IM-UH 1511 **Introduction to Digital Humanities**

# HOMEWORK 7a   ¶

# Best Korean Movies in IMDb: Descriptive Statistics

## 25 points totally

```
In [1]: import pandas as pd
        import numpy as np
        import networkx as nx
        import pygraphviz
        from networkx.drawing.nx_agraph import graphviz_layout
        from networkx.drawing.nx_agraph import to_agraph
        import matplotlib.pyplot as plt
        import matplotlib as mpl
        from pylab import hist
        import random
        from collections import Counter
        import operator
        import itertools
        from wordcloud import WordCloud
        import warnings
        warnings.filterwarnings("ignore", category=RuntimeWarning)
        warnings.filterwarnings("ignore", category=UserWarning)
        warnings.simplefilter('ignore')
```

In [2]: 
```python
allfilms = pd.read_csv('300KoreanFilms.csv', sep=',', encoding="utf-8")
films=allfilms.sample(100)
films = films.reset_index(drop=True)
print(len(films))
films
```

100

Out[2]:

|    | TITLE | YEAR | DIRECTOR | STARS | GENRE |
|----|-------|------|----------|-------|-------|
| 0 | Parasite | 2019 | Joon-ho Bong (D) | Kang-ho Song, Sun-kyun Lee, Yeo-jeong Jo | Comedy, Drama, Thriller |
| 1 | Tazza: The High Rollers | 2006 | Dong-hoon Choi | Seung-woo Cho, Yun-shik Baek, Hye-su Kim, Hae-... | Comedy, Crime |
| 2 | Hwasango | 2001 | Tae-gyun Kim | Hyuk Jang, Min-a Shin, Su-ro Kim, Sang-Woo Kwon | Action, Comedy, Fantasy |
| 3 | As One | 2012 | Hyeon-seong Moon | Ji-won Ha, Doona Bae, Han Yeri, Yoon-young Choi | Drama, Sport |
| 4 | Perfect Number | 2012 | Eun-jin Pang (D) | Seung-bum Ryoo, Yo-won Lee, Jin-woong Cho, Yoo... | Drama, Thriller |
| ... | ... | ... | ... | ... | ... |
| 95 | Lost Flower Eo Woo-dong | 2015 | Soo Sung Lee | Do-bin Baek, Eun-pi Kang, Wook Han Yeo, Kyung-... | Drama |

In [3]: 
```python
titles=films.TITLE.tolist()
print(len(titles),len(set(titles)))
```

100 99

In [4]:
```python
year=films.YEAR.tolist()
year=list(set(year))
print(len(year),min(year),max(year))
sorted(year)
```

```
21 1998 2019
```

Out[4]:
```
[1998,
 1999,
 2000,
 2001,
 2002,
 2003,
 2004,
 2005,
 2006,
 2007,
 2008,
 2009,
 2010,
 2011,
 2012,
 2013,
 2014,
 2015,
 2016,
 2017,
 2019]
```

In [5]:
```python
director=films.DIRECTOR.tolist()
director=sorted(list(set(director)))
print(len(director))
director
```

```
76
```

In [6]:
```python
s_films=pd.DataFrame()
for i in range(len(films)):
    d=dict(films.loc[i])
    sl= films.STARS[i].split(",")
    gl=films.GENRE[i].split(", ")
    for s in sl:
        for t in gl:
            d["ACTOR"]=s.strip()
            d['UGENRE']=t
            s_films=s_films.append(d,ignore_index=True)
s_films['YEAR'] = s_films['YEAR'].astype(int)
s_films.rename(columns={'GENRE':'MGENRE','UGENRE':'GENRE'}, inplace=True)
s_films=s_films[["YEAR","TITLE","DIRECTOR","ACTOR","GENRE"]]
df=s_films
print(len(df))
df.sort_values(by="YEAR").head(20)
```

920

Out[6]:

|     | YEAR | TITLE | DIRECTOR | ACTOR | GENRE |
|-----|------|-------|----------|-------|-------|
| 728 | 1998 | Christmas in August | Jin-ho Hur | Suk-kyu Han | Drama |
| 735 | 1998 | Christmas in August | Jin-ho Hur | Ji-hye Oh | Romance |
| 734 | 1998 | Christmas in August | Jin-ho Hur | Ji-hye Oh | Drama |
| 733 | 1998 | Christmas in August | Jin-ho Hur | Goo Shin | Romance |
| 732 | 1998 | Christmas in August | Jin-ho Hur | Goo Shin | Drama |
| 730 | 1998 | Christmas in August | Jin-ho Hur | Eun-ha Shim | Drama |
| 729 | 1998 | Christmas in August | Jin-ho Hur | Suk-kyu Han | Romance |
| 731 | 1998 | Christmas in August | Jin-ho Hur | Eun-ha Shim | Romance |
| 217 | 1999 | Memento Mori | Tae-yong Kim | Yeong-jin Lee | Drama |
| 361 | 1999 | Memento Mori | Kyu-dong Min | Yeong-jin Lee | Drama |
| 360 | 1999 | Memento Mori | Kyu-dong Min | Yeong-jin Lee | Romance |

In [7]:
```python
df.to_csv('Random100KoreanFilms.csv',encoding='utf-8')
```

In [8]:
```python
actor=df.ACTOR.tolist()
actor=sorted(list(set(actor)))
print(len(actor))
actor
```

291

In [9]:
```python
unique_genre=df.GENRE.tolist()
unique_genre=sorted(list(set(unique_genre)))
print(len(unique_genre))
unique_genre
```

17

Out[9]:
```
['Action',
 'Adventure',
 'Animation',
 'Biography',
 'Comedy',
 'Crime',
 'Drama',
 'Family',
 'Fantasy',
 'History',
 'Horror',
 'Music',
 'Mystery',
 'Romance',
 'Sci-Fi',
 'Sport',
 'Thriller']
```

In [10]:
```python
title=df.TITLE.tolist()
title=sorted(list(set(title)))
print(len(title))
title
```
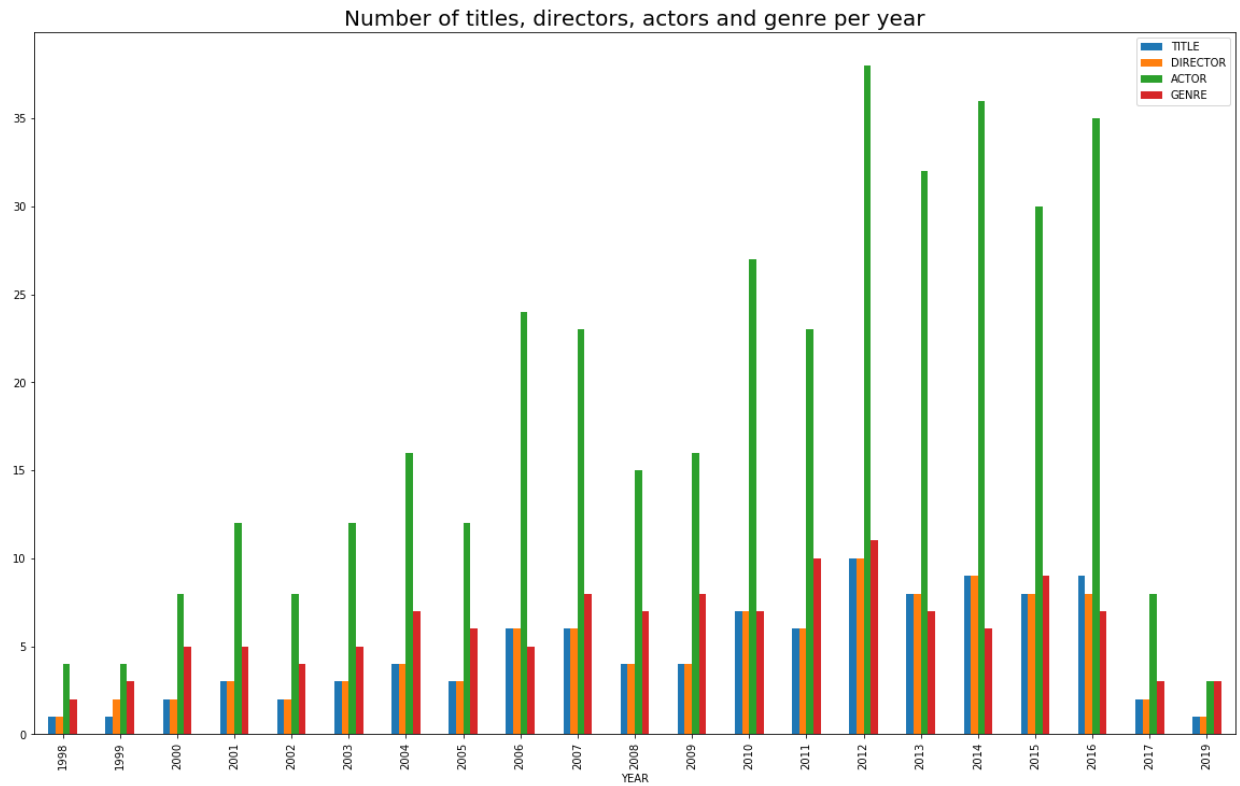
99

# Grouping per Year

In [11]:
```python
gdf=df.groupby("YEAR").nunique()[["TITLE","DIRECTOR","ACTOR","GENRE"]]
gdf = gdf.reset_index()
gdf
```
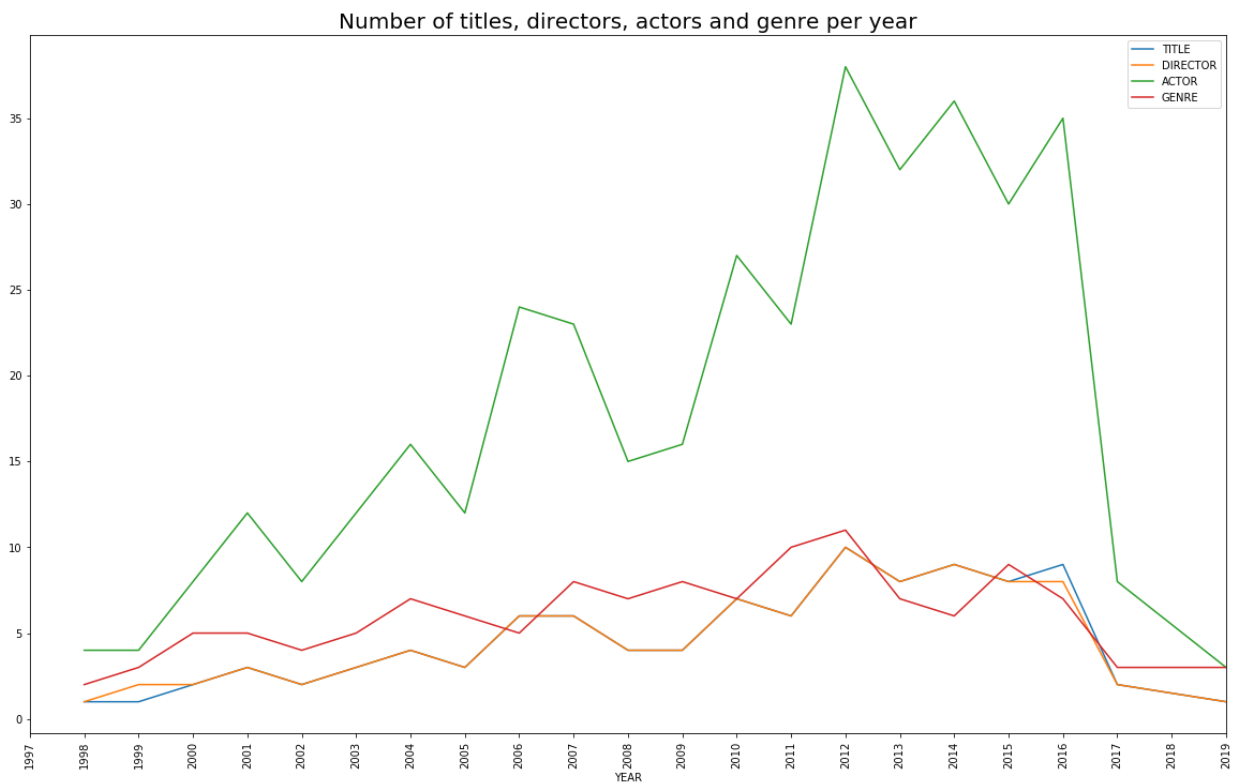
Out[11]:

|     | YEAR | TITLE | DIRECTOR | ACTOR | GENRE |
| --- | ---- | ----- | -------- | ----- | ----- |
| 0   | 1998 | 1     | 1        | 4     | 2     |
| 1   | 1999 | 1     | 2        | 4     | 3     |
| 2   | 2000 | 2     | 2        | 8     | 5     |
| 3   | 2001 | 3     | 3        | 12    | 5     |
| 4   | 2002 | 2     | 2        | 8     | 4     |
| 5   | 2003 | 3     | 3        | 12    | 5     |
| 6   | 2004 | 4     | 4        | 16    | 7     |
| 7   | 2005 | 3     | 3        | 12    | 6     |
| 8   | 2006 | 6     | 6        | 24    | 5     |
| 9   | 2007 | 6     | 6        | 23    | 8     |
| 10  | 2008 | 4     | 4        | 15    | 7     |
| 11  | 2009 | 4     | 4        | 16    | 8     |
| 12  | 2010 | 7     | 7        | 27    | 7     |
| 13  | 2011 | 6     | 6        | 23    | 10    |
| 14  | 2012 | 10    | 10       | 38    | 11    |
| 15  | 2013 | 8     | 8        | 32    | 7     |
| 16  | 2014 | 9     | 9        | 36    | 6     |
| 17  | 2015 | 8     | 8        | 30    | 9     |
| 18  | 2016 | 9     | 8        | 35    | 7     |
| 19  | 2017 | 2     | 2        | 8     | 3     |
| 20  | 2019 | 1     | 1        | 3     | 3     |

In [12]:
```python
ax=gdf.plot(x='YEAR', y=["TITLE", "DIRECTOR", "ACTOR", "GENRE"], kind="bar"
ax.set_title('Number of titles, directors, actors and genre per year', font
```

```
In [13]: ax=gdf.plot(x='YEAR', y=["TITLE", "DIRECTOR", "ACTOR", "GENRE"], kind="line
         plt.xticks(np.arange(1997, 2020, step=1),rotation='vertical');
         ax.set_title('Number of titles, directors, actors and genre per year', font
```
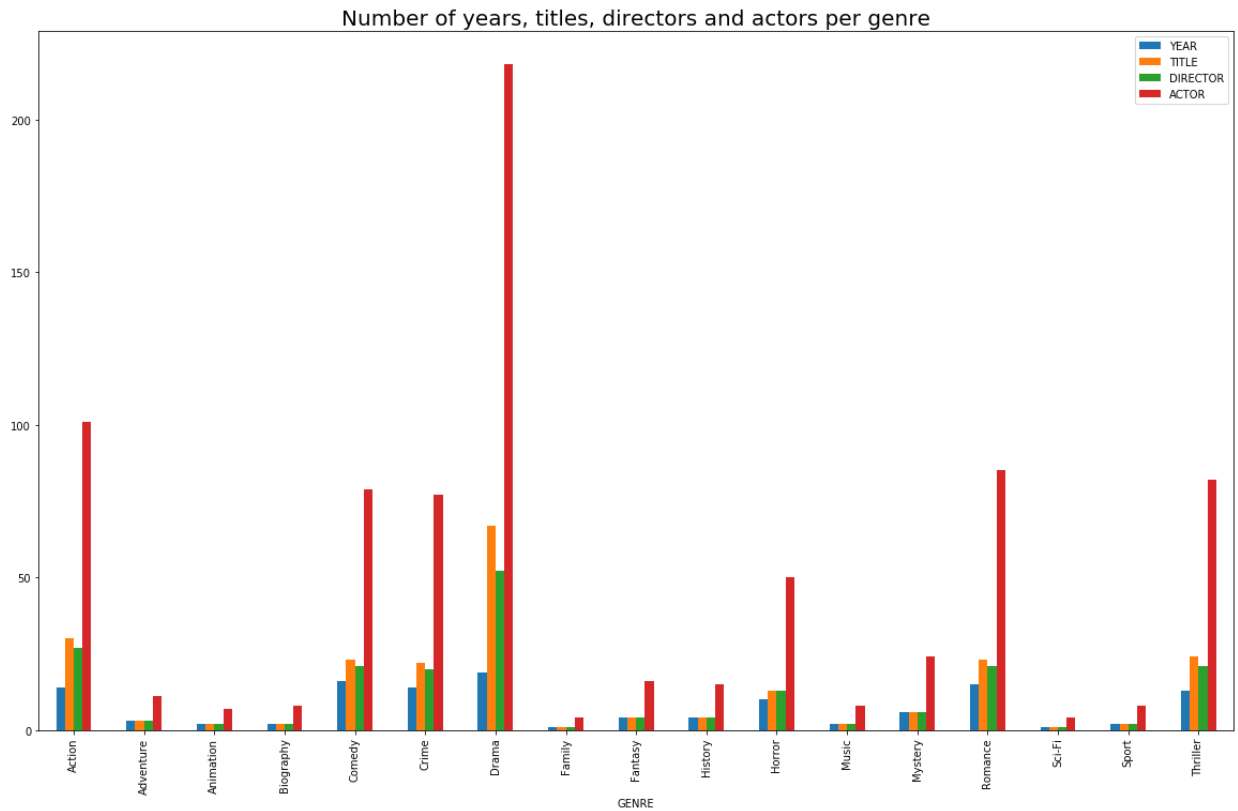


