	<pre>import pandas as pd tweets = pd.read_csv('data/trumptweets-1515775693.tweets.csv', low_memory = False) tweets.head() status_id created_at user_id screen_name text source display_text_width reply_to_status_id 2009-05-</pre>
	1 x9273573134835712 2010-11-29 15:52:46 x25073877 realDonaldTrump realDonaldTrump that Congratulations to Evan Lysacek for being nomi Twitter Web Client Twitter Web Client NaN
	2010-10- 28 x29014512646 28 x25073877 realDonaldTrump Web talked abou Twitter Web Client Tomorrow night's episode Twitter 42010-11-24 17:20:54 x25073877 realDonaldTrump of The Apprentice del
5	4 x5775731054 2009-11-16 21:06:10 x25073877 realDonaldTrump Partners with TV1 on New Reality Twitter Twitter Web Client Towns x 68 columns
	<pre>/hat is the range of dates on this dataset? tweets[['created_at']].sort_values('created_at').iloc[[0, -1]])</pre>
Н	our analysis will pertain to Donald Trump tweets from 2009 to 2018. Now often does Donald Trump retweet vs. spit out a tweet of his own? Len (tweets [tweets.is_retweet]) / len (tweets) * 100 L.5627855967830377
H th	The only retweets 1.6% of the time. This means we have a good enough dataset to analyze his own writing. Let's omit retweets from the dataset. We're also only worried about the text, so let's just keep the text and the date, and drop everything else. tweets = tweets[~tweets.is_retweet][['created_at', 'text']] tweets.head() text
3	2010-11-29 15:52:46 Congratulations to Evan Lysacek for being nomi 2010-10-28 18:53:40 I was on The View this morning. We talked abou 2010-11-24 17:20:54 Tomorrow night's episode of The Apprentice del 2009-11-16 21:06:10 Donald Trump Partners with TV1 on New Reality
If be	you look closely, even if we take out the retweets, there are some weird tweets (looks like not a first-person type of writing). Se elow for examples. tweets.iloc[0].text 'Read a great interview with Donald Trump that appeared in The New York Times Magazine: http://tinyurl.com
k at	tweets.iloc[4].text 'Donald Trump Partners with TV1 on New Reality Series Entitled, Omarosa's Ultimate Merger: http://tinyurl.c5m3lc" fter some research, it is clear that even the non-retweets may not come from himself and to get to the heart of the question, want to only get the tweets that represent Donald Trump as himself, no one else. Let's also avoid all the media links. Trump also
SO	ometimes retweets other accounts using "RT @account" method. Let's avoid this as well. tweets = tweets[~(tweets.text.str.contains('http') tweets.text.str.contains('RT @'))] print(f'we have {tweets.shape[0]} tweets after preprocessing') tweets.head(2) we have 24669 tweets after preprocessing created_at text
2	2010-11-29 15:52:46 Congratulations to Evan Lysacek for being nomi 2010-10-28 18:53:40 I was on The View this morning. We talked abou 20kay, now we're ready.
y b	Question 1 What distribution does Donald Trump's word frequency follow? For now, just plot, and find a distribution that might fit this. No arameters are needed yet, just a visual approximation. Answer 1 We first have to tokenize the text. This means to split the text into individual word(s). We are going to use unigram-approach, whi
m	means that each word is one word. Let's also get rid of punctuations. # importing an useful tokenizer from nltk import word_tokenize # this is regex! import re
8	<pre># we split tweets into individual words while also getting rid of anything at is not a letter or space. tweets['words'] = tweets['text'].apply(lambda x: word_tokenize(re.sub(r'[^\w\s]', '', x))) tweets[['words']].head(1).words.values array([list(['Congratulations', 'to', 'Evan', 'Lysacek', 'for', 'being', 'nominated', 'SI', 'sportsman', 'the', 'year', 'Hes', 'a', 'great', 'guy', 'and', 'he', 'has', 'my', 'vote', 'EvanForSI'])],</pre>
4	<pre>from collections import Counter all_words = [word for templist in tweets[['words']].values.tolist() for words in templist for word in worden(all_words) 421566</pre> are more than 400,000 words (not unique) in his tweets from 2000 to 2017
[word_frequency = Counter(all_words) sorted(word_frequency, key=word_frequency.get, reverse=True)[:10] ['the', 'to', 'realDonaldTrump', 'a', 'is', 'and', 'I', 'of', 'you', 'in'] makes sense that the top 10 words include "the", "to", his own Twitter name, and "a".
	<pre># plotting tools import matplotlib.pyplot as plt import seaborn as sns import numpy as np sns.set_style("white") word_frequency = Counter(all_words)</pre>
	<pre>word_frequency = word_frequency.most_common() frequency_values = [x[1] for x in word_frequency] fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.set(font_scale = 2) sns.scatterplot(x = np.arange(0, len(word_frequency))[:500], y = frequency_values[:500], ax= axes) plt.title('The first 500 words')</pre>
	axes.set_xlabel('Most common word index') axes.set_ylabel('Word frequency') plt.show() The first 500 words 12500 10000
Mord froguesia	7500 5000 2500
Т	0 100 200 300 400 500 Most common word index his fits the Zipf's law. $f(k;s,N) = \frac{\frac{1}{k^s}}{\sum_{1}^{N} \frac{1}{n^s}}$
C	where k is the rank, N is the number of elements, and s is the parameter that characterizes the distribution. $oldsymbol{Q}$ uestion 2
ра Д	obviously, we're going to need a more rigorous way to find the distribution. First, find the maximum likelihood estimator for the arameter of the Zipf distribution. Zipf distribution only requires 1 parameter, s , the shape parameter. Answer 2 he probability mass function using Zipf distribution (from equation 1) can be used to calculate the log-likelihood $L(s)$. $L(s) = Log(\frac{1}{s})$
B ₀	$L(s) = Log(rac{1}{k^s*HarmonicNumber(N,s)})$ ecause of the Harmonic Numbers, it will be extremely difficult to find an analytical solution for the estimator. Therefore, we use ptimization methods to find an estimate. from scipy.optimize import minimize_scalar from scipy.stats import zipf, norm def neg log likelihood(s):
	<pre>def neg_log_likelihood(s): N = sum(frequency_values) probs = np.log(frequency_values**(-1 * s) / np.sum(np.arange(1, N+1)**(-1 * s))) summed = sum(probs * frequency_values) return -1 * summed frequency_values = np.array(frequency_values) s_best = minimize_scalar(neg_log_likelihood, [0, 6.0])</pre>
1 TI	<pre>s_best = minimize_scalar(neg_log_likelihood, [0, 6.0]) print(s_best.x) 1.0472055615708566 he best estimate for the s parameter was 1.0472. sns.set_style('white') fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50)</pre>
	<pre>sns.set(font_scale = 2) plt.yscale('log') plt.xscale('log') sns.scatterplot(x = np.arange(1, len(word_frequency)+1), y =frequency_values, ax= axes, label = 'Twitter plt.title('Log-Log Plot against MVUE') axes.set_xlabel('Most common word index') axes.set_ylabel('Word frequency')</pre>
	<pre># here, we are fitting the estimate we found using Scipy's built-in zipf PMF function. fitted = np.array([zipf.pmf(n, s_best.x) for n in np.arange(1, len(frequency_values) + 1)]) * sum(frequency) p1 = sns.scatterplot(x = np.arange(1, len(word_frequency) + 1), y = fitted, ax = axes, label = 'Shape Est p1.legend(frameon = True) p1t.tight_layout() p1t.show()</pre>
rd froguopou	Log-Log Plot against MVUE 10 Twitter Data Shape Estimator
	Twitter Data Shape Estimator 10 10 10 10 10 10 10 Most common word index he estimate we found is a pretty good estimate to the 400,000 tweet data.
U	Question 3 Ising the MLE from question 2, bootstrap 10000 times and calculate this estimator 1000 times. What is the mean? the variance? Tooes this match with analytical solution?
P _y	Answer 3 ython's Scikit-Learn has a great bootstrap API that we will use. (We also could have used Pandas resample feature, which actuated ould have been easier since our data is already in dataframes, but oh well. too lazy to change it.) from sklearn.utils import resample
1	<pre>list_to_sample_from = [x[0] for x in tweets[['words']].values] len(resample(list_to_sample_from, replace=True, n_samples=100, random_state=42)) loo he following step will take a while since we are optimizing to find estimator for s 10000 times. We'll use a time-tracking library to</pre>
	<pre>eep track of the progress. from tqdm.notebook import tqdm tqdm() Dit [00:00, ?it/s]</pre>
	<pre>global word_frequency_temp, frequency_values_temp, temp_words def neg_log_likelihood_temp(s): N = sum(frequency_values_temp) probs = np.log(frequency_values_temp**(-1 * s) / np.sum(np.arange(1, N+1)**(-1 * s))) summed = sum(probs * frequency_values_temp) return -1 * summed bootsamples = []</pre>
	<pre>for i in tqdm(range(0, 10000)): sample = resample(list_to_sample_from, replace=True, n_samples=1000) temp_words = [word for words in sample for word in words] word_frequency_temp = Counter(temp_words) word_frequency_temp = word_frequency_temp.most_common() frequency_values temp = [x[1] for x in word_frequency_temp]</pre>
	<pre>frequency_values_temp = np.array(frequency_values_temp) s_best = minimize_scalar(neg_log_likelihood_temp, [0, 6.0]) del word_frequency_temp, frequency_values_temp, temp_words bootsamples.append(s_best.x)</pre>
	<pre>et's plot our estimators. sns.set_style('white') fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.set(font_scale = 2) sns.histplot(bootsamples, ax = axes)</pre>
	plt.title('Bootstrap of the shape estimators') plt.tight_layout() plt.show() Bootstrap of the shape estimators
441.00	400
	1.220 1.225 1.230 1.235 1.240 print (f'The mean of the estimator from the bootstrap population is: {np.round(np.mean(bootsamples), 3)}') print (f'The variance of the estimator from the bootstrap population is: {np.round(np.var(bootsamples), 5)} The mean of the estimator from the bootstrap population is: 1.232
N G U	Indice that when we only look at 1000 tweets at a time, it does not come close to knowing the full picture of 400,000 tweets. Question 4 Using our hypothesis testing knowledge from class, determine if Donald Trump's tweets follow the distribution with the shape stimator from the bootstrap in question 3. You will need to set up a null/alternative hypotheses and a log-likelihood ratio statistic
A	lot the p-values from the 1000 bootstraps 100 times. Answer 4 he null hypothesis is that the s parameter from the bootstrapped distribution fits the overall, underlying distribution of 400,000 weets. The alternative hypothesis is that the s parameter does not fit the underlying distribution.
E w	$H_o: s=s_{bootstrap}$ $H_a: s eq s_{bootstrap}$ $H_a: s \neq s_{bootstra$
L	$f(heta_o)=\Pi_i^nf(s_{bootstrap},i)$ where f is the Zipf probability mass function. $f(heta_a)=\Pi_i^nf(s_{data},i)$ where s_{data} represents the fitted s parameter for the sample, and not the bootstrapped parameter.
si	$D=\sum_{i}^{n}-2log(rac{f(s_{bootstrap},i)}{f(s_{data},i)})pprox \chi^{2}(n)$ ince the null hypothesis has 1 parameter and the alternative hypothesis has 1 parameter (1-1 = 0). from scipy.stats import chi2 bootstrap_s = 1.232 test statistics = [1]
	<pre>bootstrap_s = 1.232 test_statistics = [] p_values = [] global word_frequency_temp, frequency_values_temp, temp_words def neg_log_likelihood_temp(s): N = sum(frequency_values_temp) probs = np.log(frequency_values_temp**(-1 * s) / np.sum(np.arange(1, N+1)**(-1 * s))) summed = sum(probs * frequency_values_temp) return -1 * summed</pre>
	<pre>for i in tqdm(range(0, 100)): sample = resample(list_to_sample_from, replace=True, n_samples=1000) temp_words = [word for words in sample for word in words] word_frequency_temp = Counter(temp_words) word_frequency_temp = word_frequency_temp.most_common() frequency_values_temp = [x[1] for x in word_frequency_temp]</pre>
	<pre>frequency_values_temp = np.array(frequency_values_temp) s_best = minimize_scalar(neg_log_likelihood_temp, [0, 6.0]) log0 = np.log(np.array([zipf.pmf(n, bootstrap_s) for n in np.arange(1, len(frequency_values_temp) + 1 log1 = np.log(np.array([zipf.pmf(n, s_best.x) for n in np.arange(1, len(frequency_values_temp) + 1)]) D = np.sum(-2*(log0 - log1))</pre>
	<pre>test_statistics.append(D) p_values.append(1 - chi2.cdf(D, 1)) del word_frequency_temp, frequency_values_temp, temp_words sns.set_style('white')</pre>
	<pre>fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.set(font_scale = 2) sns.histplot(p_values, ax = axes) plt.title('Bootstrap of the p values.') plt.tight_layout()</pre>
	Bootstrap of the p values. 40
2	0.0 0.2 0.4 0.6 0.8 1.0
se ha	seems half of the p-values are below 0.05, the alpha level, while the other half are above. I believe this is flawed because we shee a descending count of p-values as it deviates from the alpha level. Because we gave a 1000 resampling parameter, we shoul ave seen the bootstrap distribution fit the overall, underlying population. Question 5
di ar	low that we have a distribution that describes Donald Trump's "speaking" (writing) quality, get the top 20% of words in this istribution . By doing so, we avoid very common words (called stop words) such as "the", "for", "and", but we also avoid words re very rarely used and therefore do not represent the level of speaker well. Give us 10 words that fall in this section. Answer 5 We first put the words in a dataframe and use Panda's percentile function to get the top 20%.
	<pre>all_words = [word for templist in tweets[['words']].values.tolist() for words in templist for word in word_frequency = Counter(all_words) word_frequency = word_frequency.most_common() frequency_values = [x[1] for x in word_frequency] word_frequency = [x[0] for x in word_frequency]</pre>
(<pre>percentile = pd.DataFrame({'words':word_frequency, 'frequency':frequency_values}) percentile.head() words frequency the 12968 to 9506</pre>
3	1 to 9506 2 realDonaldTrump 7663 3 a 7257 4 is 6753 percentile.quantile([1, 0.8])
	frequency 1.0 12968.0 0.8 3.0 he 80th percentile for frequency values is 3. We'll just grab the data points with frequency higher than or equal to 3.
	he 80th percentile for frequency values is 3. We'll just grab the data points with frequency higher than or equal to 3. percentile = percentile.loc[percentile.frequency >= 3] percentile.head(10) words frequency the 12968
3	1 to 9506 2 realDonaldTrump 7663 3 a 7257 4 is 6753
7	5 and 6611 6 I 5678 7 of 5648 8 you 5597 9 in 5378
G	Question 6 In an an addition 15378 In an an addition 35378 In an an an an addition and a state of the st
W	Answer 6 We use a Python library called py-readability-metrics (https://github.com/cdimascio/py-readability-metrics). from readability import `Readability ecause one tweet is too short to be tested, we combine 10 tweets together (considering one tweet to be around 10 words).
	readinglevel = (tweets.groupby(tweets.index // 10).agg(' '.join)) readinglevel.head(2) created_at text 2010-11-29 15:52:46 2010-10-28 18:53:40 2010-1 Congratulations to Evan Lysacek for being nomi
	<pre>1 2010-12-10 14:42:15 2010-12-16 14:30:23 2010-1 Congratulations to Brandy as our new Apprentic grade_levels = [] grade_scores = [] for index, row in readinglevel.iterrows(): # we put a try statement in case there are still tweets than are <100 words. try: r = Readability(row['text'])</pre>
	<pre>try: r = Readability(row['text']) fk = r.flesch_kincaid() grade_scores.append(fk.score) grade_levels.append(fk.grade_level) except: pass grade_levels = [int(x) for x in grade_levels]</pre>
	<pre>sns.set_style('white') fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.histplot(sorted(grade_levels, reverse = False), ax = axes, bins = 30) plt.title('Grade Levels of Trump\'s Tweets') sns.set(font_scale = 2) axes.set_xlim([0, 15]) plt.tight_layout() plt.show()</pre>
	Grade Levels of Trump's Tweets 800 600
4	200 0.0 2.5 5.0 7.5 10.0 12.5 15.0
7 D	np.mean(grade_levels) 7.338757396449704 conald Trump's average reading level is 7th grade, which is around 12 years old.
C	Question 7 Falculate and present a 95% confidence interval of the grade results per tweet. Answer 7 We make use of bootstrapping method again to get the sample mean and use a normal distribution citing the Central Limit Theor
	<pre>Ve make use of bootstrapping method again to get the sample mean and use a normal distribution citing the Central Limit Theor def bootstrap(): temp = tweets.sample(1000) readinglevel = (temp.groupby(temp.index // 10).agg(' '.join)) for index, row in readinglevel.iterrows(): # we put a try statement in case there are still tweets than are <100 words. try: r = Readability(row['text'])</pre>
	<pre>r = Readability(row['text']) fk = r.flesch_kincaid() grade_scores.append(int(fk.score)) grade_levels.append(int(fk.grade_level)) except: pass return np.mean(grade_levels)</pre>
	<pre>return np.mean(grade_levels) averaged_grades = [] for i in tqdm(range(0, 1000)): averaged_grades.append(bootstrap())</pre>
	sns set stylo(!white!)
	<pre>sns.set_style('white') fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.histplot(sorted(averaged_grades, reverse = False), ax = axes, bins = 20) plt.title('Grade Levels of Trump\'s Tweets') sns.set(font_scale = 2) plt.tight_layout() plt.show()</pre> Grade Levels of Trump's Tweets
	<pre>fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.histplot(sorted(averaged_grades, reverse = False), ax = axes, bins = 20) plt.title('Grade Levels of Trump\'s Tweets') sns.set(font_scale = 2) plt.tight_layout()</pre>
1	fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.histplot(sorted(averaged_grades, reverse = False), ax = axes, bins = 20) plt.title('Grade Levels of Trump\'s Tweets') sns.set(font_scale = 2) plt.tight_layout() plt.show() Grade Levels of Trump's Tweets 6 6 6 7 7 7 7 7 7 7 7 7 7
T	fig, axes = plt.subplots(1, 1, figsize=(8, 5), dpi = 50) sns.histplot(sorted(averaged_grades, reverse = False), ax = axes, bins = 20) plt.title('Grade_Levels of Trump\'s Tweets') sns.set(font_scale = 2) plt.tight_layout() plt.show() Grade Levels of Trump's Tweets 150 6 100 50
TI	fig, $ax = plt.subplots(1, 1, figsize=(8, 5), dpi = 50)$ sns.histplot(sorted(averaged_grades, reverse = False), $ax = axes$, $bins = 20)$ plt.title('Grade Levels of Trump\'s Tweets') sns.set(font_scale = 2) plt.tight_layout() plt.show() Grade Levels of Trump's Tweets Grade Levels of Trump's Tweets ### Additional Content of the Content of Trump's Tweets ### Additional Content of Trump's Tweets ###