06/14/2021		23	M	Outside/Open Areas	202106
202106	13	0	44482	970319	D	06/13/2021	L 00:00:00	32	F	Outside/Open Areas	202106
202106	13/	0	44483	970319	D	06/13/2021	L 00:00:00	42	M	Outside/Open Areas	202106
202106	M	0	44484	970319	D	06/13/2021	L 00:00:00	29	M	Outside/Open Areas	202106
202106	<b>A</b>	0	44485	970319	D	06/15/2021	L 00:00:00	34	F	Outside/Open Areas	202106
202106		0		970319	D	06/17/2021		28	M	Outside/Open Areas	202106
202106		0			D	06/20/2021			M	Outside/Open Areas	202106
202106		0		970319	D	06/20/2021			M	Outside/Open Areas	202106
202106		0		970319	D	06/15/2021		35	M	Outside/Open Areas	202106
202106		0	44491	970319	D	06/17/2021		36	M	Outside/Open Areas	202106
202106		0		970319	D	06/17/2021			M	Outside/Open Areas	202106
202106	13	0			D			35	F		202106
		<u> </u>		970320		06/13/2021				Outside/Open Areas	
202106	12	0	44480	970320	D	06/12/2021			M	Outside/Open Areas	202106
202106		0		970320	D	06/15/2021		37	M	Outside/Open Areas	202106
20210		0		970487	D	06/02/2021			M	Outside/Open Areas	202106
202106		0		970487	D	06/02/2021			M	Outside/Open Areas	202106
202106		0			D	06/07/2021			M	Outside/Open Areas	202106
202106	6	0	43771	970490	D	06/06/2021	L 00:00:00	47	M	Outside/Open Areas	202106
					1.7		- T			1.7	
					V		DT		V	V	
Redund				I	V		DT	I	V	V	The event can
(Will rer	nove after			I	V		DT	I	V	V	
(Will rer verifying	nove after equivale	nt		I	V		DT		V	V	occur prior to
(Will rer	nove after	nt		I	V				V	V	occur prior to the fatality, so
(Will rer verifying to FATA	nove after equivale	nt		I	V		Datatyp	25:	V	V	occur prior to the fatality, so this is not
(Will-rer verifying to-FATA	nove after equivale	nt		I	V		Datatype I - Int		V	V	occur prior to the fatality, so
(Will rer verifying to FATA	nove after equivale	nt		I	V		Datatyp		V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt		I	V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifying to FATA	nove after equivale	nt		I	V		Datatype I - Int	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifying to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	the fatality, so this is not
(Will rer verifying to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifying to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt		I	V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to-FATA	nove after equivale	nt		I	V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to-FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not
(Will rer verifyin to FATA	nove after equivale	nt			V		Datatype I - Int V - Varci	ar	V	V	occur prior to the fatality, so this is not

## After running transform.py, both tables are fully populated in mysql:



#### Data is returning meaningful queries for analysis:

