

```
[36]: import imageio as ii
import visvis as vv
import numpy as np
import matplotlib
import matplotlib.pyplot as pyplot
import pandas as pd
from FIL import Image
import csv
import os
%matplotlib inline
pd.options.mode.chained_assignment = None

# Set working directories
directory = r'C:\Users\jeric\OneDrive\Documents\classFiles\DSC680\Project_1-new\puppeteer_scraper\images\'
directory2 = r'C:\Users\jeric\OneDrive\Documents\classFiles\DSC680\Project_1-new\puppeteer_scraper\thumbs\'
directory3 = r'C:\Users\jeric\OneDrive\Documents\classFiles\DSC680\Project_1-new\puppeteer_scraper\thumbs_min\'
```

```
print("Rows in df: "+str(len(df)))
```

df.head()						
Rows in df: 4686						
Out[97]:	hash	site_nm	date	country	img_url	id
0	8dbdf2f936a5fca98539799183a36a2b	Marble	June 11, 2022	Netherlands	https://assets.awwwards.com/awards/media/cache...	8dbdf2f936a5fca98539799183a36a2b
1	75556269230285f6a74414104	POLA: M...	June 12, 2022	Netherlands	https://assets.awwwards.com/awards/media/cache...	75556269230285f6a74414104

			Day	2022			
2	ea4e13c773862ebf76812971cdd07b25	Zulu Longines	June 9, 2022	France	https://assets.awwwards.com/awards/media/cache...	ea4e13c773862ebf76812971c	
3	33e8c7a0c2c17fbc7baffb80f808bbe	Ana Blagojevic	June 8, 2022	Italy	https://assets.awwwards.com/awards/media/cache...	33e8c7a0c2c17fbc7baffb80f	
4	caec15565ed1b1bf6efec6585191a0330	Kim Kneipp Folio	June 7, 2022	United States	https://assets.awwwards.com/awards/media/cache...	caec15565ed1b1bf6efec658519	

### Bring in the Nominee's csv data

```
In [98]: df2 = pd.read_csv(r'C:\Users\jeric\OneDrive\Documents\classFiles\DSC680\Project_1-new\puppeteer_scrape\data\nominees.csv')
df2.columns = ['hash', 'site_nm', 'date', 'country', 'img_url', 'img_full_nm']
df2['winner'] = 0

print("Rows in df: "+str(len(df2)))

df2.head()

Rows in df: 13083
```

			Z2				
1	3f2dfa1d6327bbdd1058ff43834941597	Cobo© Jun-22	Netherlands	<a href="https://assets.awwwards.com/awards/media/cache...">https://assets.awwwards.com/awards/media/cache...</a>	3f2dfa1d6327bbdd1058ff43834941597		
2	2777b6d0afb87db9ffa9f84d89704145	Paul McCartney © Jun-22	United Kingdom	<a href="https://assets.awwwards.com/awards/media/cache...">https://assets.awwwards.com/awards/media/cache...</a>	2777b6d0afb87db9ffa9f84d89704145		
3	36301ba2bh96c74r9r9b9deec17ch79a	Callista © Jun-22	Netherlands	<a href="http://assets.awwwards.com/awards/media/cache...">http://assets.awwwards.com/awards/media/cache...</a>	36301ba2bh96c74r9r9b9deec17ch79a		

```
4 fe15f5752fdbdd501cddcc17012a05f7 .PEAM Jun-10-22 Germany https://assets.awwwards.com/awards/media/cache... fe15f5752fdbdd501cddcc17012a05f7

Comparing CSV Files

In [39]: df['date'] = pd.to_datetime(df['date']) #winners
df2['date'] = pd.to_datetime(df2['date']) #nominees

print("Winner's record count: "+str(len(df)))
print("Winner's min date: "+str(df['date'].min()))
print("Winner's max date: "+str(df['date'].max()))
print("")
print("Nominee's record count: "+str(len(df2)))
print("Nominee's min date: "+str(df2['date'].min()))
print("Nominee's min date: "+str(df2['date'].max()))

Nom_Min = df2['date'].min()
Nom_Max = df2['date'].max()

print(Nom_Min)
print(Nom_Max)
```

```
Winner's min date: 2
Winner's max date: 2
```

```

Nominee's min date: 2019-12-05 00:00:00
Nominee's min date: 2022-06-10 00:00:00
2019-12-05 00:00:00
2022-06-10 00:00:00

