

Preamble

Title: Powerball Winners Data (2021)

Author: Ashish Acharya

Email: ashishacharya0@gmail.com

Last Update: 2021-09-29

Introduction

This report analyses the Powerball winners numbers in 2020. Data for this report was obtained from information published by lotto.net.

▼ Dataset

The dataset for this report was built from informaton published daily at [Powerball Numbers from 2020](#). Web scraping techniques were employed to compile the winning numbers.

About Powerball: Powerball is an America lottery game in which a player picks 5 numbers from [1-69] and 1 number from [1-26] for a given drawing. They have options for powerplay which adds to the jackpot. The more numbers a player matches to the jackpot number, the bigger the payout.

I have chosen 2020 as the data to work with because it shows a complete picture of how the drawing pool looks like at the end of a year.

The script below automatically extracts Winning numbers. The information is saved in a json file (powerball_numbers_data.json) with the following structure:

```
[  
  {"date": "December 30th 2020", "numbers": [3,43,45,61,65]},  
  ...  
  {"date": "January 1st 2020", "numbers": [49,53,57,59,62]}  
]
```

```
1 # CS390Z - Introduction to Data Mining - Fall 2021  
2 # Instructor: Thyago Mota  
3 # Description: data collection  
4 # Author: Ashish Acharya  
5  
6 from bs4 import BeautifulSoup  
7 from google.colab import drive
```

```

8  from datetime import datetime, timedelta
9  import requests
10 import json
11 import time
12
13 # definitions/parameters
14 DATA_FOLDER = '/content/drive/MyDrive/Colab_Datasets/powerball_numbers_data/'
15 DATASET_NAME = 'powerball_numbers_data.json'
16 BASE_URL = 'https://www.lotto.net/powerball/numbers/2020'
17 HEADERS = {"User-Agent": "Mozilla/5.0 (X11; CrOS x86_64 12871.102.0) AppleWebKit/537
18 records = [] #dict for json dump
19 winning_numbers_list = []
20 # Google drive mount
21 # drive.mount('/content/drive')
22
23 result = requests.get(BASE_URL)
24 # 200 means its successful
25 if result.status_code == 200:
26     soup = BeautifulSoup(result.content, 'html.parser')
27     #to check if you actually got anything
28     #print(soup)
29
30     #the information for 1 drawing - 5 winning numbers & the powerball number
31     # are stored in the results-vsmall archive list.
32     # we find all of them and extract data from them
33     divs = soup.find_all('div', class_="results-vsmall archive-list")
34     for div in divs:
35         #print to check what you are working with
36         #print(div)
37
38         #There is an unwanted information "Day", the day on which the numbers were drawn
39         #we use extract to remove it and just get Month, Date and Year
40         date = div.find('div', class_='date')
41         unwanted = date.find('span')
42         unwanted.extract()
43         #there is a newline character before the date so we strip it
44         date = date.text.strip()
45         #print("Date:", date)
46
47         #this holds the 5 winning numbers
48         winning_numbers = div.find_all('li', class_='ball ball')
49
50         #printing the 5 winning numbers
51         for numbers in winning_numbers:
52             winning_number = numbers.text
53             winning_numbers_list.append(int(winning_number))
54             #print("Winning Number:", winning_number)
55         #print(winning_numbers_list)
56
57         #printing the powerball number
58         powerball_number = div.find('li', class_='ball power-play')
59         span = powerball_number.find('span') #span holds the power ball number

```

```

59     span = powerball_number.find( span ) #span holds the powerball number
60     span = int(span.text)
61     #print("Powerball Number:", span.text)
62
63     #write numbers to dictionary
64     records.append({'date': date, 'numbers': winning_numbers_list[0:6]})
65     #powerball numbers have been omitted since their range is 1-26
66     #which would need a separate handling of their own
67     #just need duplicate code to append numbers[6] separately in later parts
68     #of the code for statistics and graph
69     #records[date].append(span)
70
71     #reset list
72     winning_numbers_list.clear()
73
74
75     print(records)
76
77     with open(DATA_FOLDER + DATASET_NAME, 'w') as json_file:
78         json.dump(records, json_file, ensure_ascii=False, indent=2)
79