## Grouping per Genre

```
In [14]: ggdf=df.groupby("GENRE").nunique()[["YEAR","TITLE","DIRECTOR","ACTOR"]]
         ggdf = ggdf.reset_index()
         ggdf
```
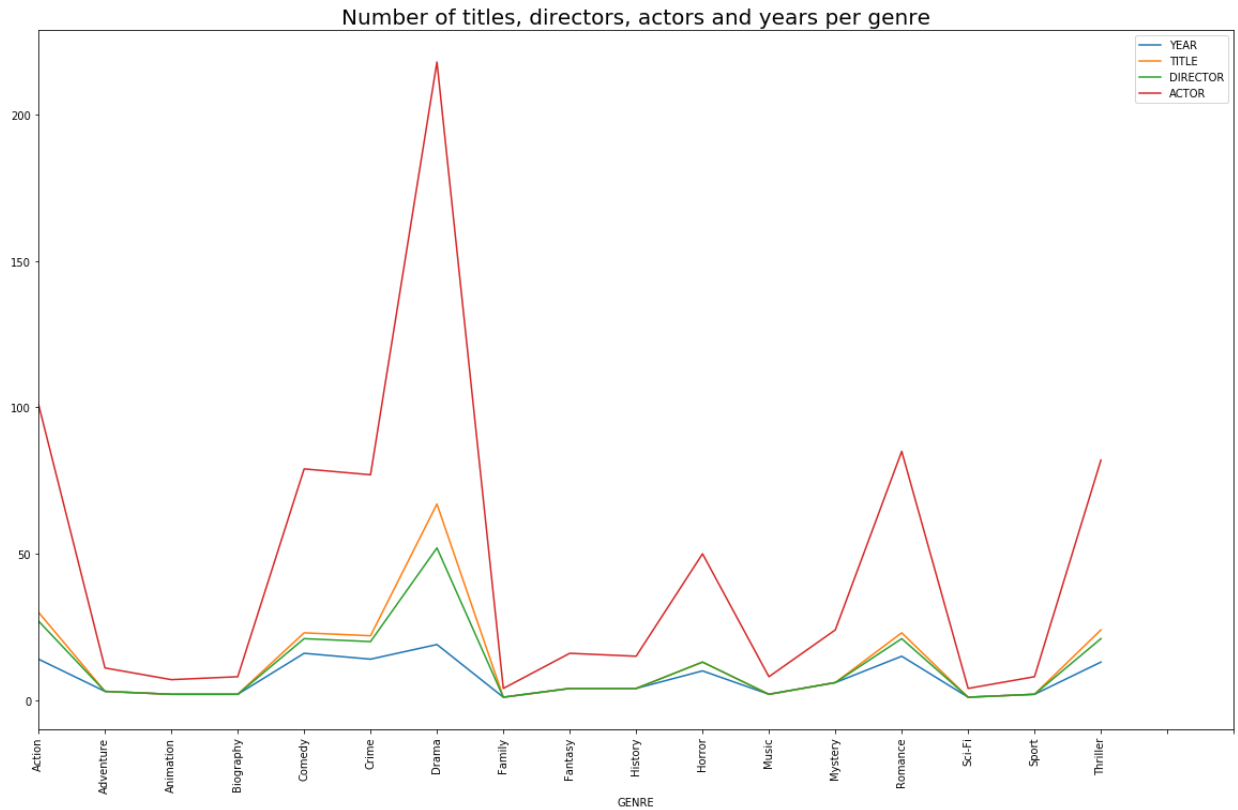
Out[14]:

|    | GENRE     | YEAR | TITLE | DIRECTOR | ACTOR |
|----|-----------|------|-------|----------|-------|
| 0  | Action    | 14   | 30    | 27       | 101   |
| 1  | Adventure | 3    | 3     | 3        | 11    |
| 2  | Animation | 2    | 2     | 2        | 7     |
| 3  | Biography | 2    | 2     | 2        | 8     |
| 4  | Comedy    | 16   | 23    | 21       | 79    |
| 5  | Crime     | 14   | 22    | 20       | 77    |
| 6  | Drama     | 19   | 67    | 52       | 218   |
| 7  | Family    | 1    | 1     | 1        | 4     |
| 8  | Fantasy   | 4    | 4     | 4        | 16    |
| 9  | History   | 4    | 4     | 4        | 15    |
| 10 | Horror    | 10   | 13    | 13       | 50    |
| 11 | Music     | 2    | 2     | 2        | 8     |
| 12 | Mystery   | 6    | 6     | 6        | 24    |
| 13 | Romance   | 15   | 23    | 21       | 85    |
| 14 | Sci-Fi    | 1    | 1     | 1        | 4     |
| 15 | Sport     | 2    | 2     | 2        | 8     |
| 16 | Thriller  | 13   | 24    | 21       | 82    |

