


```
#URL of the webpage to scrape
url = "https://www.baseball-reference.com/teams/CHS/attend.shtml"

#Read the HTML table into a list of DataFrames
tables = pd.read_html(url)

#Select the table of interest (index 0 in this case)
table = tables[0]

#Save the DataFrame as a CSV file
file_path = "C:/Users/andre/OneDrive/Project/MLB_data/white_sox_statistics.csv"
table.to_csv(file_path, index=False, header=True)

print("Data saved to", file_path)

Data saved to C:/Users/andre/OneDrive/Project/MLB_data/white_sox_statistics.csv

In [116]:
path = "C:/Users/andre/OneDrive/Project/MLB_data/"
white_sox_df = pd.DataFrame(pd.read_csv(path + 'white_sox_statistics.csv'))
white_sox_df.head()
```

	Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium
0	2023	Chicago White Sox	AL Central	33	45	4	NaN	758109.0	19439.0	11th of 15	\$157,571,666	104	104	Guaranteed Rate Field
1	2022	Chicago White Sox	AL Central	81	81	2	NaN	2009359.0	24807.0	8th of 15	\$163,958,334	103	103	Guaranteed Rate Field
2	2021	Chicago White Sox	AL Central	93	69	1	Lost ALDS (3-1)	1596385.0	19708.0	5th of 15	\$115,546,333	101	102	Guaranteed Rate Field
3	2020	Chicago White Sox	AL Central	35	25	3	Lost ALWC (2-1)	NaN	NaN	3rd of 15	\$119,066,333	100	100	Guaranteed Rate Field
4	2019	Chicago White Sox	AL Central	72	89	3	NaN	1649775.0	20622.0	11th of 15	\$80,846,333	99	98	Guaranteed Rate Field


```
In [117]:
white_sox_df = white_sox_df[white_sox_df['Year'] > 2004]
white_sox_df.head()
```

	Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium
0	2023	Chicago White Sox	AL Central	33	45	4	NaN	758109.0	19439.0	11th of 15	\$157,571,666	104	104	Guaranteed Rate Field
1	2022	Chicago White Sox	AL Central	81	81	2	NaN	2009359.0	24807.0	8th of 15	\$163,958,334	103	103	Guaranteed Rate Field
2	2021	Chicago White Sox	AL Central	93	69	1	Lost ALDS (3-1)	1596385.0	19708.0	5th of 15	\$115,546,333	101	102	Guaranteed Rate Field
3	2020	Chicago White Sox	AL Central	35	25	3	Lost ALWC (2-1)	NaN	NaN	3rd of 15	\$119,066,333	100	100	Guaranteed Rate Field
4	2019	Chicago White Sox	AL Central	72	89	3	NaN	1649775.0	20622.0	11th of 15	\$80,846,333	99	98	Guaranteed Rate Field
5	2018	Chicago White Sox	AL Central	62	100	4	NaN	1608817.0	19862.0	12th of 15	\$75,092,000	99	98	Guaranteed Rate Field
6	2017	Chicago White Sox	AL Central	67	95	4	NaN	1629470.0	20117.0	13th of 15	\$97,842,000	98	97	Guaranteed Rate Field
7	2016	Chicago White Sox	AL Central	78	84	4	NaN	1746293.0	21559.0	12th of 15	\$113,416,000	96	96	U.S. Cellular Field
8	2015	Chicago White Sox	AL Central	76	86	4	NaN	175810.0	21677.0	13th of 15	\$112,889,700	97	96	U.S. Cellular Field
9	2014	Chicago White Sox	AL Central	73	89	4	NaN	1650821.0	20381.0	13th of 15	\$87,475,500	98	97	U.S. Cellular Field
10	2013	Chicago White Sox	AL Central	63	99	5	NaN	1768413.0	21832.0	10th of 15	\$81,041,900	105	104	U.S. Cellular Field
11	2012	Chicago White Sox	AL Central	85	77	2	NaN	1965955.0	24271.0	9th of 14	\$116,208,000	104	104	U.S. Cellular Field
12	2011	Chicago White Sox	AL Central	79	83	3	NaN	2001117.0	24705.0	7th of 14	\$127,789,000	106	106	U.S. Cellular Field
13	2010	Chicago White Sox	AL Central	81	84	2	NaN	2194378.0	27209.0	7th of 14	\$107,195,000	103	103	U.S. Cellular Field
14	2009	Chicago White Sox	AL Central	79	83	3	NaN	2284163.0	28001.0	6th of 14	\$101,081,000	105	105	U.S. Cellular Field
15	2008	Chicago White Sox	AL Central	89	74	1	Lost ALDS (3-1)	2500640.0	30496.0	5th of 14	\$121,189,332	105	104	U.S. Cellular Field
16	2007	Chicago White Sox	AL Central	72	90	4	NaN	2684395.0	33141.0	5th of 14	\$108,671,833	104	104	U.S. Cellular Field
17	2006	Chicago White Sox	AL Central	90	72	3	NaN	2957414.0	36511.0	3rd of 14	\$102,750,667	104	104	U.S. Cellular Field
18	2005	Chicago White Sox	AL Central	99	63	1	Won WS (4-0)	2342833.0	28924.0	7th of 14	\$75,178,000	103	103	U.S. Cellular Field

```
In [118]:
#Remove the dollar sign from the 'Estimated Payroll' column:
white_sox_df['Est. Payroll'] = white_sox_df['Est. Payroll'].replace({'$':''}, regex = True)
white_sox_df['Est. Payroll'] = white_sox_df['Est. Payroll'].replace({' ','.'}, regex = True)
white_sox_df.head()
```

```
Out[118]:
Year      Tm      Lg  W  L  Finish  Playoffs  Attendance  Attend/G  Rank  Est. Payroll  PPF  BPF  Stadium
0  2023  Chicago White Sox  AL Central  33  45  4      NaN      758109.0  19439.0  11th of 15  157571666  104  104  Guaranteed Rate Field
1  2022  Chicago White Sox  AL Central  81  81  2      NaN      2009359.0  24807.0  8th of 15  163958334  103  103  Guaranteed Rate Field
2  2021  Chicago White Sox  AL Central  93  69  1  Lost ALDS (3-1)  1596385.0  19708.0  5th of 15  115546333  101  102  Guaranteed Rate Field
3  2020  Chicago White Sox  AL Central  35  25  3  Lost ALWC (2-1)  NaN      NaN      3rd of 15  119066333  100  100  Guaranteed Rate Field
4  2019  Chicago White Sox  AL Central  72  89  3      NaN      1649775.0  20622.0  11th of 15  80846333  99  98  Guaranteed Rate Field
```

```
In [119]:
white_sox_df['Est. Payroll'] = white_sox_df['Est. Payroll'].astype(float)
white_sox_df.dtypes
```

```
Out[119]:
Year      int64
Tm        object
Lg        object
W         int64
L         int64
Finish    int64
Playoffs  object
Attendance float64
Attend/G  float64
Rank      object
Est. Payroll float64
PPF       int64
BPF       int64
Stadium   object
dtype: object
```

Create binary feature that defines a successful year as wins greater than or equal to 50% of total games played