```

```

[{'date': 'December 30th 2020', 'numbers': [3, 43, 45, 61, 65]}, {'date': 'December

```

▼ Summary Statistics

```

1  # CS390Z - Introduction to Data Mining - Fall 2021
2  # Instructor: Thyago Mota
3  # Description: summary statistics
4
5  import json
6  from google.colab import drive
7  from datetime import datetime, timedelta
8  import numpy as np
9
10 # definitions/parameters
11 DATA_FOLDER = '/content/drive/MyDrive/Colab_Datasets/powerball_numbers_data/'
12 DATASET_NAME = 'powerball_numbers_data.json'
13
14 # Google drive mount
15 # drive.mount('/content/drive')
16
17 with open(DATA_FOLDER + DATASET_NAME, 'rt') as json_file:
18     records = json.load(json_file)
19
20 numbers = []
21 for record in records:
22     for num in record['numbers']:

```

```

23     numbers.append(num)
24     #print(numbers)
25     numbers_array = np.array(numbers)
26
27     print('*** Summary Statistics ***')
28     print(f'#records: {len(records)}')
29     # [1-69] is expected and it is what we get
30     print(f'Numbers range: [{np.min(numbers_array)},{np.max(numbers_array)}]')
31     print('Numbers mean: {:.2f}'.format(np.mean(numbers_array)))
32     print('Numbers median: {:.2f}'.format(np.median(numbers_array)))
33     print('Numbers std: {:.2f}'.format(np.std(numbers_array)))

*** Summary Statistics ***
#records: 105
Numbers range: [1,69]
Numbers mean: 35.25
Numbers median: 36.00
Numbers std: 19.92

```

▼ Visualizations

```

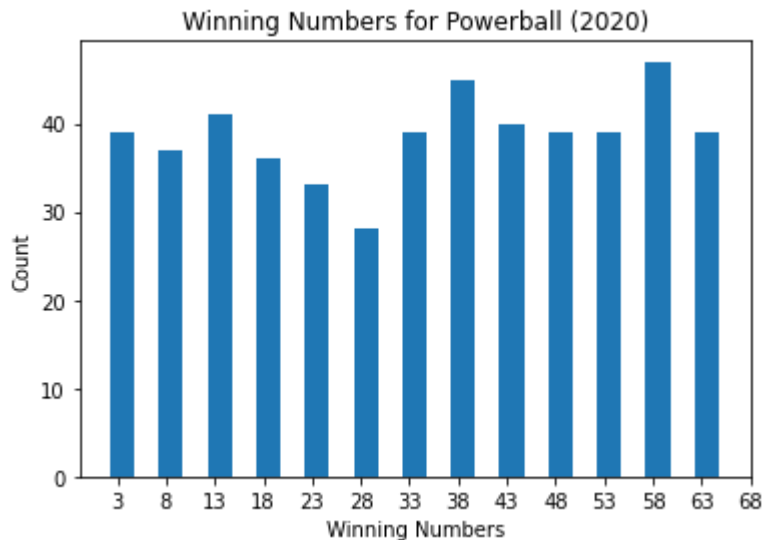
1  # CS390Z - Introduction to Data Mining - Fall 2021
2  # Instructor: Thyago Mota
3  # Description: histogram
4
5  from google.colab import drive
6  import matplotlib.pyplot as plt
7
8  # definitions/parameters
9  DATA_FOLDER = '/content/drive/MyDrive/Colab_Datasets/powerball_numbers_data/'
10 DATASET_NAME = 'powerball_numbers_data.json'
11
12 # Google drive mount
13 # drive.mount('/content/drive')
14
15 with open(DATA_FOLDER + DATASET_NAME, 'rt') as json_file:
16     records = json.load(json_file)
17
18 numbers = []
19 for record in records:
20     for num in record['numbers']:
21         numbers.append(num)
22 #print(numbers)
23
24 bins = list(range(1, 69, 5))
25 counts, bins, _ = plt.hist(
26     ...numbers,
27     ...bins=bins,
28     ...rwidth=0.5