In [15]:
```python
ax=ggdf.plot(x='GENRE', y=["YEAR","TITLE", "DIRECTOR", "ACTOR"], kind="bar"
ax.set_title('Number of years, titles, directors and actors per genre', fon
```



Number of years, titles, directors and actors per genre

```
In [16]: ax=ggdf.plot(x='GENRE', y=["YEAR", "TITLE", "DIRECTOR", "ACTOR"], kind="lin
         ax.set_xticks(np.arange(0, 19, step=1))
         ax.set_xticklabels(unique_genre,rotation='vertical')
         ax.set_title('Number of titles, directors, actors and years per genre', fon
```



## Grouping per Director

```
In [17]: dgdf=df.groupby("DIRECTOR").nunique()[["YEAR","TITLE","ACTOR",'GENRE']]
         dgdf = dgdf.reset_index()
         dgdf.sort_values(by="TITLE",ascending=False)
```
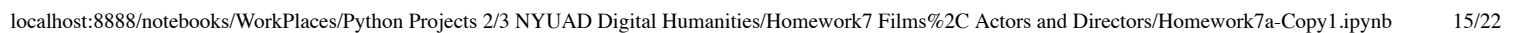
Out[17]:

|    | DIRECTOR | YEAR | TITLE | ACTOR | GENRE |
|----|----------|------|-------|-------|-------|
| 43 | Ki-duk Kim (D) | 5 | 5 | 20 | 5 |
| 4  | Chan-wook Park | 3 | 4 | 15 | 7 |
| 36 | Jee-woon Kim | 4 | 4 | 14 | 8 |
| 63 | Tae-yong Kim | 3 | 3 | 12 | 3 |
| 39 | Joon-ho Bong (D) | 3 | 3 | 10 | 5 |
| ... | ... | ... | ... | ... | ... |
| 33 | Jae-eun Jeong | 1 | 1 | 4 | 1 |
| 34 | Jae-rim Han | 1 | 1 | 4 | 3 |
| 37 | Ji-hoon Kim | 1 | 1 | 4 | 3 |
| 1  | Byeong-heon Lee | 1 | 1 | 4 | 3 |
| 75 | Yun-hyeon Jang | 1 | 1 | 4 | 3 |

76 rows × 5 columns

```python
In [18]: close = dgdf
         mpl.rcParams['font.size'] = 12.0
         fig = plt.figure(figsize=(20,20))
         plt.pie(
             close['TITLE'],
             labels=close['DIRECTOR'],
             shadow=False,
             startangle=90,
             # with the percent listed as a fraction
             autopct='%1.1f%%',
             rotatelabels=True,labeldistance=1.01,textprops={'fontsize':9},pctdistan
             )

         # View the plot drop above
         plt.axis('equal')

         # View the plot
         ss="Number of Titles per Director"
         plt.suptitle(ss,fontsize=25)
         plt.tight_layout(rect=[0, 0, 1, 0.93])
         plt.show()
```
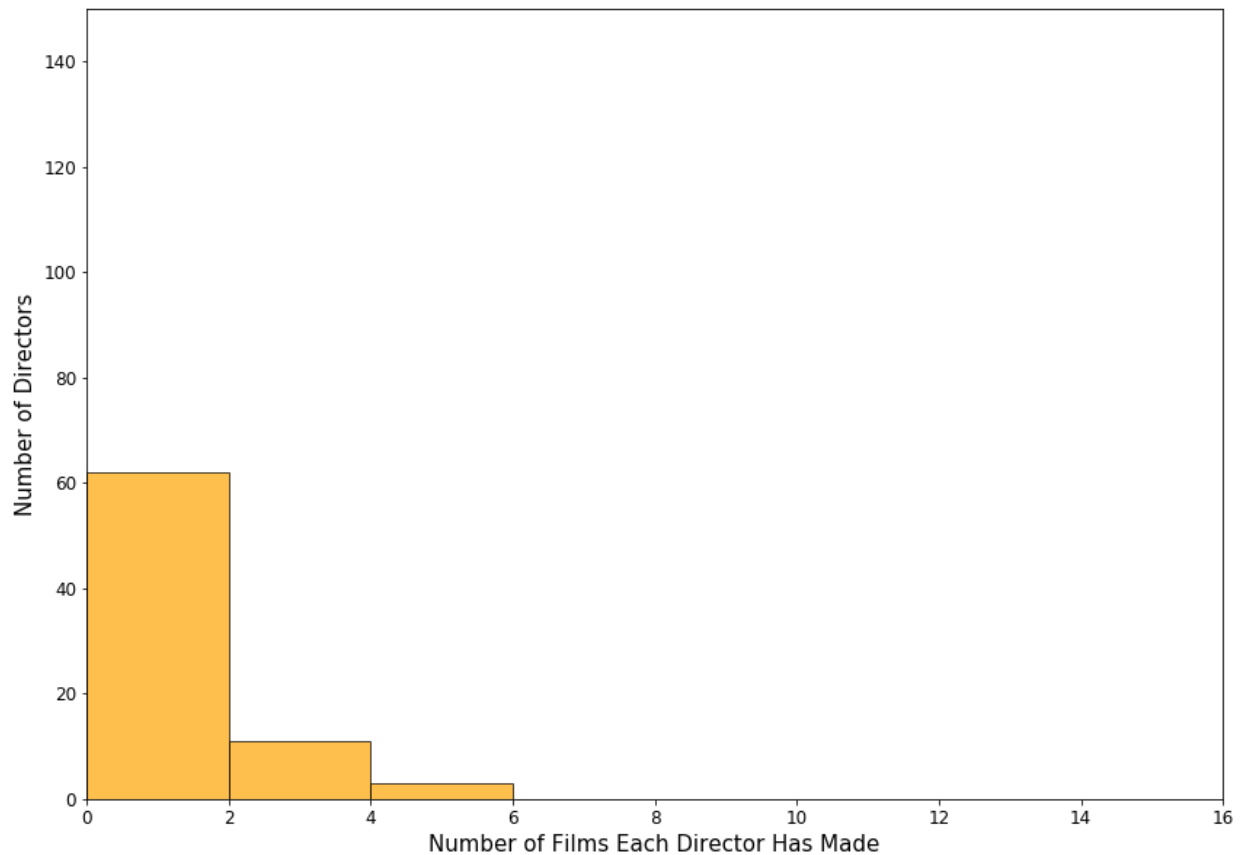
## Number of Titles per Director