```

### Trimming Winner's dataframe to match date range of nominee's dataframe

```

In [100]: df2_trunc = df[(df.date <= Nom_Max) & (df.date>= Nom_Min)]
          len(df2_trunc)

```

```
In [101]: print("Winner(1) count in nominee's dataframe")
df3 = pd.DataFrame()
df3 = df2.assign(winner=df2_trunc['site_nm'].isin(df2['site_nm']).astype(int))
df3.groupby(['winner'], dropna=False, as_index=False).size()
```

```

Winner(1) count in nominee's dataframe
Out[101]:
  winner  size
0      0.0    14
1      1.0   892
2      NaN  12177

```

### Setting NaNs to zero

```

In [102]: print("Winner(1) count in dataframe after checking against winner's file")
df3['winner'] = df3['winner'].fillna(0)
df3.groupby(['winner'], dropna=False, as_index=False).size()

```

```
Out[102]:
```

	0	1
0	0.0	12191
1	1.0	892

```
In [103]: df3 = df3[df3.winner != 1]
df3.groupby(['winner'], dropna=False, as_index=False).size()

Out[103]:
```

winner	size
0	0.0 12191

### Combining datasets into 1 dataframe

```
In [104]: print("Winner Dataframe length: "+str(len(df2_trunc)))
print("Non-winning Nominees Dataframe length: "+str(len(df3)))
print("Should add to: "+ str(len(df2_trunc)+len(df3)))
```

```
len(cmb_df)
```

```
Non-winning Nominees Dataframe length:12191
Should add to: 13097

Out[104_]: 13097

In [111]:
cmb_df['winner'].value_counts()

Out[111]:
0.0    12191
1.0      906
Name: winner, dtype: int64
```

```
In [112]: n=0
tot = 0
for i, row in cmb_df.iterrows():
    tot = tot + 1
    filename = row[5]
    try:
        img = Image.open(directory+filename)
```

```

        cmb_df.at[i, '404_error'] = 0

    except:
        cmb_df.at[i, '404_error'] = 1
        n = n+1

print("Total Images: "+ str(tot))
print("Error Images: "+str(n))

Total Images: 13097

```

```
In [113]: cmb_df_2 = cmb_df[(cmb_df['404_error']==0)]
len(cmb_df_2[(cmb_df_2['404_error']==1)])

Out[113]: 0

In [114]: len(cmb_df_2)
```

```
Out[114]: 9815
```

```
In [120]: cmb_df_2['winner'].value_counts()
```

```
Out[120]: 0.0    8918  
         1.0    897  
         Name: winner, dtype: int64
```

**Creating 50 X 50 copies of images**

```
#####
### ONLY RUN THIS ONCE
#####
```

```

102         filenames = os.listdir(directory)
103         filename = row[5]
104         try:
105             img = Image.open(directory+filename)
106             if img is not None:
107                 image = image.open(directory+filename)
108                 image = image.convert('RGB')
109                 image.thumbnail((50,50))
110                 image.save(directory+row[0]+'_'+row[1]+'.jpg')
111         except:
112             print("=skipped")

```

### Add unique color count from 50X50 images to dataframe

```
#Checking that all rows in dataframe were updated
missing_color_count = len(cmb_df_2[cmb_df_2['unique_color_count_50'].isnull() == True])
print("Missing Color Count: " + str(missing_color_count))

Missing Color Count: 0

In [125]: print("Max Color Count:" +str(max(cmb_df_2['unique_color_count_50'])))
          print("Min Color Count:" +str(min(cmb_df_2['unique_color_count_50'])))

Max Color Count:1817.0
```

```
In [126]: cmb_df_2['winner'].value_counts()

Out[126]: 0.0    8918
          1.0     897
          Name: winner, dtype: int64
```

