	<pre>print(outage[['OUTAGE.START.DATE', 'OUTAGE.START.TIME', 'OUTAGE.RESTORATION.DATE', 'OUTAGE.RESTORATION.TIME']].head().to_markdown(index=</pre>	
:[]:	# Drop the old time-keeping values from the dataframe outage = outage.drop(columns=['OUTAGE.START.DATE', 'OUTAGE.START.TIME', 'OUTAGE.RESTORATION.DATE', 'OUTAGE.RESTORATION.TIME']) print(outage[['OUTAGE.RESTORATION', 'OUTAGE.START']].head().to_markdown(index=False)) outage[['OUTAGE.RESTORATION', 'OUTAGE.START']].dtypes OUTAGE.RESTORATION OUTAGE.START	
[]:]: # open output in a text editor to view in full # Finds the datatype for each column in the outage dataframe print(outage.dtypes.to_markdown(index=True))	
	DEMAND.LOSS.NW	
[]:	POPDEN_RURAL float64	s might be aggregated.
[]:	2014 5 Minnesota intentional attack 2010 10 Minnesota severe weather 2012 6 Minnesota severe weather 2015 7 Minnesota 2015 7 Minnesota severe weather 2015 7 Minnesota 2015 2015 7 Minnesota 2015 20	nd away from 2011 seem to follow a bell-like shape.
[]:[Looking into Extent of Outages fig = outage['OUTAGE.DURATION'].hist(title='Outage Duration Distribution [in Minutes]')	Megawatt)')
[]:	• Note: the DEMAND.LOSS.MW column is missing about half its values. Reference the <i>Plot Null Counts of Columns with Missing Data</i> section fig = outage['CUSTOMERS.AFFECTED'].hist(title='Distribution of Customers Affected by Each Outage')	national Dateline. Index values of +0.5 or higher indicate El Niño. Values of -0.5 or lower indica
[]:	Niña." (https://www.climate.gov/news-features/understanding-climate/climate-variability-oceanic-ni%C3%B1o-index#:~:text=The%20ONI%20is%20the%20rolling,or%20low/lina." (https://www.climate.gov/news-features/understanding-climate/climate-variability-oceanic-ni%C3%B1o-index#:~:text=The%20ONI%20is%20the%20rolling,or%20low/lina. It seems as thought outages tend to be most frequent around the -0.3 anomaly level. This is strange as -0.3 indicates neither an El Niño nor a La Niña weather event, and it to the frequency of power outages? Does is it really relevant or was this distribution random? Could this distribution indicate that outages are more frequent during La Niña of the frequency of power outages. Duration: ', outages outages are more frequent during La Niña of the frequency of power outages. Duration: ', outages outages are more frequent during La Niña of the frequency of power outages. Duration: ', outages outages are more frequent during La Niña of the frequency of power outages. Duration: ', outages outages. Duration: ', out	I it's not at 0.0 (between the two). This plot brings up many questions. What relevance does -0 exercise (<= -0.5).
[]:	# Groups rows by state and aggregates by the count of each group. # Takes those counts and sorts them # Creates a horizontal bar plot of outage observation counts for each state fig = outage.groupby(by='U.SSTATE').count().OBS.sort_values().plot(kind='barh', width=890, height=1000, title='Outages Counts by State fig.write_html('./assets/outage-counts-state.html', include_plotlyjs='cdn') This plot shows the number of outages that were obseved in each state. It is not very interesting as it stands, the results speak for themselves. More outages happen in more adjusted for square footage. **Creates a series with the number of outages in each state outage_dist_by_state = outage.groupby(by='U.SSTATE').count().OBS.sort_values() # Creates a series of each state's mean population from 2000-2016 pop_by_state = outage.groupby(by='U.SSTATE').mean().POPULATION # Creates a horizontal bar plot displaying the number of outages per 10k people in each state fig = (outage_dist_by_state' / pop_by_state * 10_000).sort_values().plot(kind='barh', width=800, height=1000, title='Outages per 10k Peop fig.write_html('./assets/outages_per_capita_state.html', include_plotlyjs='cdn') Ah! Much more interesting. Our first bivariate plot! It seems as though Delaware is has by far the most outages per person. More than double DC (the second highest). Why Side Note: Texas and California (the top two of the previous plot) are middle of the pack. This is a good example for the importance of bivariate analysis. Also, note that the **Counts of outages in each NERC region* outages_by_NERC = outage.groupby(by='NERC.REGION').count().OBS.sort_values() **Counts of outages in each NERC region* outages_by_NERC = outage.groupby(by='NERC.REGION').count().OBS.sort_values()	hore populous states. It would be interesting to see the same plot adjusted for population and opple by State') thy might Delaware's number of outages be so disproportionately high given their population?