In [19]:
```python
x=dgdf.TITLE.tolist()
fig = plt.figure(figsize=(14,10))
plt.hist(x,bins = [0,2,4,6,8,10,12,14,16],edgecolor='black',color="orange",
plt.axis([0, 16, 0, 150])
plt.xlabel("Number of Films Each Director Has Made",fontsize=15)
plt.ylabel("Number of Directors",fontsize=15)
plt.xticks([0,2,4,6,8,10,12,14,16])
plt.suptitle("Histogram of Directors in Films",fontsize=15)
plt.show()
```

Histogram of Directors in Films
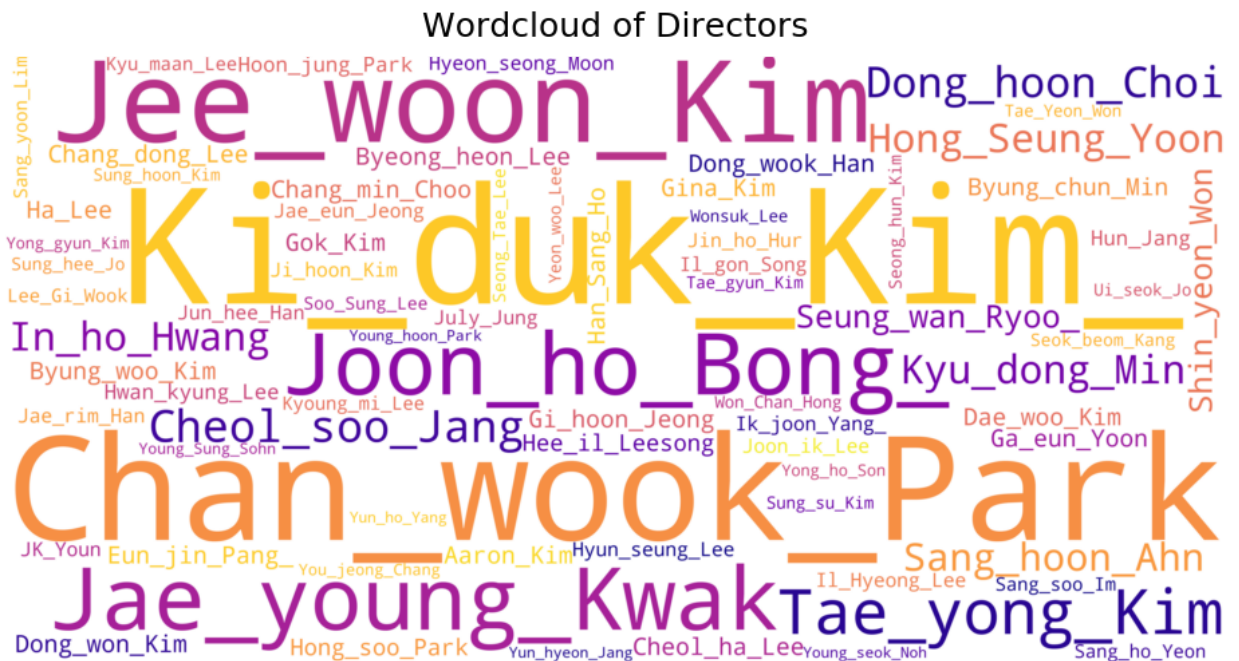
```
In [20]: subsetd = dgdf[['DIRECTOR', 'TITLE']]
         tuplesd = [tuple(x) for x in subsetd.values]

         t=[]
         for (i,j) in tuplesd:
             for k in range(j):
         #         print(i.replace(" ","_").replace("-","_"))
                 t.append(i.replace(" ","_").replace("-","_"))
         ttd=' '.join(t)

         wordcloud = WordCloud(collocations=False,background_color="white",colormap=
         fig = plt.figure(figsize=(13,13))
         default_colors = wordcloud.to_array()
         plt.imshow(default_colors, interpolation="bilinear")
         plt.axis("off")
         ss="Wordcloud of Directors"
         plt.suptitle(ss,fontsize=25)
         plt.tight_layout(rect=[0, 0, 1, 1.4])
         plt.show()
```