### Question 1: How many Unique colors are there in the images?

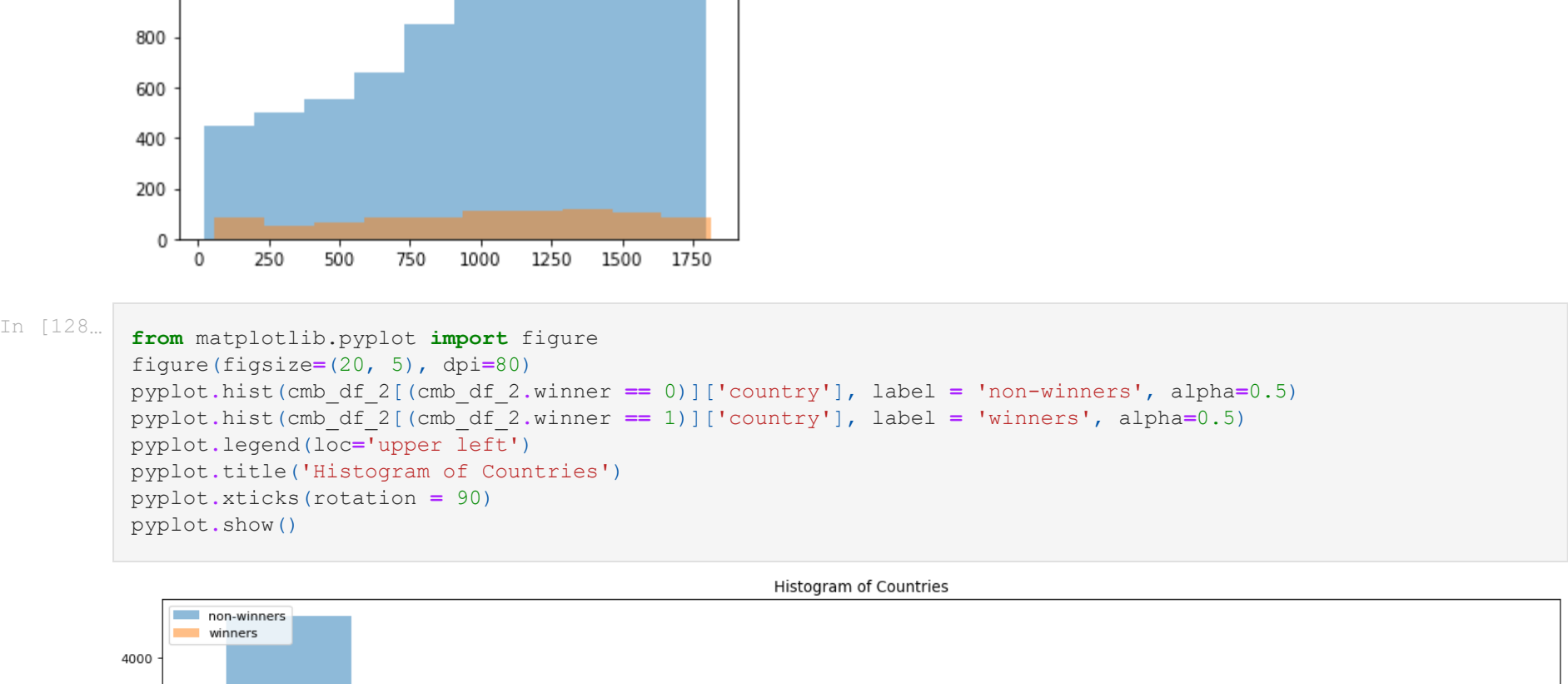
Helps answer the question: Do more colors have an impact on likelihood to win an award?

```
In [127... from matplotlib import pyplot
```

```
pyplot.hist(cmb_df_2[(cmb_df_2.winner == 0)]['unique_color_count_50'],
```

```
pyplot.legend(loc='upper left')
pyplot.title("Histogram of Unique Color Count 50 X 50 ")
pyplot.show()
```

The histogram displays the frequency of unique color counts for two groups: non-winners and winners. The x-axis represents the unique color count, and the y-axis represents the frequency. The 'non-winners' group (blue bars) shows a higher frequency of unique color counts, peaking around 1300. The 'winners' group (orange bars) shows a lower frequency of unique color counts, peaking around 1000.



The bar chart displays the annual number of publications from 1970 to 2020. The y-axis is labeled 'Number of publications' and ranges from 0 to 3000 in increments of 1000. The x-axis is labeled 'Year' and lists years from 1970 to 2020. The data shows a sharp increase in publications starting in the late 1970s, reaching a peak of approximately 2500 publications in 2010, followed by a decline to around 1000 publications by 2020.