```

```

29 )
30 xticks.=.[x+2.for.x.in.bins]
31 axes.=.plt.gca().#.get.a.reference.to.the.plot's.axes
32 axes.set_xticks(xticks)
33 plt.xlabel('Winning.Numbers')
34 plt.ylabel('Count')
35 plt.title('Winning.Numbers.for.Powerball.(2020)')
36 plt.show()

```



```

1 # CS390Z - Introduction to Data Mining - Fall 2021
2 # Instructor: Thyago Mota
3 # Description: box plot
4
5 from google.colab import drive
6 import matplotlib.pyplot as plt
7
8 # definitions/parameters
9 DATA_FOLDER = '/content/drive/MyDrive/Colab_Datasets/powerball_numbers_data/'
10 DATASET_NAME = 'powerball_numbers_data.json'
11
12 # Google drive mount
13 # drive.mount('/content/drive')
14
15 with open(DATA_FOLDER + DATASET_NAME, 'rt') as json_file:
16     records = json.load(json_file)
17
18 numbers = []
19 for record in records:
20     for num in record['numbers']:
21         numbers.append(num)
22
23 bp = plt.boxplot(
24     numbers,
25     vert=False
26 )
27 for median in bp['medians']:

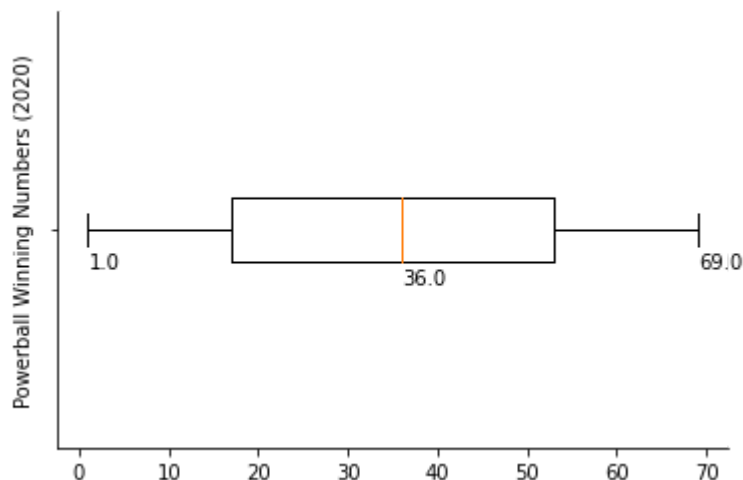
```

```

27     for median in bp['medians']:
28         xy = median.get_xydata()[0]
29         xy[1] -= .05
30         plt.annotate(str(xy[0]), xy=xy)
31
32     for cap in bp['caps']:
33         xy = cap.get_xydata()[0]
34         xy[1] -= .05
35         plt.annotate(str(xy[0]), xy=xy)
36
37     min_whisker = bp['caps'][0].get_xydata()[0][0]
38     max_whisker = bp['caps'][1].get_xydata()[0][0]
39
40     outliers = []
41     for record in records:
42         for num in record['numbers']:
43             if num < min_whisker or num > max_whisker:
44                 outliers.append(num)
45     print('*** Outliers ***')
46     for outlier in outliers:
47         print(outlier)
48
49     axes = plt.gca()
50     axes.spines['right'].set_visible(False)
51     axes.spines['top'].set_visible(False)
52     axes.set_yticklabels([''])
53     plt.ylabel('Powerball Winning Numbers (2020)')
54
55     plt.show()

```

*** Outliers ***



```

1  # CS390Z - Introduction to Data Mining - Fall 2021
2  # Instructor: Thyago Mota
3  # Description: time series
4
5  #Another interesting way to look at the data is counting how many even/odd
6  #numbers were there in a given winning drawing.
7  #Here we have a plot of winning number sets in which there were