Wordcloud of Directors

## Grouping per Actor

In [21]:
```python
agdf=df.groupby("ACTOR").nunique()[["YEAR","DIRECTOR","TITLE",'GENRE']]
agdf = agdf.reset_index()
agdf.sort_values(by="TITLE",ascending=False)
```
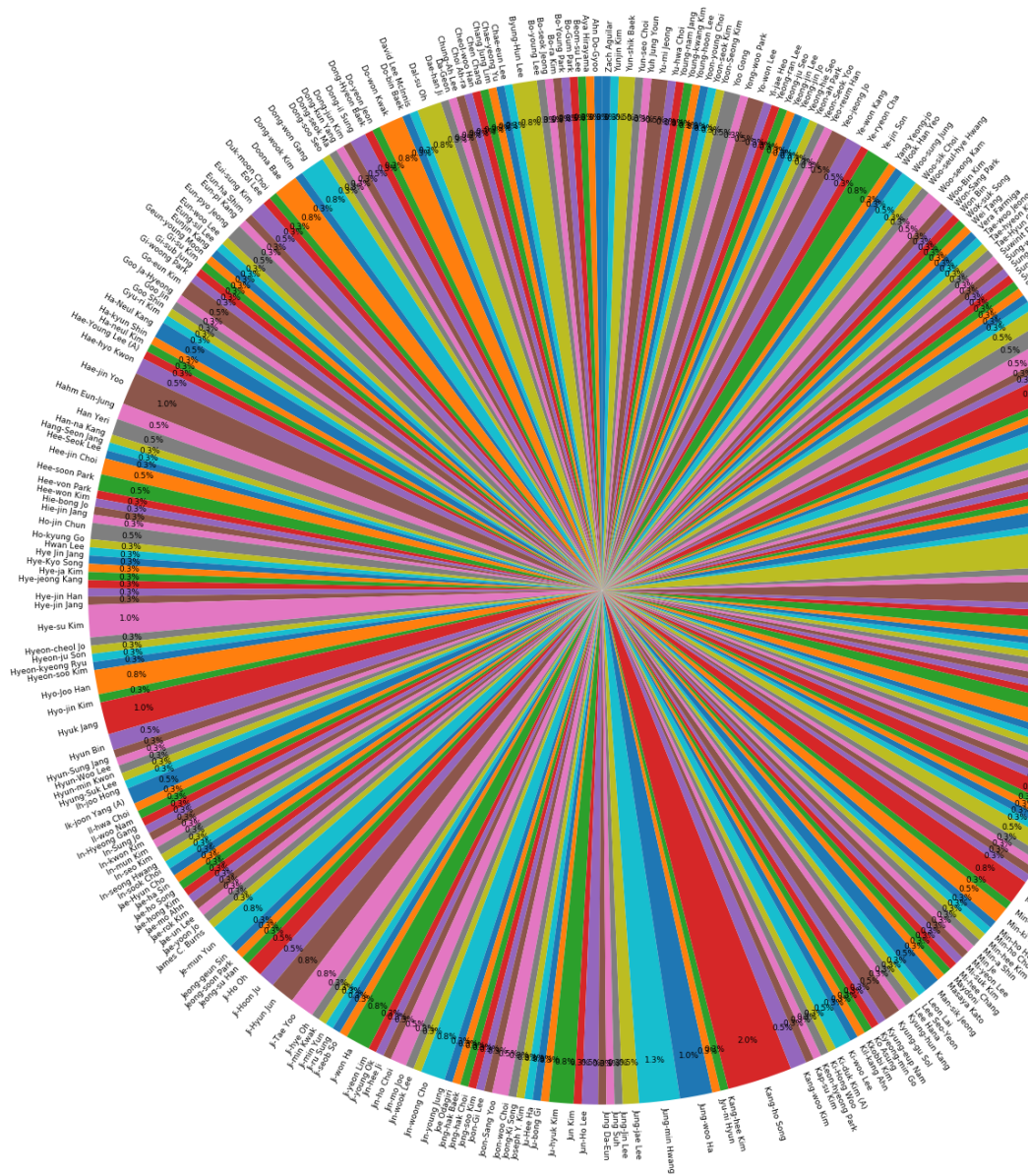
Out[21]:

|     | ACTOR | YEAR | DIRECTOR | TITLE | GENRE |
|-----|-------|------|----------|-------|-------|
| 153 | Kang-ho Song | 7 | 6 | 8 | 8 |
| 149 | Jung-min Hwang | 4 | 5 | 5 | 4 |
| 218 | Seung-ryong Ryu | 5 | 5 | 5 | 7 |
| 55  | Hae-jin Yoo | 4 | 3 | 4 | 5 |
| 76  | Hye-su Kim | 4 | 4 | 4 | 6 |
| ... | ... | ... | ... | ... | ... |
| 111 | Jeong-soon Park | 1 | 1 | 1 | 2 |
| 112 | Jeong-su Han | 1 | 1 | 1 | 2 |
| 117 | Ji-hye Oh | 1 | 1 | 1 | 2 |
| 118 | Ji-min Kwak | 1 | 1 | 1 | 1 |
| 290 | Zach Aguilar | 1 | 1 | 1 | 3 |