Year	Number of publications
1970	500
1971	500
1972	500
1973	500
1974	500
1975	500
1976	500
1977	500
1978	500
1979	500
1980	500
1981	500
1982	500
1983	500
1984	500
1985	500
1986	500
1987	500
1988	500
1989	500
1990	500
1991	500
1992	500
1993	500
1994	500
1995	500
1996	500
1997	500
1998	500
1999	500
2000	500
2001	500
2002	500
2003	500
2004	500
2005	500
2006	500
2007	500
2008	500
2009	500
2010	2500
2011	2500
2012	2500
2013	2500
2014	2500
2015	2500
2016	2500
2017	2500
2018	2500
2019	2500
2020	1000



### Getting Most Frequent Color by Clustr

```
In [134]: from colorthief import ColorThief

n = 0
n2 = 0
for i, row in cmb_df_2.iterrows():
    filename = row[0]
    try:
        img = Image.open(directory3+filename+'.jpg')
        if img is not None:
            img = img.convert('RGB')
```

```
colors = color_chief.get_colors()
colors = list(colors)
colors[0] = round((colors[0] + 1) / 2, 2)
```

```

        colors[1] = round((colors[1]/256), 0)
        colors[2] = round((colors[2]/256), 0)
        hex_val = matplotlib.colors.to_hex([colors[0], colors[1], colors[2]])
        cmb_df_2.at[i, 'dominant_color'] = hex_val
        n2 = n2+1
    except:
        n = n+1

print("errors: "+str(n))
print("Records Updated: "+ str(n2))

errors: 0
Records Updated: 9815

```

```

colors_win = avg_colors_win['Average Color']
avg_colors_non.plot.bar(x='Average Color', y='counts', color= colors_non, edgecolor= 'black', title= 'Nominees Average Image Color')
avg_colors_win.plot.bar(x='Average Color', y='counts', color= colors_win, edgecolor= 'black', title= 'Winners Average Image Color')

Out[196]:
<AxesSubplot:title='center':Winners Average Image Color', xlabel='Average Color'>

```

Year	Number of Publications
2000	100
2001	3200
2002	100
2003	300
2004	100
2005	300
2006	100
2007	300
2008	100
2009	300
2010	100
2011	300
2012	100
2013	300
2014	100
2015	300
2016	100
2017	300
2018	100
2019	300

A bar chart titled "Winners Average Image Color". The x-axis is labeled "Average Color" and has seven categories: #ff0000, #000000, #0000ff, #ffff00, #00ff00, #ff00ff, and #00ff00. The y-axis is labeled "counts" and ranges from 300 to 400. There are two bars: a black bar for #ff0000 with a count of approximately 400, and a white bar for #000000 with a count of approximately 360.

Average Color	counts
#ff0000	~400
#000000	~360
#0000ff	0
#ffff00	0
#00ff00	0
#ff00ff	0
#00ff00	0

Symbol	Occurrences
#000000	220
#ffff	0
#ff0000	50
#ff00	35
#0000ff	30
#00ffff	15
#ff00ff	10
#0000ff	5

```
dom_non = cmb_df_2[(cmb_df_2.winner == 0)]
dom_win = cmb_df_2[(cmb_df_2.winner == 1)]

dom_colors_non = dom_non['dominant_color'].value_counts().rename_axis('Dominant Color').reset_index(names='count')
dom_colors_win = dom_win['dominant_color'].value_counts().rename_axis('Dominant Color').reset_index(names='count')
colors_non = dom_colors_non['Dominant Color']
colors_win = dom_colors_win['Dominant Color']

# Create a new column for the average color across the two groups
average_color = (colors_non + colors_win) / 2
```

```
dom_colors_win.plot.bar(x='Dominant Color',y='counts',color= colors_win, edgecolor = 'black', title= 'Winners')
Out[195]: <AxesSubplot:title=('center': 'Winners Dominant Image Color'), xlabel='Dominant Color'>
```

Dominant Color	counts
Blue	3900
Green	3700
Orange	3600
Purple	3500
Red	3400
Yellow	3300

Variant	Number of Reads (approx.)
000000	2200
#ttttt	100
0000tt	400
000000	300
#tttt0	300
0000tt	250
#ttttt	100
000000	50
000000	50

Dominant Color	counts
Black	420
White	370

Dominant Color	Count
#000000	150
#ffffff	150
#ff0000	45
#00ff00	30
#0000ff	25
#ff00ff	10
#00ffff	10
#0000ff	5

Save work to CSV to save time for future work