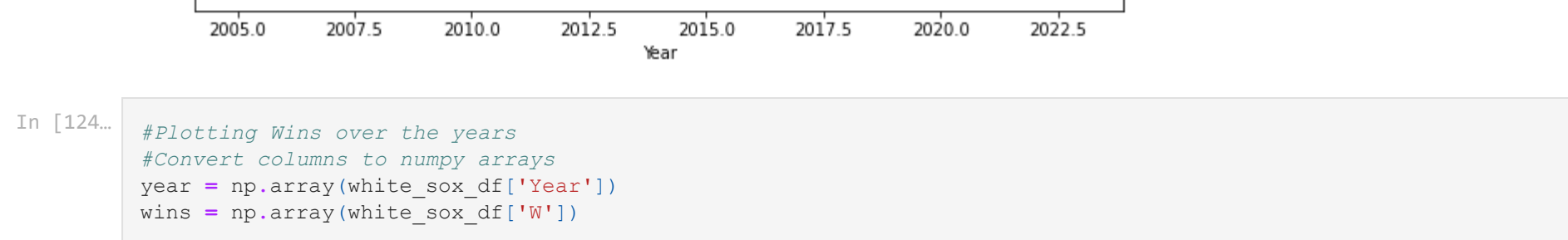
```
In [120]:
#Add a column to define a binary output that has a 1 if the team has a winning record for the year, and 0 if white_sox_df['success'] = (white_sox_df['W'] / (white_sox_df['W'] + white_sox_df['L'])) >= 0.5).astype(int)
white_sox_df.head()
```

	Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium	success
0	2023	Chicago White Sox	AL Central	33	45	4	NaN	758109.0	19439.0	11th of 15	157571666	104	104	Guaranteed Rate Field	0
1	2022	Chicago White Sox	AL Central	81	81	2	NaN	2009359.0	24807.0	8th of 15	163958334	103	103	Guaranteed Rate Field	1
2	2021	Chicago White Sox	AL Central	93	69	1	Lost ALDS (3-1)	1596385.0	19708.0	5th of 15	115546333	101	102	Guaranteed Rate Field	1
3	2020	Chicago White Sox	AL Central	35	25	3	Lost ALWC (2-1)	NaN	NaN	3rd of 15	119066333.0	100	100	Guaranteed Rate Field	1
4	2019	Chicago White Sox	AL Central	72	89	3	NaN	1649775.0	20622.0	11th of 15	80846333.0	99	98	Guaranteed Rate Field	0

```
In [121]:
#Count the number of successes and failures
success_counts = white_sox_df['success'].value_counts()

#Create a bar plot
plt.bar(success_counts.index, success_counts.values)
plt.xticks(success_counts.index, ['Failure', 'Success'])
plt.xlabel('Success')
plt.ylabel('Count')
plt.title('Success Counts')

#Display the plot
plt.show()
```



```
In [122]:
white_sox_df.describe()
```

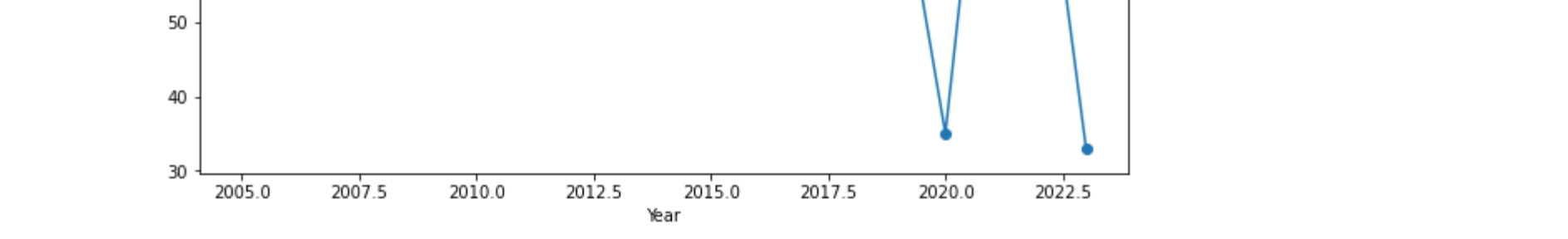
	Year	W	L	Finish	Attendance	Attend/G	Est. Payroll	PPF	BPF	success
count	19.000000	19.000000	19.000000	19.000000	1.800000e+01	18.000000	1.900000e+01	19.000000	19.000000	19.000000
mean	2014.000000	74.421053	77.789474	3.000000	1.959231e+06	24630.166667	1.087894e+08	101.789474	101.473684	0.421053
std	5.627314	17.359966	18.304762	1.20185	4.997730e+05	5086.213237	2.436409e+07	3.137213	3.372576	0.502757
min	2005.000000	33.000000	25.000000	1.000000	7.581090e+05	19439.000000	7.509200e+07	99.000000	96.000000	0.000000
25%	2009.500000	69.500000	73.000000	2.000000	1.650303e+06	20441.250000	9.26875e+07	99.000000	98.000000	0.000000
50%	2014.000000	83.000000	83.000000	3.000000	1.867184e+06	23051.500000	1.086718e+08	103.000000	103.000000	0.000000
75%	2018.500000	86.500000	89.000000	4.000000	2.261771e+06	27922.750000	1.186372e+08	104.000000	104.000000	1.000000
max	2023.000000	99.000000	100.000000	5.000000	2.957414e+06	36511.000000	1.635958e+08	106.000000	106.000000	1.000000

```
In [123]:
#Plotting
#Convert columns to numpy arrays
year = np.array(white_sox_df['Year'])
attendance = np.array(white_sox_df['Attendance'])

plt.figure(figsize=(10,6))

#Plotting Attendance over the years
plt.plot(year, attendance, marker='o')
plt.xlabel('Year')
plt.ylabel('Attendance')
plt.title('Attendance of Chicago White Sox over the Years')

#Display the plot
plt.show()
```

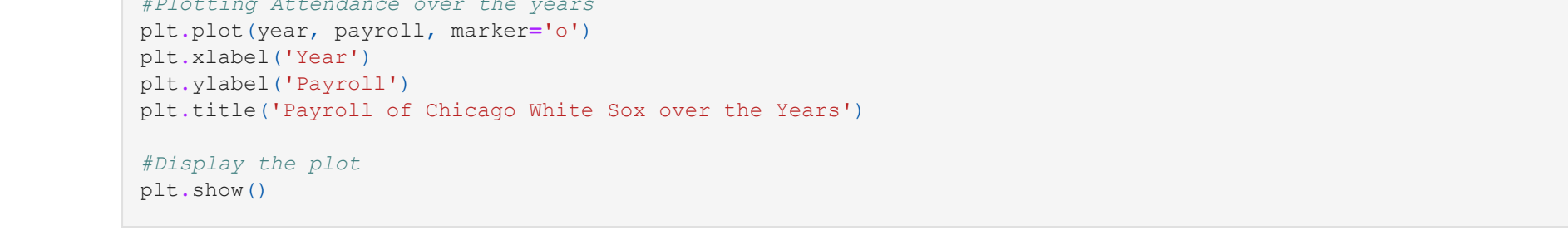


```
In [124]:
#Plotting Wins over the years
#Convert columns to numpy arrays
wins = np.array(white_sox_df['Year'])
wins = np.array(white_sox_df['W'])

plt.figure(figsize=(10,6))

#Plotting Attendance over the years
plt.plot(year, wins, marker='o')
plt.xlabel('Year')
plt.ylabel('Wins')
plt.title('Wins by Chicago White Sox over the Years')

#Display the plot
plt.show()
```



```
In [125]:
#Plotting Payroll over the years
#Convert columns to numpy arrays
year = np.array(white_sox_df['Year'])
payroll = np.array(white_sox_df['Est. Payroll'])

plt.figure(figsize=(10,6))

#Plotting Attendance over the years
plt.plot(year, payroll, marker='o')
plt.xlabel('Year')
plt.ylabel('Payroll')
plt.title('Payroll of Chicago White Sox over the Years')

#Display the plot
plt.show()
```

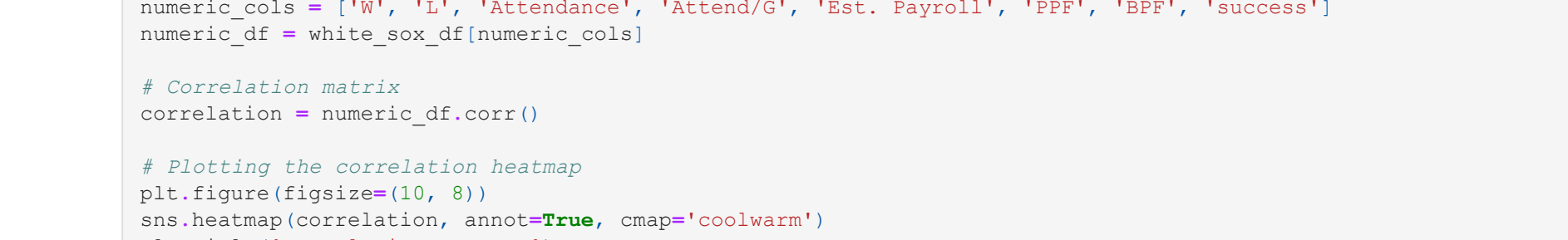


```
In [126]:
# Select numeric columns for correlation analysis
numeric_cols = ['W', 'L', 'Attendance', 'Attend/G', 'Est. Payroll', 'PPF', 'BPF', 'success']
numeric_df = white_sox_df[numerics_cols]

# Correlation matrix
correlation = numeric_df.corr()

# Plotting the correlation heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation, annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')

# Display the plot
plt.show()
```



```
In [127]:
#Extract the year from the "Date" column in the DataFrame
cws_df['Year'] = pd.to_datetime(cws_df['Date']).dt.year

#Merge df and white_sox_df DataFrames on the "Year" column
merged_cws_df = pd.merge(cws_df, white_sox_df[['Year', 'success']], on='Year', how='left')

#Add a new column 'Year_Successful' based on 'success' field
merged_cws_df['Year_Successful'] = merged_cws_df['success'].fillna(0).astype(int)

merged_cws_df.head()
```

	Title	Author	Date	Content	Year	success	Year_Successful
0	MLB Trade Rumors Podcast: Exciting Youth Movements in Cincinnati and Pittsburgh, Bad Central D...	Danrigh McDonald	2023-06-21	The latest episode of the MLB Trade Rumors Podcast is now live on Spotify, Apple Podcasts, and w...	2023	0	0
1	AL Central Notes: Buxton, Crochet, Tigers	Danrigh McDonald	2023-06-21	Twins outfielder Byron Buxton has dealt with many injuries throughout his career, which has led...	2023	0	0
2	White Sox Claim Touki Toussaint From Guardians	Danrigh McDonald	2023-06-20	Toussaint off waivers from the Guardians, per anno...	2023	0	0
3	White Sox Recall Jose Rodriguez For MLB Debut	Anthony Franco	2023-06-19	The White Sox announced a handful of transactions before tonight's series opener with the Ranger...	2023	0	0
4	White Sox Place Mike Clevinger On Injured List	Anthony Franco	2023-06-16	The White Sox placed starter Mike Clevinger on the 15-day injured list, retroactive to June 15, ...	2023	0	0

```
In [128]:
len(merged_cws_df)
```

1000

Cleaning the Data

```
In [129]:
#Identify any noise in the data
RE_SUSPICIOUS = re.compile(r'[!@<>{}~\|\\\/\|']

def impurity(text, min_len=10):
    """Returns the share of suspicious characters in a text"""
    if text == None or len(text) < min_len:
        return 0
    return len(RE_SUSPICIOUS.findall(text))/len(text)

merged_cws_df['Content'].apply
```

```
bound method Series.apply of 0 The latest episode of the MLB Trade Rumors Podcast is now live on Spotif
7 Apple Podcasts, and w...
2 The White Sox have claimed right-hander Touki Toussaint off waivers from the Guardians, per anno...
3 The White Sox announced a handful of transactions before tonight's series opener with the Ranger...
4 The White Sox placed starter Mike Clevinger on the 15-day injured list, retroactive to June 15, ...
995 The White Sox have outrighted catcher Dustin Garneau to Triple-A Charlotte, Darryl Van Schouwen o...
996 White Sox GM Rick Rahn addressed the media yesterday regarding the state of his organization's c...
997 White Sox outfielder Avisail Garcia is set to undergo right knee surgery, he told reporters incl...
998 By the end of the 2017 season, the list of pitchers closing out games for their respective teams...
999 It's often difficult to feel positive about a team when it is finishing out a season that won't h...
Name: Content, dtype: object
```

```
In [131]:
pd.options.display_max_colwidth = 100 ##
#Add new column to cws_df frame
merged_cws_df['Impurity'] = merged_cws_df['Content'].apply(impurity, min_len=10)

#Get the top 3 records
merged_cws_df[['Content', 'Impurity']].sort_values(by='Impurity', ascending=False).head(3)
```

	Content	Impurity
333	The White Sox aren't planning to make any coaching changes, manager Tony La Russa told Darryl Van... 0.003578	0.003578
875	Eloy Jimenez's career-opening season with the White Sox included "an understanding" that Jimenez... 0.002797	0.002797
450	Andre Engel will begin the season the injured list as White Sox manager Tony La Russa told Darryl Van... 0.002782	0.002782

The above shows the highest impurity levels for the White Sox's manager site. Note, these are extremely low...well below 1% of all characters in each blog post are suspicious.

Character Normalization and Tokenization

```
In [132]:
#Example function to normalize the text in the "Contents" column
def normalize_text(text):
    """Normalize the text
    text = text.lower()

    #Remove punctuation
    text = text.translate(str.maketrans("", "", string.punctuation))

    #Remove special characters and digits
    text = re.sub(r'[\W_]+', '', text)

    #Tokenize the text
    tokens = word_tokenize(text)

    #Remove stopwords
    stop_words = set(stopwords.words("english"))
    tokens = [token for token in tokens if token not in stop_words]

    #Perform stemming
    stemmer = PorterStemmer()
    tokens = [stemmer.stem(token) for token in tokens]

    return tokens
```

```
In [133]:
#Apply normalization to the "Contents" column
merged_cws_df['Contents_Normalized'] = merged_cws_df['Content'].apply(normalize_text)

#Print the head of the DataFrame with the normalized contents
merged_cws_df[['Content', 'Contents_Normalized']].head()
```

	Content	Contents_Normalized
0	The latest episode of the MLB Trade Rumors Podcast is now live on Spotify, Apple Podcasts, and w...	[latest, episod, mlb, trade, rumor, podcast, live, spotifi, appl, podcast, where, get, podcast...]
1	Twins outfielder Byron Buxton has dealt with many injuries throughout his career, which has led...	[twi, outfiel, byron, buxton, deal, mani, injuri, throughout, career, led, twi, use, exclus...]
2	The White Sox have claimed right-hander Touki Toussaint off waivers from the Guardians, per anno...	[white, sox, claim, righthand, touki, toussaint, waiver, guardiann, club, design, a...]
3	The White Sox announced a handful of transactions before tonight's series opener with the Ranger...	[white, sox, announc, hand, transact, toni, seri, open, ranger, notabl, recal, infiel, prosp...]
4	The White Sox placed starter Mike Clevinger on the 15-day injured list, retroactive to June 15, ...	[white, sox, place, starter, mike, cleving, day, injur, list, retroact, june, due, biop, inflam...]