```

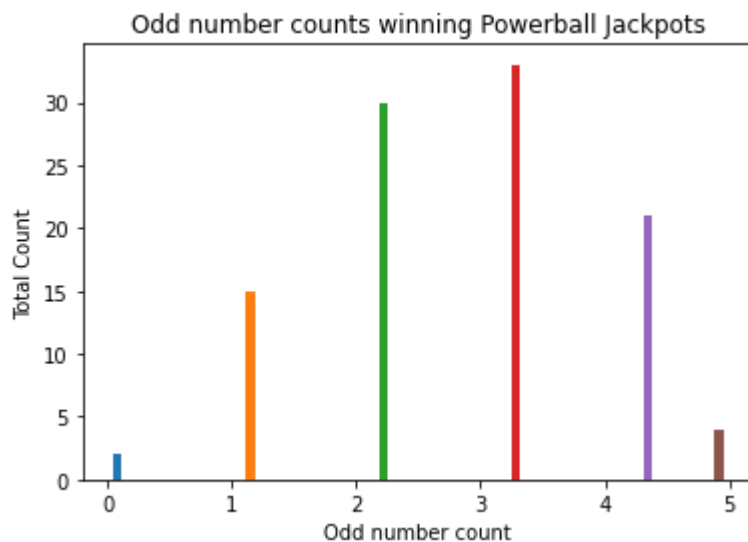
```
7 # where we have a plot of winning number 3003 in which there were
8 # 0 odd numbers, 1 odd numbers, .. all 5 odd numbers.
9
10 from google.colab import drive
11 import matplotlib.pyplot as plt
12 from datetime import datetime, timedelta
13
14 # definitions/parameters
15 DATA_FOLDER = '/content/drive/MyDrive/Colab_Datasets/powerball_numbers_data/'
16 DATASET_NAME = 'powerball_numbers_data.json'
17
18 # Google drive mount
19 # drive.mount('/content/drive')
20
21 with open(DATA_FOLDER + DATASET_NAME, 'rt') as json_file:
22     records = json.load(json_file)
23
24 odd0 = 0
25 odd1 = 0
26 odd2 = 0
27 odd3 = 0
28 odd4 = 0
29 odd5 = 0
30
31 odd = 0
32 for record in records:
33     for num in record['numbers']:
34         if num % 2 == 1:
35             odd += 1
36         if odd == 0:
37             odd0 += 1
38         if odd == 1:
39             odd1 += 1
40         if odd == 2:
41             odd2 += 1
42         if odd == 3:
43             odd3 += 1
44         if odd == 4:
45             odd4 += 1
46         if odd == 5:
47             odd5+=1
48
49     odd = 0
50
51 oddlist = []
52 oddlist.append([0]*odd0)
53 oddlist.append([1]*odd1)
54 oddlist.append([2]*odd2)
55 oddlist.append([3]*odd3)
56 oddlist.append([4]*odd4)
57 oddlist.append([5]*odd5)
58 #print(oddlist)
```

```

59
60 #verify that the sum of these match the total records
61 # print(odd0)
62 # print(odd1)
63 # print(odd2)
64 # print(odd3)
65 # print(odd4)
66 # print(odd5)
67
68 bins = list(range(0, 5))
69 counts, bins, _ = plt.hist(
70     oddlist
71 )
72
73 plt.xlabel('Odd number count')
74 plt.ylabel('Total Count')
75 plt.title('Odd number counts winning Powerball Jackpots')
76 plt.show()

```

/usr/local/lib/python3.7/dist-packages/numpy/core/_asarray.py:83: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a common outcome of np.array() of a list of nested arrays) is deprecated. If you would like to see this warning, you should explicitly use object(). Otherwise, create the array as follows: np.array([], dtype=object, copy=False, order=order)



✓ 0s completed at 5:14 AM

● ✕