[]:	# Hbar plot on the results from the previous line fig = outages_by_NERC.plot(kind='barh', width=1100, height=500, title='Outages Counts by NERC Weather Region (North American Electric Refig.write_html('./assets/outage-counts-NERC.html', include_plotlyjs='cdn') fig One can see that, again the most populous weather regious seem to populate the top of the plot. Unsurprising ! # Finds the number of outages in each NERC region per 1M people NERC_region_pop = outage.groupby(bys='NERC.REGION', 'U.SSTATE'])\mean()reset_index().groupby('NERC.REGION'), 'D.SSTATE'])\mean()reset_index().groupby('NERC.REGION'), 'D.SSTATE'])\mean()reset_index().groupby('NERC.RE	
[]:	Lets look quickly at the regular weather regions to get a point of comparison. **Counts of outages in each weather region outages_by_NERC = outage.groupby(by='CLIMATE.REGION').count().OBS.sort_values()	th=800)
	It seems as though severe weather is to blame for most outages, but it is interesting to see that a substantial portion of outages are intentional attacks. It would be interesting also look into the number of intentional attacks that happen per capita (mean of POPULATION of each group of observations). # Finds the mean of the populations in each cause category group mean_pop_cause = outage.groupby(by='CAUSE.CATEGORY').mean()['POPULATION']	eaches much closer to the severe weather bar. This may indicate a relationship between per c age was observed, not the affected population, making this metric very shaky.
[]:	South Dakota 0 0 0 2 0 0 0 0 North Dakota 0 1 0 0 1 0 0 1 0 0	er system operability disruption :
[]:]: # normalizes the pivot table created above cause_by_state_norm = cause_by_state.apply(lambda x: x / cause_by_state.sum(axis=1)) # plots the normalized data fig = cause_by_state_norm.plot(kind='barh', width=800, height=1300, title='Distribution of Outage Causes for each State') fig.write_html('./assets/cause-state-norm.html', include_plotlyjs='cdn') fig This confirms our results from above; we are mostly seeing severe weather and intentional attack as the primary causes. Now, let's scale these values to understand the da]: # scales the pivot table by state population to give per capita view of the data cause_by_state_scaled = cause_by_state.apply(lambda x: x / pop_by_state * 10_000) # calculates the total of each row and adds it into the df as total for sorting cause_by_state_scaled = cause_by_state_scaled.assign(total = cause_by_state_scaled.sum(axis=1)) # sorts the df using total and then drops total from the df cause_by_state_scaled = cause_by_state_scaled.sort_values('total').drop(columns=['total'])	data per capita.
[]:	<pre># plots the scaled pivot table fig = cause_by_state_scaled.plot(kind='barh', width=800, height=1300, title='Outages per 10k People by Cause for each State') fig.write_html('./assets/cause-state-per-capita.html', include_plotlyjs='cdn') fig Ah! This plot immediately gives us more context for Delaware's outlandish number of outages per capita. There are a disproportionate number of intentional attack. Might states with high rates of intentional attack have more outages per person overall? Let's explore this possible trend.]: # adds the total back in in case one wants to sort the data descending df = cause_by_state_scaled.assign(</pre>	outages in Delaware. In fact, many of the higher ranking states tend to have high instances o
[]:	fig = px.scatter(data_frame=df, x='total', y='intentional attack', color=df.index, color_discrete_sequence=px.colors.qualitative.G10, title = 'Intentional Attack Outages (per 10k people) compared to total Outages (per 10k people)') fig.write_html('./assets/intentional-attack-scatter-total-outages.html', include_plotlyjs='cdn') fig Correlation between total outages per capidta and severe: 0.9542374010674626 It seems like we're onto something. Hoverver, there could still be confounders. Below! created the same graph but removed Delaware since it seems to be an outlier.]: # adds the total back in in case one wants to sort the data descending and removes Delaware from the data df = cause_by_state_scaled.assign(total = cause_by_state_scaled.sum(axis=1)).sort_values(by='total', ascending=True).drop(['Delaware']) # Calculates the correlation and prints it print('Correlation between total outages per capidta and severe:', df['total'].corr(df['intentional attack'])) # Scatter plots the intentional attack by total	
[]:	fig = px.scatter(data_frame=df, x='total', y='intentional attack', color=df.index, color_discrete_sequence=px.colors.qualitative.G10, title = 'Intentional Attack Outages (per 10k people) compared to total Outages (per 10k people) [Delaware Removed]') fig.write_html('./assets/intentional-attack-scatter-total-no-delaware.html', include_plotlyjs='cdn') fig Correlation between total outages per capidta and severe: 0.8246248841759625 We can see that the correlation still holds, albeit, to a lesser degree. Assessment of Missingness Plot Null Counts of Columns with Missing Data]: # counts all null values in each column outages null_counts = outage.isna().sum().sort_values()	
	# plots hbar of missing value conts for each column with missing values fig = null_counts[null_counts > 0].plot(
[]: [Alt Hypothesis: The distribution of CAUSE.CATEGORY.DETAIL is different when CAUSE.CATEGORY name is missing as opposed to when it is not missing **Variables* for columns we are testing to enable hot switching indep_var = 'CAUSE.CATEGORY' q_var = 'CAUSE.CATEGORY' q_var = 'CAUSE.CATEGORY.DETAIL' **generate pivot table to display missingness of q_var with relation to indep_var	
	<pre>fig = df_indep.plot(kind='barh', title='Observed Distribution of Cause Category Detail Conditional on Cause Category', barmode='group', height=500) fig.write_html('./assets/observed-dist-dependent.html', include_plotlyjs='cdn') fig]: print(df_indep.to_markdown()) CAUSE.CATEGORY</pre>	
[]:	<pre>Simulation]: # number of times to repeat permutations and calculate TVD n_repetitions = 500 shuffled = outage.copy() tvds = []</pre>	
	<pre>for _ in range(n_repetitions): # Create permutation shuffled[indep_var] = np.random.permutation(shuffled[indep_var]) # Computing and storing the TVD. pivoted = (shuffled .assign(missing = shuffled[q_var].isna()) .pivot_table(index=indep_var, columns='missing', aggfunc='size') .apply(lambda x: x / x.sum())) tvd = pivoted.diff(axis=1).iloc[:, -1].abs().sum() / 2 tvds.append(tvd)]: # plotting the imerical distribution of the TVD fig = px.histogram(pd.DataFrame(tvds), x=0, nbins=50, histnorm='probability', title='Empirical Distribution of the TVD')</pre>	
[]:	<pre>fig.add_vline(x=observed_tvd, line_color='red') fig.add_annotation(text=f'span style="color:red">span style=</pre>	
[]:	Testing Independent Null Hypothesis: The distribution of CAUSE.CATEGORY.DETAIL is the same when CAUSE.CATEGORY name is missing and when it is not missing Alt Hypothesis: The distribution of CAUSE.CATEGORY.DETAIL is different when CAUSE.CATEGORY name is missing as opposed to when it is not missing indep_var = 'YEAR'	
[]:[df_inde*Darh', title*'Observed Distribution of Hurricane Names Conditional on Cause Category Detail', barmode*'group', height=1100) imissing False True YEAR 2000 0.015992 0.019108 2001 0.001881 0.027601 2002 0.012230 0.008493 2003 0.036689 0.014862 2004 0.057385 0.021231 2005 0.043274 0.019108 2006 0.046096 0.038217 2007 0.042333 0.023355 2008 0.070555 0.076433 2009 0.048918 0.055202	
	2010 0.052681 0.106157 2011 0.145814 0.242038 2012 0.139229 0.055202 2013 0.111947 0.072187 2014 0.085607 0.044586 2015 0.052681 0.133758 2016 0.036689 0.042463]: observed_tvd = df_indep.diff(axis=1).iloc[:, -1].abs().sum() / 2 observed_tvd]: 0.277548419826913	
	<pre>Simulation in_repetitions = 500 shuffled = outage.copy() tvds = [] for _ in range(n_repetitions): shuffled[indep_var] = np.random.permutation(shuffled[indep_var]) # Computing and storing the TVD. pivoted = (shuffled .assign(missing = shuffled[q_var].isna()) .pivot_table(index=indep_var, columns='missing', aggfunc='size') .apply(lambda x: x / x.sum())) tvd = pivoted.diff(axis=1).iloc[:, -1].abs().sum() / 2 tvds.append(tvd) fig = px.histogram(pd.DataFrame(tvds), x=0, nbins=50, histnorm='probability',</pre>	
	The p-value is, again less than our significance level of 0.05. This means we reject the null hypothesis; making CAUSE CATEGORY,DETAIL dependent on YEAR. I actually move on to the next section. Hypothesis Testing Null: number of outages by state per capita comes from the same distribution as the proportion of outages that are caused by intentional attack by state Alt: number of outages by state per capita and the proportion of outages that are caused by intentional attack by state come from different distributions Significance level: 0.05 # number of outages by state per capita x = outage.groupby(by='U.SSTATE').count().08S / outage.groupby(by='U.SSTATE').mean().POPULATION # proportion of outages that are caused by intentional attack by state cause_by_state = pd.pivot_table(outage, columns=['CAUSE.CATEGORY'], index=['U.SSTATE'], values='OBS', aggfunc='count').fillna(0) y = cause_by_state': severe weather'] / cause_by_state.sum(axis=1) # put X and Y into a df df = pd.DataFrame().assign(out_per_cap=X, prop_attack=Y) # normalize the df df = df / df.sum(axis=0) # plot the distributions fig = df.plot(kind='barh', height=1000, barmode='group', title='Observed Distributions of Outages by State Per Capita and Prop Outages C fig. write_html('./assets/hyp-test-observed.html', include_plotlyjs='cdn') fig	
[]:	Fig	
	<pre>n_repetitions = 500 shuffled = outage.copy() tvds = [] for _ in range(n_repetitions): shuffled['U.SSTATE'] = np.random.permutation(shuffled['U.SSTATE']) # number of outages by state per capita X = shuffled.groupby(by='U.SSTATE').count().0BS / shuffled.groupby(by='U.SSTATE').mean().POPULATION # proportion of outages that are caused by intentional attack by state cause_by_state = pd.pivot_table(shuffled, columns=['CAUSE.CATEGORY'], index=['U.SSTATE'], values='OBS', aggfunc='count').fillna(0) Y = cause_by_state['severe weather'] / cause_by_state.sum(axis=1) # put X and Y into a df df = pd.DataFrame().assign(out_per_cap=X, prop_attack=Y)</pre>	

Analyzing Power Outages