291 rows × 5 columns

```
In [22]: aclose = agdf
         mpl.rcParams['font.size'] = 12.0
         fig = plt.figure(figsize=(20,20))
         plt.pie(
             aclose['TITLE'],
             labels=aclose['ACTOR'],
             shadow=False,
             startangle=90,
             # with the percent listed as a fraction
             autopct='%1.1f%%',
             rotatelabels=True,labeldistance=1.01,textprops={'fontsize':9},pctdistan
             )

         # View the plot drop above
         plt.axis('equal')

         # View the plot
         ss="Number of Titles per Actor"
         plt.suptitle(ss,fontsize=25)
         plt.tight_layout(rect=[0, 0, 1, 0.93])
         plt.show()
```
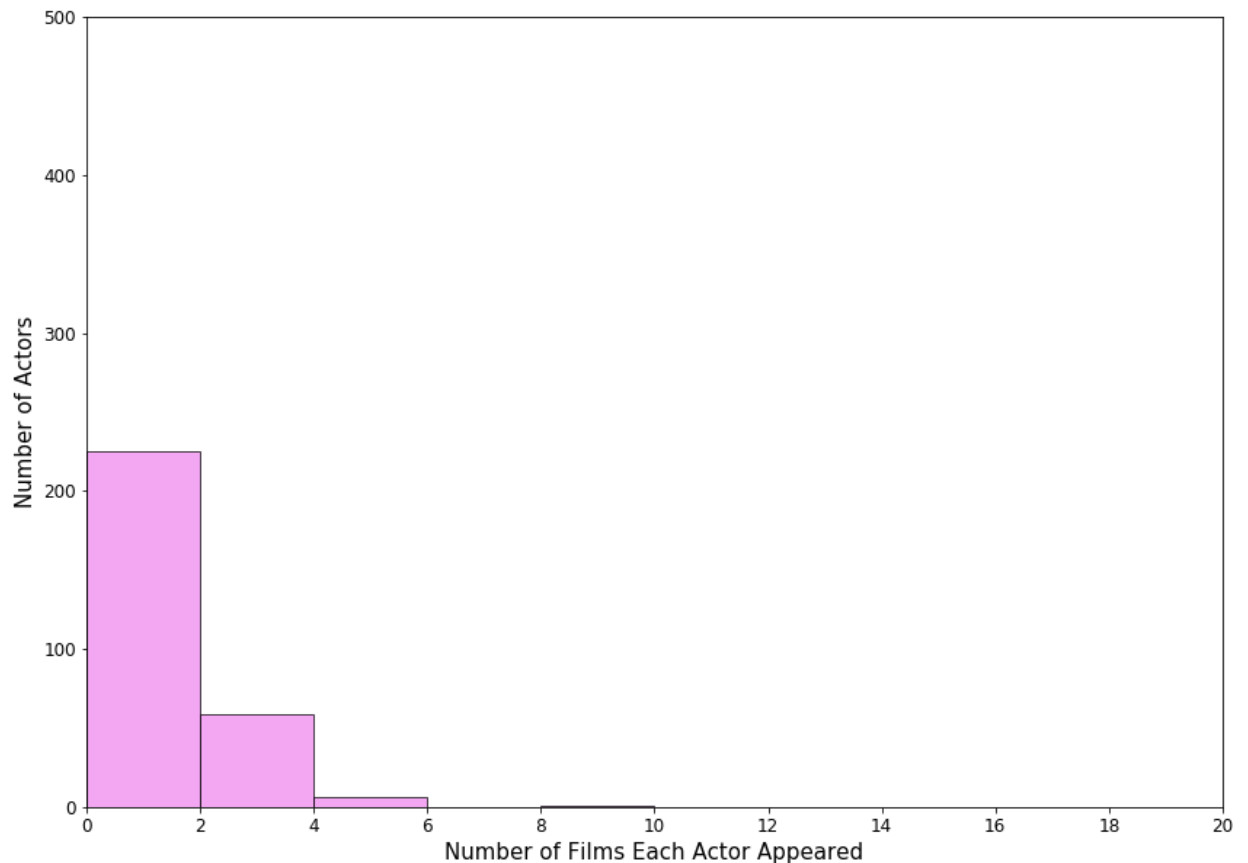
## Number of Titles per Actor

```
In [23]: x=agdf.TITLE.tolist()
         fig = plt.figure(figsize=(14,10))
         plt.hist(x,bins = [0,2,4,6,8,10,12,14,16,18,20],edgecolor='black',color="vi
         plt.axis([0, 20, 0, 500])
         plt.xlabel("Number of Films Each Actor Appeared",fontsize=15)
         plt.ylabel("Number of Actors",fontsize=15)
         plt.xticks([0,2,4,6,8,10,12,14,16,18,20])
         plt.suptitle("Histogram of Actors in Films",fontsize=15)
         plt.show()
```

Histogram of Actors in Films

In [24]:
```python
subset = agdf[['ACTOR', 'TITLE']]
tuples = [tuple(x) for x in subset.values]

t=[]
for (i,j) in tuples:
    for k in range(j):
#         print(i.replace(" ","_").replace("-","_"))
        t.append(i.replace(" ","_").replace("-","_"))
tt=' '.join(t)

wordcloud = WordCloud(collocations=False,background_color="white",colormap=
fig = plt.figure(figsize=(13,13))
default_colors = wordcloud.to_array()
plt.imshow(default_colors, interpolation="bilinear")
plt.axis("off")
ss="Wordcloud of Actors"
plt.suptitle(ss,fontsize=25)
plt.tight_layout(rect=[0, 0, 1, 1.4])
plt.show()
```



Wordcloud of Actors

In [ ]: