COMP4447 Final Project

Peter Strimbu

May 16, 2022

# METROPOLITAN MUSEUM OF ART DATA ANALYSIS

## COLUMNS

1. Is Highlight
2. Object ID
3. AccessionYear
4. Culture
5. Reign
6. Artist Role
7. Artist Display Bio
8. Artist Nationality
9. Artist Gender
10. Object Date
11. Medium
12. Geography Type
13. County
14. Subregion
15. Excavation
16. Rights and Reproduction
17. Metadata Date
18. Tags AAT URL

# COLUMN TYPES / RANGE / CLEANING TASKS

## OBJECT ID

- int64
- 1 through 860873
- No cleanup, ok as-is

## ACCESSIONYEAR

- 4 digit year
- Ex: nan, 1889.0, 1928, 19171917.0, 1956-08-24
- Clean: keep only first 4 digits will fix all 4 issues

## CULTURE

- category
- Ex: nan, 'reign of Amenhotep III', 'reign of Xerxes I'
- remove duplicate spaces
- possibly remove 'reign of', 'reigns of', 'or later', 'possibly', 'or slightly later', ',early', ',probably', (anything after comma), (anything in parens), (question marks), (space at end), (split on slash,dash/'and'/'or' and process individually), 'to xyz' and group by remaining names

## ARTIST ROLE

- object
- Ex: 'nan', 'Maker', 'Designer|Manufacturer'
- Clean: split on vert bar?

## ARTIST DISPLAY BIO

- object
- Ex: '1785–1844', 'nan', 'British, London 1873–1952 Hailsham, Sussex|British, Wiltshire 1877–1952 Oxford', '||||Female|Female'
- Issues: not sure what the vert bars separate
- None - don't use or use as-is

## OBJECT DATE

- Numeric
- Ex: '1853', '1901', '1909–27', '1782-1784', 'December 1, 1925', 'after 1773'
- Medium
- object
- Ex: 'Gold' 'Silver' 'Bronze or copper' ... 'Overlay for 23.112.2889, graphite and ink on glazed linen tracing paper'
- Don't use or use as-is

## AccessionYear

- 4 digit year
- Ex: nan, 1889.0, 1928, 19171917.0, 1956-08-24
- Clean: keep only first 4 digits will fix all 4 issues

```
1  df_clean['AccessionYear'] = df_clean['AccessionYear'].str[:4]
2  #df_clean['AccessionYear'].dropna()
3
```
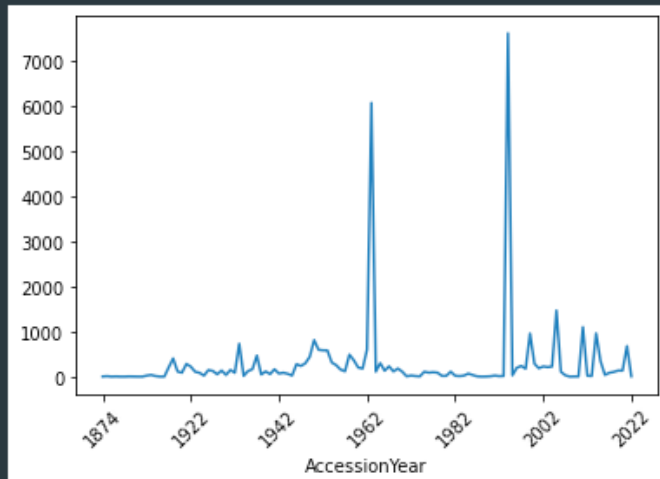
[111]  ✓  0.1s

```
1  # plot accession years
2  counts = df_clean['AccessionYear'].dropna().groupby([df['AccessionYear']]).count()
3  # counts.sort_values()
4  counts.plot(rot=45)
```

[110]  ✓  0.2s

... `<AxesSubplot:xlabel='AccessionYear'>`

</>

## Culture

- category
- Ex: nan, 'reign of Amenhotep III', 'reign of Xerxes I'
- remove duplicate spaces
- possibly remove 'reign of', 'reigns of', 'or later', 'possibly', 'or slightly later', ',early', ',probably', (anything after comma), (split on slash,dash/'and'/'or' and process individually), 'to xyz' and group by remaining names

```
1   df_clean['Culture'].astype(str)
2   df_clean['Culture'].replace('reign. of', '', regex=True, inplace=True)
3   df_clean['Culture'].replace('or.*later', '', regex=True, inplace=True)
4   df_clean['Culture'].replace('early', '', regex=False, inplace=True)
5   df_clean['Culture'].replace('or', '', regex=False, inplace=True)
6   df_clean['Culture'].replace('probably', '', regex=False, inplace=True)
7   df_clean['Culture'].replace('possibly', '', regex=False, inplace=True)
8   df_clean['Culture'].replace('Possibly', '', regex=False, inplace=True)
9   df_clean['Culture'].replace(';', '', regex=False, inplace=True)
10  df_clean['Culture'].replace('\(.*\)', '', regex=True, inplace=True)
11  df_clean['Culture'].replace('\?', '', regex=True, inplace=True)
12  df_clean['Culture'] = df_clean['Culture'].dropna().str.strip()
13  df_clean['Culture'] = df_clean['Culture'].dropna().str.lower()
14  # for culture in df_clean['Culture'].unique():
15  #     print(culture)
16
```

[100]   ✓  3.1s

```
1   # generate a WordCloud for "Culture"
2
3   stopwords = set(STOPWORDS)
4
5   word_list = lines_to_words(df_clean['Culture'].dropna().tolist())
6
7   for i in range(len(word_list)):
8       word_list[i] = word_list[i].lower()
9
10  word_string = ''
11  word_string += " ".join(word_list)+" "
12
13  # wordcloud = WordCloud(collocations=False, width=800, height=800, background_color='white', stopwords=stopwords
14  # plt.figure(figsize=(10,10), dpi=100)
15  # plt.imshow(wordcloud, interpolation='bilinear')
16  # plt.axis("off")
17  # plt.tight_layout(pad=0)
18
19  du_mask = np.array(Image.open('University-of-Denver-logo.png'))
20  colors = ImageColorGenerator(du_mask)
21  wordcloud = WordCloud(collocations=False, stopwords=stopwords, mask=du_mask, mode='RGB', background_color=None,
22  plt.figure(figsize=(15,15))
23  plt.imshow(wordcloud)
24  plt.axis('off')
25  # plt.title('')
26  plt.show()
```

# Artist Role

- object
- Ex: 'nan', 'Maker', 'Designer|Manufacturer'
- Clean: split on vert bar?

```python
artist_roles = lines_to_words(df_clean['Artist Role'].dropna().tolist())
artist_roles_clean = []
for i in range(len(artist_roles)):
    artist_roles_clean.extend(artist_roles[i].lower().split('|'))

artist_roles_df = pd.DataFrame(artist_roles_clean, columns=['Artist Role'])

#calculate sum of values by group
df_groups = artist_roles_df.groupby(['Artist Role'])['Artist Role'].sum()

artist_roles_df = artist_roles_df['Artist Role'].dropna().groupby([artist_roles_df['Artist Role']]).count().rese
artist_roles_df['percentage'] = artist_roles_df['counts'] / artist_roles_df['counts'].sum()
top10 = artist_roles_df.sort_values(by='counts', ascending=False).head(10)
top10.plot.bar(x='Artist Role',y='percentage', rot=90)
```
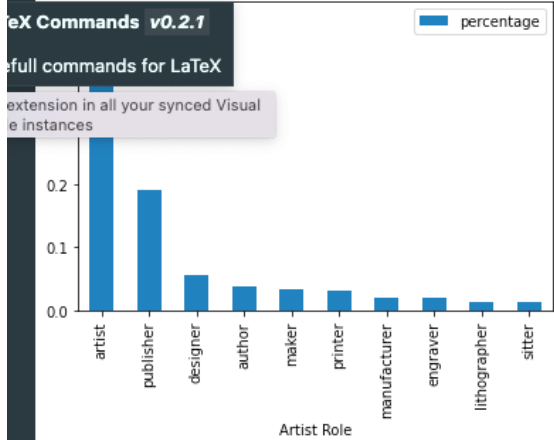
✓  10.9s                                                                    Python

<AxesSubplot:xlabel='Artist Role'>

## Object Date

- Numeric
- Ex: '1853', '1901', '1909–27', '1782-1784', 'December 1, 1925', 'after 1773'
- Clean: select first 4 digit number

```python
df_clean['Object Date'] = df_clean['Object Date'].str.extract('([0-9]{4})')
object_dates_df = df_clean['Object Date'].dropna().groupby([df_clean['Object Date']]).count().reset_index(name='

object_dates_df.plot(kind='line',x='Object Date', y='count')
```

[168]  ✓  1.2s                                                                              Python

...  `<AxesSubplot:xlabel='Object Date'>`