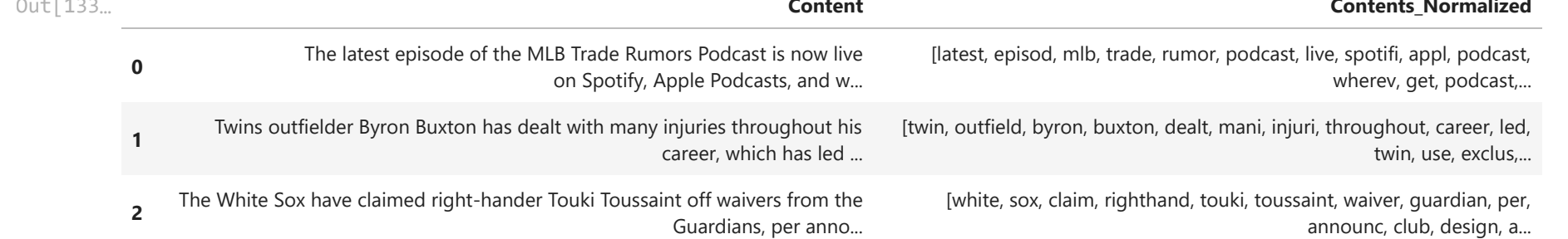
```
In [134]:
merged_cws_df['length'] = merged_cws_df['Contents_Normalized'].str.len()
merged_cws_df.head()
```

	Title	Author	Date	Content	Year	success	Year_Successful	Impurity	Contents_Normalized	length
0	MLB Trade Rumors Podcast: Exciting Youth Movements in Cincinnati and Pittsburgh, Bad Central D...	Danrigh McDonald	2023-06-21	The latest episode of the MLB Trade Rumors Podcast is now live on Spotify, Apple Podcasts, and w...	2023	0	0	0.0	[latest, episod, mlb, trade, rumor, podcast, live, spotifi, appl, podcast, where, get, podcast...]	124
1	AL Central Notes: Buxton, Crochet, Tigers	Danrigh McDonald	2023-06-21	Twins outfielder Byron Buxton has dealt with many injuries throughout his career, which has led...	2023	0	0	0.0	[twi, outfiel, byron, buxton, deal, mani, injuri, throughout, career, led, twi, use, exclus...]	377
2	White Sox Claim Touki Toussaint From Guardians	Danrigh McDonald	2023-06-20	Toussaint off waivers from the Guardians, per anno...	2023	0	0	0.0	[white, sox, claim, righthand, touki, toussaint, waiver, guardiann, club, design, a...]	143
3	White Sox Recall Jose Rodriguez For MLB Debut	Anthony Franco	2023-06-19	The White Sox announced a handful of transactions before tonight's series opener with the Ranger...	2023	0	0	0.0	[white, sox, announc, hand, transact, toni, seri, open, ranger, notabl, recal, infiel, prosp...]	165
4	White Sox Place Mike Clevinger On Injured List	Anthony Franco	2023-06-16	The White Sox placed starter Mike Clevinger on the 15-day injured list, retroactive to June 15, ...	2023	0	0	0.0	[white, sox, place, starter, mike, cleving, day, injur, list, retroact, june, due, biop, inflam...]	229

Analyze descriptive statistics for text in dataframe.

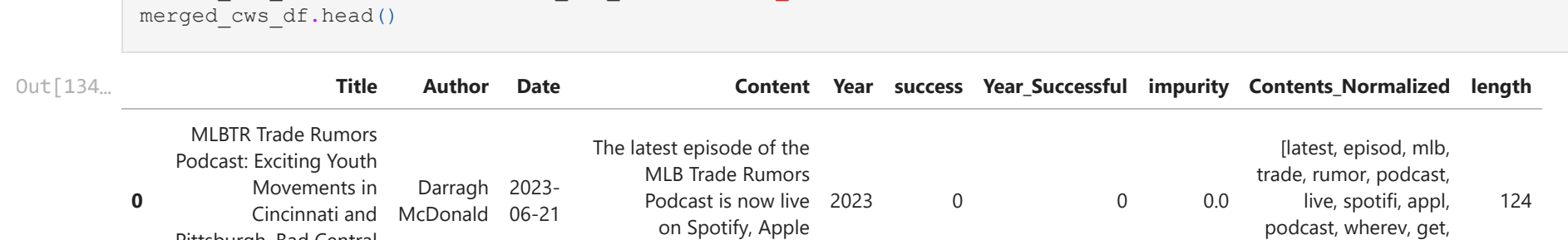
```
In [135]:
merged_cws_df['length'].plot(kind='box', vert=False, figsize=(6, 1))

#Axes: >
```



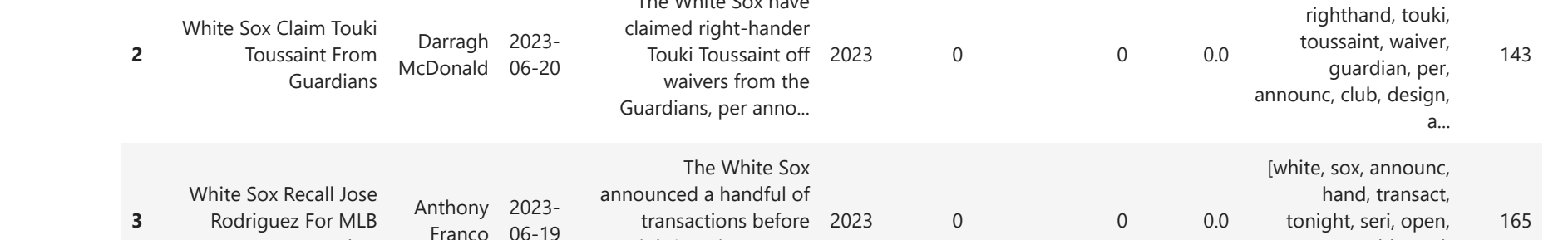
```
In [136]:
merged_cws_df['Title_length'] = merged_cws_df['Title'].str.len()
merged_cws_df['Title_length'].plot(kind='box', vert=False, figsize=(6, 1))

#Axes: >
```



```
In [137]:
merged_cws_df['length'].plot(kind='hist', bins=30, figsize=(8, 2))

#Axes: ylabel='Frequency'>
```



```
In [138]:
# Extract the month from the date and create a new column
merged_cws_df['Month'] = merged_cws_df['Date'].dt.month

In [139]:
# Plot the average post length
freq_df = pd.DataFrame.from_dict(counter, orient='index', columns=['freq'])
freq_df = freq_df.query('freq >= min_freq')
freq_df.index.name = 'token'

return freq_df.sort_values('freq', ascending=False)
```

```
In [140]:
def count_words(df, column='Contents_Normalized', preprocess=None, min_freq=2):
    # process tokens and update counter
    def update(doc):
        tokens = doc if preprocess is None else preprocess(doc)
        counter.update(tokens)

    # create counter and run through all data
    counter = Counter()
    tqdm.pandas() # initialize tqdm for progress bar
    df[column].apply(update)

    #
```


	Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium		
In [153].	0	2023	Los Angeles Angels	AL West	42	36	2	NaN	1181985.0	32833.0	5th of 15	\$224,228,095	103	103	Angel Stadium of Anaheim	
	1	2022	Los Angeles Angels	AL West	73	89	3	NaN	2457461.0	30339.0	5th of 15	\$177,063,095	103	103	Angel Stadium of Anaheim	
	2	2021	Los Angeles Angels	AL West	77	85	4	NaN	1515689.0	18484.0	6th of 15	\$188,408,395	103	103	Angel Stadium of Anaheim	
	3	2020	Los Angeles Angels	AL West	26	34	4	NaN	NaN	NaN	8th of 15	\$181,524,762	102	101	Angel Stadium of Anaheim	
	4	2019	Los Angeles Angels	AL West	72	90	4	NaN	3023012.0	37321.0	2nd of 15	\$158,078,584	100	100	Angel Stadium of Anaheim	
	5	2018	Los Angeles Angels	AL West	80	82	4	NaN	3020276.0	37287.0	2nd of 15	\$166,649,666	98	98	Angel Stadium of Anaheim	
	6	2017	Los Angeles Angels	AL West	80	82	2	NaN	3019585.0	37279.0	3rd of 15	\$181,125,500	97	97	Angel Stadium of Anaheim	
	7	2016	Los Angeles Angels	AL West	74	88	4	NaN	3016142.0	37236.0	3rd of 15	\$139,712,000	95	95	Angel Stadium of Anaheim	
	8	2015	Los Angeles Angels of Anaheim	AL West	85	77	3	NaN	3012765.0	37195.0	2nd of 15	\$131,525,500	94	94	Angel Stadium of Anaheim	
	9	2014	Los Angeles Angels of Anaheim	AL West	98	64	1	Lost ALDS (3-0)	3095953.0	38221.0	2nd of 15	\$128,667,000	94	95	Angel Stadium of Anaheim	
	10	2013	Los Angeles Angels of Anaheim	AL West	78	84	3	NaN	3019505.0	37278.0	4th of 15	\$116,532,500	94	94	Angel Stadium of Anaheim	
	11	2012	Los Angeles Angels of Anaheim	AL West	89	73	3	NaN	3061770.0	37800.0	3rd of 14	\$141,073,500	93	94	Angel Stadium of Anaheim	
	12	2011	Los Angeles Angels of Anaheim	AL West	86	76	2	NaN	3166332.0	39090.0	3rd of 14	\$138,543,166	92	92	Angel Stadium of Anaheim	
	13	2010	Los Angeles Angels of Anaheim	AL West	80	82	3	NaN	3250814.0	40134.0	2nd of 14	\$104,963,866	96	96	Angel Stadium of Anaheim	
	14	2009	Los Angeles Angels of Anaheim	AL West	97	65	1	Lost ALCS (4-3)	3240386.0	40005.0	2nd of 14	\$118,169,000	98	99	Angel Stadium of Anaheim	
	15	2008	Los Angeles Angels of Anaheim	AL West	100	62	1	Lost ALDS (3-1)	3367474.0	41194.0	2nd of 14	\$119,216,333	102	102	Angel Stadium of Anaheim	
	16	2007	Los Angeles Angels of Anaheim	AL West	94	68	1	Lost ALDS (3-0)	3365632.0	41551.0	2nd of 14	\$109,251,333	100	101	Angel Stadium of Anaheim	
	17	2006	Los Angeles Angels of Anaheim	AL West	89	73	2	NaN	3406790.0	42059.0	2nd of 14	\$103,472,000	99	100	Angel Stadium of Anaheim	
	18	2005	Los Angeles Angels of Anaheim	AL West	95	67	1	Lost ALCS (4-1)	3404686.0	42033.0	2nd of 14	\$94,867,822	97	98	Angel Stadium of Anaheim	
	<pre>#Remove the dollar and comma signs from the 'Estimated Payroll' column: angels_df['\$Est. Payroll'] = angels_df['\$Est. Payroll'].replace({'\\\$': ''}, regex = True) angels_df['\$Est. Payroll'] = angels_df['\$Est. Payroll'].replace(',', ''), regex = True angels_df.head()</pre>															
	Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium		
	0	2023	Los Angeles Angels	AL West	42	36	2	NaN	1181985.0	32833.0	5th of 15	\$224,228,095	103	103	Angel Stadium of Anaheim	
	1	2022	Los Angeles Angels	AL West	73	89	3	NaN	2457461.0	30339.0	5th of 15	\$177,063,095	103	103	Angel Stadium of Anaheim	
	2	2021	Los Angeles Angels	AL West	77	85	4	NaN	1515689.0	18484.0	6th of 15	\$188,408,395	103	103	Angel Stadium of Anaheim	
	3	2020	Los Angeles Angels	AL West	26	34	4	NaN	NaN	NaN	8th of 15	\$181,524,762	102	101	Angel Stadium of Anaheim	
	4	2019	Los Angeles Angels	AL West	72	90	4	NaN	3023012.0	37321.0	2nd of 15	\$158,078,584	100	100	Angel Stadium of Anaheim	
	<pre>angels_df['\$Est. Payroll'] = angels_df['\$Est. Payroll'].astype(float) angels_df.dtypes</pre>															
	Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium	success	
	0	2023	Los Angeles Angels	AL West	42	36	2	NaN	1181985.0	32833.0	5th of 15	\$224,228,095	103	103	Angel Stadium of Anaheim	1
	1	2022	Los Angeles Angels	AL West	73	89	3	NaN	2457461.0	30339.0	5th of 15	\$177,063,095	103	103	Angel Stadium of Anaheim	0

```
In [154]. #Remove the dollar and comma signs from the 'Estimated Payroll' column:
angels_df['Est. Payroll'] = angels_df['Est. Payroll'].replace('\\$|,','', regex = True)
angels_df['Est. Payroll'] = angels_df['Est. Payroll'].replace('\\$|,','', regex = True)
angels_df.head()
```

```
#Count the number of successes and failures
success_counts = angels_df['success'].value_counts()

#Create a bar plot
plt.bar(success_counts.index, success_counts.values)
plt.xticks(success_counts.index, ['Failure', 'Success'])
plt.xlabel('Success')
plt.ylabel('Count')
plt.title('Success Counts')

#Display the plot
plt.show()
```



Category	Count
Failure	9
Success	10

```
In [155]. angels_df['Est. Payroll'] = angels_df['Est. Payroll'].astype(float)
angels_df.dtypes
```

	Year	W	L	Finish	Attendance	Attend/G	Est. Payroll	PPF	BPF	success
count	19.000000	19.000000	19.000000	19.000000	1.800000e+01	18.000000	1.900000e+01	19.000000	19.000000	19.000000
mean	2019.000000	79.736642	72.473684	5.2526316	2.921996e+05	37074.308889	1.433150e+08	97.894737	98.157895	0.526316
std	5.627314	18.525784	15.832757	1.1722929	6.151490e+05	5527.931452	5.551089e+07	3.649882	3.513810	0.512989
min	2005.000000	26.000000	34.000000	1.181985e+06	18484.000000	9.486782e+07	92.000000	92.000000	0.000000	

```
In [156]. #Add a column to define a binary output that has a 1 if the team has a winning record for the year, and 0 if it
angels_df['success'] = angels_df['W'] / (angels_df['W'] + angels_df['L']) == 0.5).astype(int)
angels_df.head()
```

```
#Plotting the data
#Convert columns to numpy arrays
year = np.array(angels_df['Year'])
attendance = np.array(angels_df['Attendance'])

plt.figure(figsize=(10, 6))

#Plotting Attendance over the years
plt.plot(year, attendance, marker='o')
```

```
In [157]. #Count the number of successes and failures
success_counts = angels_df['success'].value_counts()

#Create a bar plot
plt.bar(success_counts.index, success_counts.values)
plt.xticks(success_counts.index, ('Failure', 'Success'))
plt.xlabel('Success')
plt.ylabel('Count')
plt.title('Success Counts')

#Display the plot
plt.show()
```



```
In [158]. angels_df.describe()
```

	Year	W	L	Finish	Attendance	Attend/G	Est. Payroll	PPF	BPF	success
Out [158].	count	19.000000	19.000000	19.000000	19.000000	1.800000e+01	1.900000e+01	19.000000	19.000000	19.000000
	mean	2014.000000	79.736842	72.473684	2.536376	2.921969e+06	37074.388889	1.433158e+08	97.894737	98.157895
	std	5.627314	18.525784	15.837257	1.172292	6.151490e+05	5527.931452	3.515893e+07	3.649882	3.531810
	min	2005.000000	26.000000	34.000000	1.000000	1.181985e+06	18484.000000	9.486782e+07	92.000000	92.000000
	25%	2009.500000	75.500000	66.000000	1.500000	3.016983e+06	37246.500000	1.173508e+08	94.500000	95.000000
	50%	2014.000000	80.000000	76.000000	3.000000	3.042391e+06	37560.500000	1.385432e+08	98.000000	98.000000
	75%	2018.500000	91.500000	83.000000	3.500000	3.248207e+06	40101.750000	1.719564e+08	101.000000	101.000000
	max	2023.000000	100.000000	90.000000	4.000000	3.406790e+06	42059.000000	2.242281e+08	103.000000	103.000000

```
In [159]. #Plotting the data
#Convert columns to numpy arrays
year = np.array(angels_df['Year'])
wins = np.array(angels_df['W'])
attendance = np.array(angels_df['Attendance'])

plt.figure(figsize=(10, 6))

#Plotting Attendance over the years
plt.plot(year, attendance, marker='o')
plt.xlabel('Year')
plt.ylabel('Attendance')
plt.title('Attendance of Los Angeles Angels over the Years')

#Display the plot
plt.show()
```

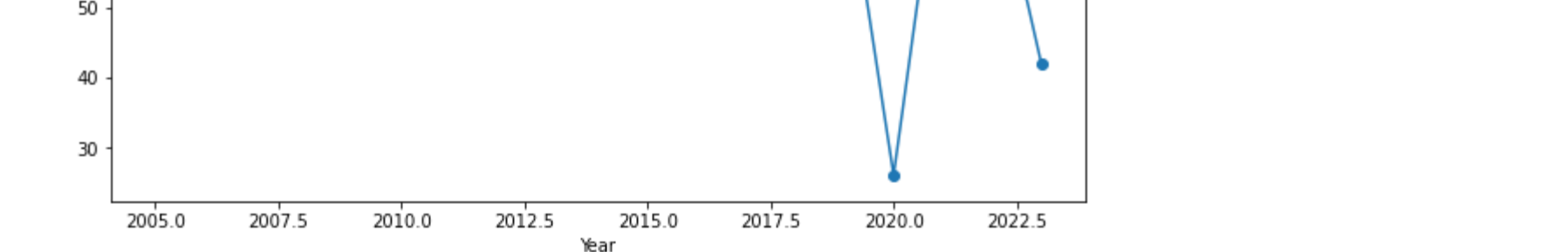


```
In [160]. #Plotting Wins over the years
#Convert columns to numpy arrays
year = np.array(angels_df['Year'])
wins = np.array(angels_df['W'])

plt.figure(figsize=(10, 6))

#Plotting Attendance over the years
plt.plot(year, wins, marker='o')
plt.xlabel('Year')
plt.ylabel('Wins')
plt.title('Wins by Los Angeles Angels over the Years')

#Display the plot
plt.show()
```



```
In [161]. #Plotting Payroll over the years
#Convert columns to numpy arrays
year = np.array(angels_df['Year'])
payroll = np.array(angels_df['Est. Payroll'])

plt.figure(figsize=(10, 6))

#Plotting Attendance over the years
plt.plot(year, payroll, marker='o')
plt.xlabel('Year')
plt.ylabel('Payroll')
plt.title('Payroll of Los Angeles over the Years')

#Display the plot
plt.show()
```

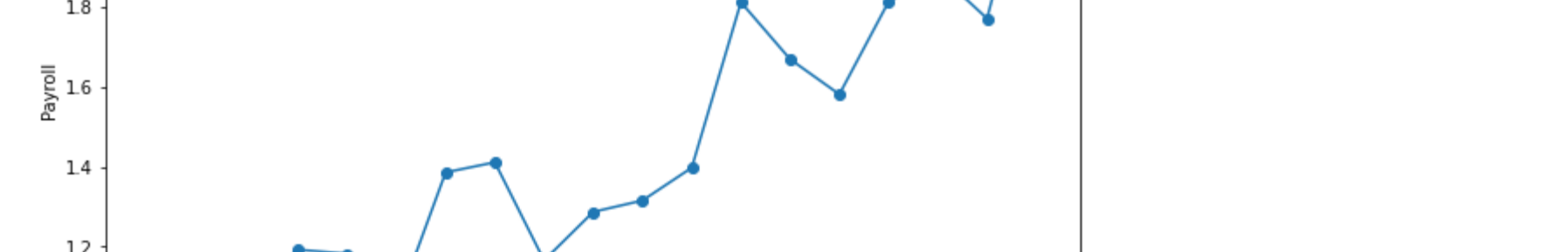


```
In [162]. #Select numeric columns for correlation analysis
numeric_cols = ['W', 'L', 'Attendance', 'Attend/G', 'Est. Payroll', 'PPF', 'BPF', 'success']
numeric_df = angels_df[numeric_cols]

#Correlation matrix
correlation = numeric_df.corr()

#Plotting the correlation heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation, annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')

#Display the plot
plt.show()
```



Add column to analyze whether blog post occurred in a successful year or not.

```
In [163]. #Extract the year from the 'Date' column in df DataFrame
laa_df['Year'] = pd.to_datetime(laa_df['Date']).dt.year

#Merge df and pd.DataFrame of dataFrames on the 'Year' column
merged_laa_df = pd.merge(laa_df, angels_df[['Year', 'success']], on='Year', how='left')

#Add a new column 'Year_Successful' based on 'success' column
merged_laa_df['Year_Successful'] = merged_laa_df['success'].fillna(0).astype(int)

merged_laa_df.head()
```

	Title	Author	Date	Content	Year	success	Year_Successful
Out [163].	0	Angels Acquire Mike Moustakas	Tim Deeks 2023-06-24	The Angels acquired infielder Mike Moustakas tonight, sending minor league righty Connor Van Sco...	2023	1	1
	1	Latest On Matt Moore	Nick Deeds 2023-06-24	Angels lefty Matt Moore is making progress in his rehab from an oblique injury that has left his...	2023	1	1
	2	Angels Designate Chris Okey For Assignment	Nick Deeks 2023-06-24	The Angels have designated catcher Chris Okey for assignment, according to Sam Blum of The Athle...	2023	1	1
	3	Angels To Promote David Fletcher	Mark Polshuk 2023-06-24	12:04PM: Walsh and infielder Michael Stefanic have been optioned to Triple-A to make room for FL...	2023	1	1
	4	Angels Acquire Eduardo Escobar	Anthony Franco 2023-06-23	The Mets and Angels pulled off an unexpected swap Friday night. New York dealt veteran infielder...	2023	1	1

```
In [164]. len(merged_laa_df)
1000
```

Cleaning the Angels' Data

```
In [165]. #Identify any noise in the data
RE_SUSPICIOUS = re.compile(r'[!@<>()\[\]\{\}']

def impurity(text, min_len=10):
    """returns the share of suspicious characters in a text"""
    if text == None or len(text) < min_len:
        return 0
    else:
        return len(RE_SUSPICIOUS.findall(text))/len(text)
```

```
In [166]. merged_laa_df['Content'].apply
bound Method Series.apply of 0
```

```
Out [166]. 0 The Angels acquired infielder Mike Moustakas tonight, sending minor league
1 Angels lefty Matt Moore is making progress in his rehab from an oblique injury that has left his...
2 The Angels have designated catcher Chris Okey for assignment, according to Sam Blum of The Athle...
3 12:04PM: Walsh and infielder Michael Stefanic have been optioned to Triple-A to make room for F...
4 The Mets and Angels pulled off an unexpected swap Friday night. New York dealt veteran infielder...
```

```
995 At 50-48 and 11 games back in the American League West, the Angels don't have much hope of conte...
996 Angels first baseman Albert Pujols exited the team's game Thursday with left hamstring tightness...
997 We'll use this post to cover the day's minor moves, both involving a pair of former Giants outfi...
998 5:03PM: Angels righty Hoo Ramirez has also been suspended for a trio of contests after also bein...
999 MONDAY: An MRI revealed a "small" calf strain, Hoonstra was among those to tweet. Trout's day-t...
Name: Content, Length: 1000, dtype: object
```

```
In [167]. pd.options.display.max_colwidth = 100 ##
#Add new columns to data frame
merged_laa_df['Impurity'] = merged_laa_df['Content'].apply(impurity, min_len=10)

#Get the top 3 records
merged_laa_df[['Content', 'Impurity']].sort_values(by='Impurity', ascending=False).head(3)
```

Cleaning the Angels' Data

```
#Identify any noise in the data
RE_SUSPICIOUS = re.compile(r'[^<()\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\|\\
```

The above shows the highest impurity levels for the Angels' data set. Note, these are extremely low...well below 1% of all characters in each blog post are suspicious.

Character Normalization and Tokenization

```
In [168]. #Example function to normalize the text in the "Contents" column
def normalize_text(text):
    #Lowercase the text
    text = text.lower()

    #Remove punctuation
    text = text.translate(str.maketrans("", "", string.punctuation))

    #Remove special characters and digits
    text = re.sub(r'[\W_]+', '', text)

    #Tokenize the text
    tokens = word_tokenize(text)

    #Remove stopwords
    stop_words = set(stopwords.words('English'))
    tokens = [token for token in tokens if token not in stop_words]

    #Perform stemming
    stemmer = PorterStemmer()
    tokens = [stemmer.stem(token) for token in tokens]

    return tokens

#Apply normalization to the "Contents" column
merged_laa_df['Contents_Normalized'] = merged_laa_df['Content'].apply(normalize_text)

#Print the head of the DataFrame with the normalized contents
merged_laa_df[['Content', 'Contents_Normalized']].head()
```

```
merged_laa_df[['impurity']].merged_laa_df[['Content']].apply(impurity, min_len=10)

#Get the top 3 records
merged_laa_df[['Content', 'impurity']].sort_values(by='impurity', ascending=False).head(3)
```

	Content	impurity
62	An early look at the trade deadline possibilities, particularly focusing on Shohei Ohtani and wh...	0.010753

```
In [169]. merged_laa_df['length'] = merged_laa_df['Contents_Normalized'].str.len()
merged_laa_df.head()
```

	Title	Author	Date	Content	Year	success	Impurity	Contents_Normalized	length
Out [169].	0	Angels Acquire Mike Moustakas	Tim Deeks 2023-06-24	The Angels acquired infielder Mike Moustakas tonight, sending minor league righty Connor Van Sco...	2023	1	1	0.0	angel, acquir, infield, mike, moustaka tonight, send, minor, league, right, connor, van, scoy...
	1	Latest On Matt Moore	Nick Deeks 2023-06-24	Angels lefty Matt Moore is making progress in his rehab from an oblique injury that has left his...	2023	1	1	0.0	angel, lefti, matt, moor, make, progress, rehab, obliqu, injuri, left, seldin, nearl, month, ...
	2	Designate Chris Okey for Assignment	Nick Deeks 2023-06-24	The Angels have designated catcher Chris Okey for assignment, according to Sam Blum of The Athle...	2023	1	1	0.0	angel, design, catchr, chri, okey, assign, accord, sam, blum, athlet, move, complet, seri...
	3	Angels To Promote David Fletcher	Mark Polshuk 2023-06-24	12:04PM: Walsh and infielder Michael Stefanic have been optioned to Triple-A to make room for FL...	2023	1	1	0.0	[pm, walsh, infield, michael, stefan, option, triple, make, room, fletcher, escobar, sam, blum...
	4	Angels Acquire Eduardo Escobar	Anthony Franco 2023-06-23	The Mets and Angels pulled off an unexpected swap Friday night. New York dealt veteran infielder...	2023	1	1	0.0	[met, angel, pull, unexpect, swap, friday, night, new, york, dealt, veteran, infield, eduardo, e...

Analyze descriptive statistics for text in dataframe.

```
In [170]. merged_laa_df['length'].plot(kind='box', vert=False, figsize=(8, 1))

<Axes: >
```



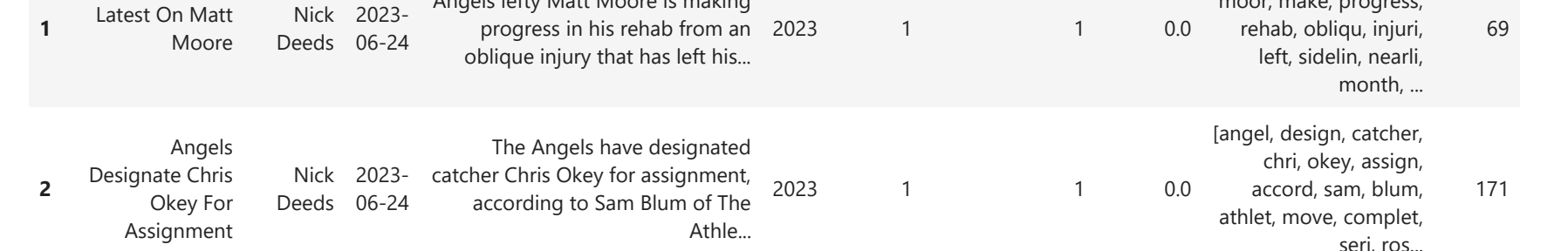
```
In [171]. merged_laa_df['Title_length'] = merged_laa_df['Title'].str.len()
merged_laa_df['Title_length'].plot(kind='box', vert=False, figsize=(8, 1))

<Axes: >
```



```
In [172]. merged_laa_df['length'].plot(kind='hist', bins=30, figsize=(8, 2))

<Axes: xlabel='Frequency'>
```



```
In [173]. #Extract the month from the date and create a new column
merged_laa_df['Month'] = merged_laa_df['Date'].dt.month

#Plot the average post length
merged_laa_df.groupby('Month').agg(['length': 'mean']) \
    .plot(title='Avg. Post Length', ylim=(0, 500), figsize=(6, 2))

<Axes: title='Center': 'Avg. Post Length', xlabel='Month'>
```



```
In [174]. #df count_words(df, column='Contents_Normalized', preprocess=None, min_freq=2):
def update(doc):
    tokens = doc if preprocess is None
```



```
Topic 0:
paul (0.4149124755426609)
coals (0.608102953805686)
cuts (0.5956734841045294)
gosselin (0.4432252252782354)
lieutenant (0.5215348253110752)

Topic 1:
kt (0.6314216353220825)
224 (0.6055153923434654)
held (0.4432252252782354)
garza (0.44960233614118694)
interests (0.4019971433366223)

Topic 2:
pattern (1.4610200879160924)
avoidance (1.2440177979133975)
simon (0.5720011132254754)
allen (0.3801362139696933)
donny (0.3541602294434306)

Topic 3:
000 (1.195462964992141)
delaplane (1.094642059657612)
eighty (0.9215887565063594)
224 (0.5592900793112495)
kay (0.508301960741038)

Topic 4:
culture (1.447942461361509)
oscar (0.788676300463056)
persorption (0.6667472727954942)
committal (0.4791995645994407)
393 (0.39497219808913914)

Topic 5:
stration (1.8615727805768956)
treat (0.5494023156369307)
fanbases (0.5188406530318043)
uncomfortable (0.5141113800589384)
eclipse (0.4834722606616162)

Topic 6:
characterize (1.0552092666550668)
425m (0.8180447661698501)
schaaf (0.6966382809408916)
pencil (0.55870086464737393)
encarnation (0.5466705816124041)

Topic 7:
velo (1.761040041407068)
etunning (0.4802126262826204)
traction (0.4122252252782354)
2010 (0.40481119021376244)
leads (0.3890107409662605)

Topic 8:
extensions (1.381403207536027)
grilled (0.6218271294747395)
dropoff (0.561263221030135)
stove (0.4441136632327407)
infamous (0.366617262383625)

Topic 9:
refusal (0.661202214205134)
concern (0.5996383782481639)
204mm (0.41254874169126693)
mm (0.35855024801390005)
clauses (0.3519597135631397)
```

Modeling

Logistic Regression Modeling

```
In [228]: !pip install scikit-learn

Requirement already satisfied: scikit-learn in c:\users\andre\downloads\new folder\envs\rlib\site-packages (1.2.2)
Requirement already satisfied: joblib>=1.1.1 in c:\users\andre\downloads\new folder\envs\rlib\site-packages (ffrom scikit-learn) (1.2.0)
Requirement already satisfied: numpy>=1.3.2 in c:\users\andre\downloads\new folder\envs\rlib\site-packages (ffrom scikit-learn) (1.7.3)
Requirement already satisfied: numpy>=1.17.3 in c:\users\andre\downloads\new folder\envs\rlib\site-packages (ffrom scikit-learn) (1.21.5)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\andre\downloads\new folder\envs\rlib\site-packages (from scikit-learn) (3.1.0)

In [229]: from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

In [230]: #Split the data into 80/20 train-test sets
X_train, X_test, y_train, y_test = train_test_split(comb_df['Content'], comb_df['success'], test_size=0.2, random_state=0)

In [231]: vectorizer = TfidfVectorizer()
X_train_vec = vectorizer.fit_transform(X_train)
X_test_vec = vectorizer.transform(X_test)

In [232]: model = LogisticRegression()
model.fit(X_train_vec, y_train)

Out[232]: * LogisticRegression
LogisticRegression()

In [233]: y_pred = model.predict(X_test_vec)

In [234]: accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)

Accuracy: 0.7758333333333334
```

Penalized Logistic Regression

This model was penalized because of our data containing numerous variables, which resulted in shrinking the coefficients of the less contributive variables toward zero.

```
In [235]: import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report
from sklearn.model_selection import GridSearchCV

In [236]: #Grid search was implemented to identify the ideal values for a model's hyperparameters.
#Set a grid of parameters to map estimator parameters to sequences of allowed values
param_grid = {'C': [0.001, 0.01, 0.1, 1, 10]}
grid_search = GridSearchCV(LogisticRegression(penalty='l2'), param_grid, cv=5)
grid_search.fit(X_train_vec, y_train)

best_C = grid_search.best_params_['C']

In [237]: model = LogisticRegression(penalty='l2', C=best_C)
model.fit(X_train_vec, y_train)

Out[237]: * LogisticRegression
LogisticRegression(C=1)

In [238]: y_pred = model.predict(X_test_vec)

In [239]: accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
print(classification_report(y_test, y_pred))

Accuracy: 0.7758333333333334

              precision    recall  f1-score   support

    0         0.79         0.78         0.78         627
    1         0.76         0.77         0.77         973

 accuracy          0.78         0.78         0.78        1200
 macro avg         0.78         0.78         0.78        1200
weighted avg         0.78         0.78         0.78        1200
```

The results from our classification model indicates that the model is able to accurately predict if a blog post is associated with a winning MLB team roughly 77%-78% of the time.

```
In [240]: coefficients = model.coef_
print(coefficients)

[[ 0.00694688 -0.07996149 -0.03090452 ... -0.05781421 -0.00909974
 -0.10394782]]

In [241]: # Retrieve feature names from vectorizer
feature_names = vectorizer.get_feature_names_out()

# Calculate feature importance
feature_importance = np.abs(model.coef_[0])

# Sort feature importance in descending order
sorted_indices = np.argsort(feature_importance)[::-1]
sorted_feature_importance = feature_importance[sorted_indices]
sorted_feature_names = feature_names[sorted_indices]

In [242]: # Select top 20 features
top_feature_importance = sorted_feature_importance[:50]
top_feature_names = sorted_feature_names[:50]

# Plot top 20 feature importance
plt.figure(figsize=(15, 11))
plt.barh(range(len(top_feature_names)), top_feature_importance, align='center')
plt.yticks(range(len(top_feature_names)), top_feature_names)
plt.xlabel('Feature Importance')
plt.ylabel('Feature')
plt.title('Logistic Regression - Top 20 Feature Importance')
plt.show()
```



The bar chart of our logistic regression indicates that the features team and player names mention show greatest importance towards the model's prediction process. It also demonstrates a strong performance when predicting if a blog post is associated with a winning team.

Overall, our data text analysis demonstrates a concise analysis between textual data from blog posts and MLB team performance. This analysis is one that could be used to evaluate team decisions and if performance of teams is trending in the right direction. This also demonstrates that our text data has the potential to predict status of on team performance and determine if blog posts are associated with a winning or losing team.

Amid some limitations noted within our data analysis that included data as far back as 2019 being used, the 2020 COVID-19 year impacting outcomes, and further cleaning the data to remove team names while removing biases, the use of additional features like attendance, revenue, and payroll could enhance the predictability rating of our modeling process.