

✓ Load Data

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 url = "https://raw.githubusercontent.com/rfordatascience/tidytuesday/master/data/age_gaps/age_gaps.csv"
6 age_gaps = pd.read_csv(url)
7
8 age_gaps.columns = age_gaps.columns.str.lower()
9
10 # parse dates
11 date_cols = [c for c in age_gaps.columns if "birthdate" in c]
12 for c in date_cols:
13     age_gaps[c] = pd.to_datetime(age_gaps[c], errors="coerce")
14
15 age_gaps.info()
16 age_gaps.head()
```

[Show hidden output](#)

Next steps:

[Generate code with age_gaps](#)

[New interactive sheet](#)

✓ Data Processing

```
1 age_gaps.isna().sum()
```

	0
movie_name	0
release_year	0
director	0
age_difference	0
couple_number	0
actor_1_name	0
actor_2_name	0
character_1_gender	0
character_2_gender	0
actor_1_birthdate	0
actor_2_birthdate	0
actor_1_age	0
actor_2_age	0

dtype: int64

1 age_gaps.describe()

	release_year	age_difference	couple_number	actor_1_birthdate	actor_2_birthdate
count	1155.000000	1155.000000	1155.000000	1155	1155
mean	2000.799134	10.424242	1.398268	1960-09-07 07:15:07.012987008	1960-09-07 07:15:07.012987008
min	1935.000000	0.000000	1.000000	1889-04-16 00:00:00	1889-04-16 00:00:00
25%	1997.000000	4.000000	1.000000	1953-05-16 00:00:00	1953-05-16 00:00:00
50%	2004.000000	8.000000	1.000000	1964-10-03 00:00:00	1964-10-03 00:00:00
75%	2012.000000	15.000000	2.000000	1973-08-07 12:00:00	1973-08-07 12:00:00
max	2022.000000	52.000000	7.000000	1996-06-01 00:00:00	1996-06-01 00:00:00
std	16.365819	8.511086	0.754419		NaN

I add a new feature:

- release_to_2026
- younger_age
- order_age
- older_gender

- meam_age
- relative_age

```

1 # copy age_gaps
2 age_gaps_copy = age_gaps.copy()
3 # drop column = 'movie_name', 'director', 'couple_numbe
4 age_gaps_copy = age_gaps_copy.drop(columns=['movie_name
5 # add a new column named "release to 2026" (2026 - rele
6 age_gaps_copy['release_to_2026'] = 2026 - age_gaps_copy
7
8 # compare actor_1_age, actor_2_age from each row and as
9 age_gaps_copy['younger_age'] = age_gaps_copy[['actor_1_
10 age_gaps_copy['order_age'] = age_gaps_copy[['actor_1_ag
11
12 age_gaps_copy['older_gender'] = np.where(
13     age_gaps_copy['actor_1_age'] > age_gaps_copy['actor
14     np.where(
15         age_gaps_copy['actor_2_age'] > age_gaps_copy['a
16         'same_age'))
17
18 # add a column named 'mean_age' = (younger_age+order_ag
19 age_gaps_copy['mean_age'] = (age_gaps_copy['younger_age
20
21 # add new column named " relative_age_gap"
22 age_gaps_copy['relative_age_gap'] = age_gaps_copy['age_
23
24 age_gaps_clean = age_gaps_copy.drop(columns=['actor_1_a
25 age_gaps_clean.head().reset_index()

```

	index	release_year	age_difference	couple_number	character_1_gend
0	0	1971	52	1	wom
1	1	2006	50	1	m
2	2	2002	49	1	m
3	3	1998	45	1	m
4	4	2010	43	1	m

Next steps: [Generate code with age_gaps_clean](#)

[New interactive sheet](#)

✓ EDA

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt

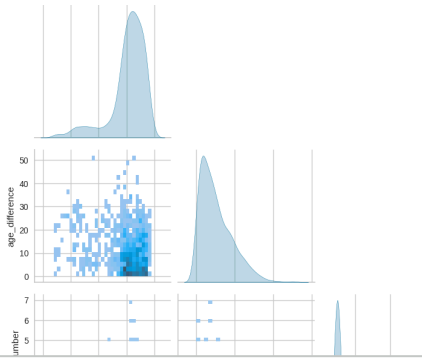
```

```

1 sns.pairplot(age_gaps_clean, kind="hist", diag_kind="kde", corner:

```

<seaborn.axisgrid.PairGrid at 0x7ecbde3dfd70>



✓ Question

Is there a threshold in age gap size beyond which gender ordering becomes systematically asymmetric, indicating a shift in how relationships are portrayed?

✓ Feature Normalization

Because the ranges of the features differ substantially and the units are not consistent, feature normalization was applied to prevent any single feature from dominating the distance calculations

```
1 from sklearn.preprocessing import StandardScaler
2
3
4 feature_cols = ['mean_age', 'relative_age_gap']
5
6 X = age_gaps_copy[feature_cols]
7
8 sc = StandardScaler()
9 X_scaled = sc.fit_transform(X)
10
11 X_scaled_df = pd.DataFrame(
12     X_scaled,
13     columns=feature_cols,
14     index=age_gaps_copy.index
15 )
16
17 X_scaled_df.head()
```

	mean_age	relative_age_gap	
0	1.693628	5.471972	
1	1.693628	4.956601	
2	1.132271	6.020999	
3	1.257017	4.588479	
4	3.003460	2.193749	

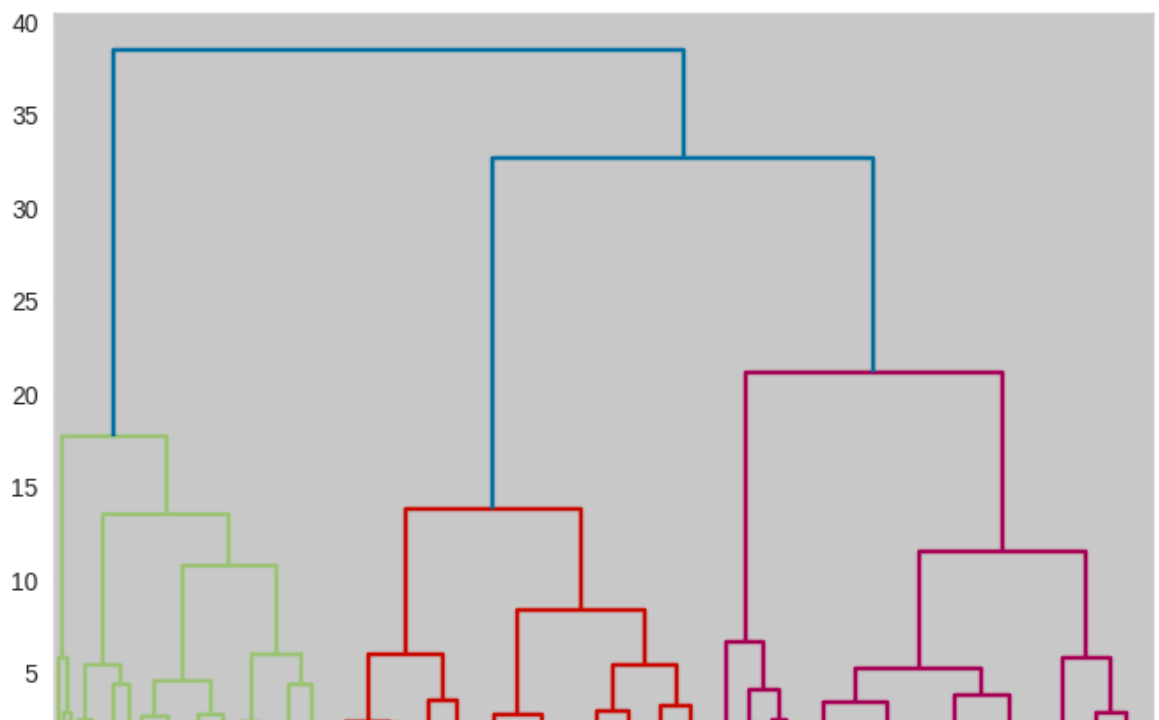
Next steps:

[Generate code with X_scaled_df](#)[New interactive sheet](#)

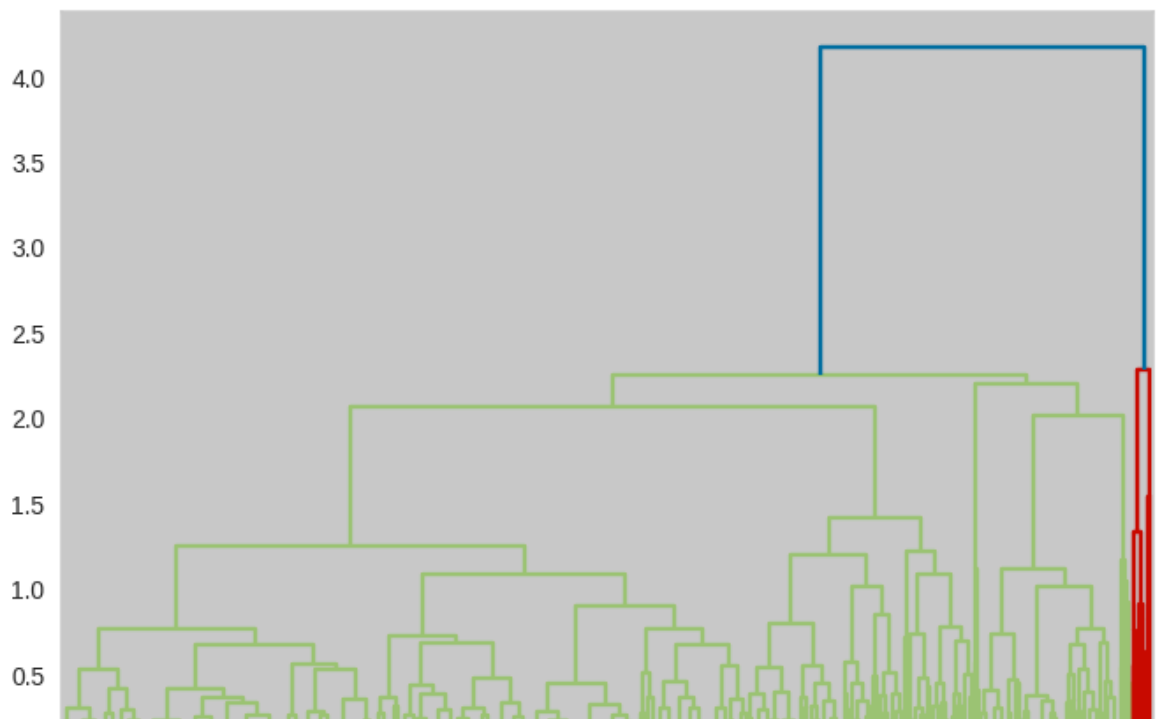
✓ Hierarchical Clustering

```
1 data_hiera = age_gaps_copy.copy()
```

```
1 from scipy.cluster.hierarchy import dendrogram
2 from scipy.cluster.hierarchy import linkage, fcluster
3
4 Z_ward = linkage(X_scaled, 'ward') #You may try different methods
5 tree_2 = dendrogram(Z_ward,)
```



```
1 Z_average = linkage(X_scaled, 'average') #You may try different m
2 tree_2 = dendrogram(Z_average,)
```



```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn.preprocessing import StandardScaler
5 from scipy.cluster.hierarchy import linkage, fcluster
6
7
8 df = age_gaps_copy.copy()
9 y_var = "age_difference"
10 feature_cols = ["mean_age", "relative_age_gap"]
11 n_clusters = 3
12 work_df = df.dropna(subset=feature_cols + [y_var]).copy()
13
14 scaler = StandardScaler()
15 X_scaled = scaler.fit_transform(work_df[feature_cols].to_numpy())
16
17 Z_ward = linkage(X_scaled, method="ward")
18 Z_avg = linkage(X_scaled, method="average")
19
20 work_df["cluster_ward"] = fcluster(Z_ward, t=n_clusters, criterion="maxclust")
21 work_df["cluster_avg"] = fcluster(Z_avg, t=n_clusters, criterion="maxclust")
22
23 fig, axes = plt.subplots(1, 2, figsize=(12, 5), sharey=True)
24
25 linkages = [
26     ("cluster_ward", "Ward linkage"),
27     ("cluster_avg", "Average linkage"),
28 ]
29
30 for ax, (col, title) in zip(axes, linkages):
31     sns.violinplot(
32         data=work_df,
33         x=col,
34         y=y_var,

```

```

35     inner="quartile",
36     cut=0,
37     scale="width",
38     color="#2c92b1",
39     ax=ax
40 )
41 ax.set_title(title)
42 ax.set_xlabel("Cluster")
43
44 axes[0].set_ylabel(y_var)
45
46 plt.suptitle(
47     f"{y_var} Distribution by Cluster Across Linkage Methods",
48     y=1.05
49 )
50 plt.tight_layout()
51 plt.show()

```

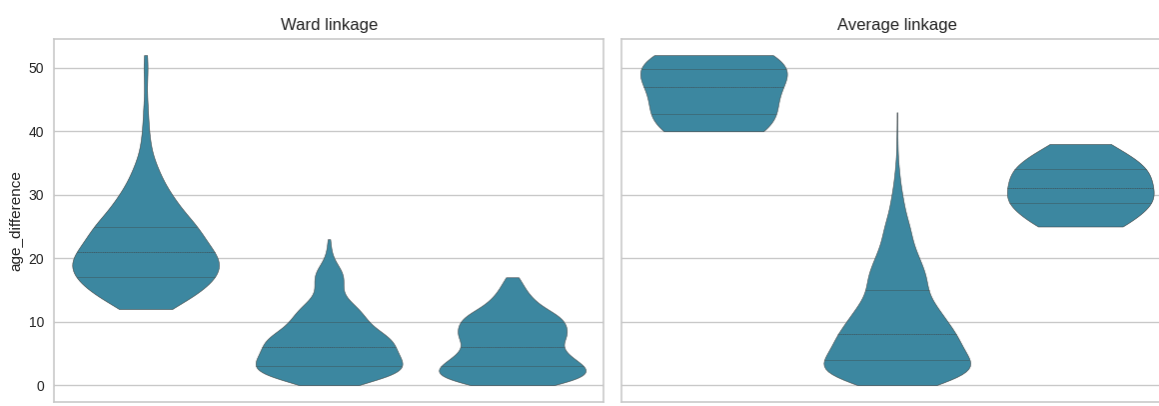
/tmp/ipython-input-3457659015.py:31: FutureWarning:

The `scale` parameter has been renamed and will be removed in v0.15.0.
sns.violinplot(

/tmp/ipython-input-3457659015.py:31: FutureWarning:

The `scale` parameter has been renamed and will be removed in v0.15.0.
sns.violinplot(

age_difference Distribution by Cluster Across Linkage Methods



Silhouette Score

Try different number of clusters to see which number of clusters has higher score.

I tried number of 2, 3 and 4. 3 has best performance.

```

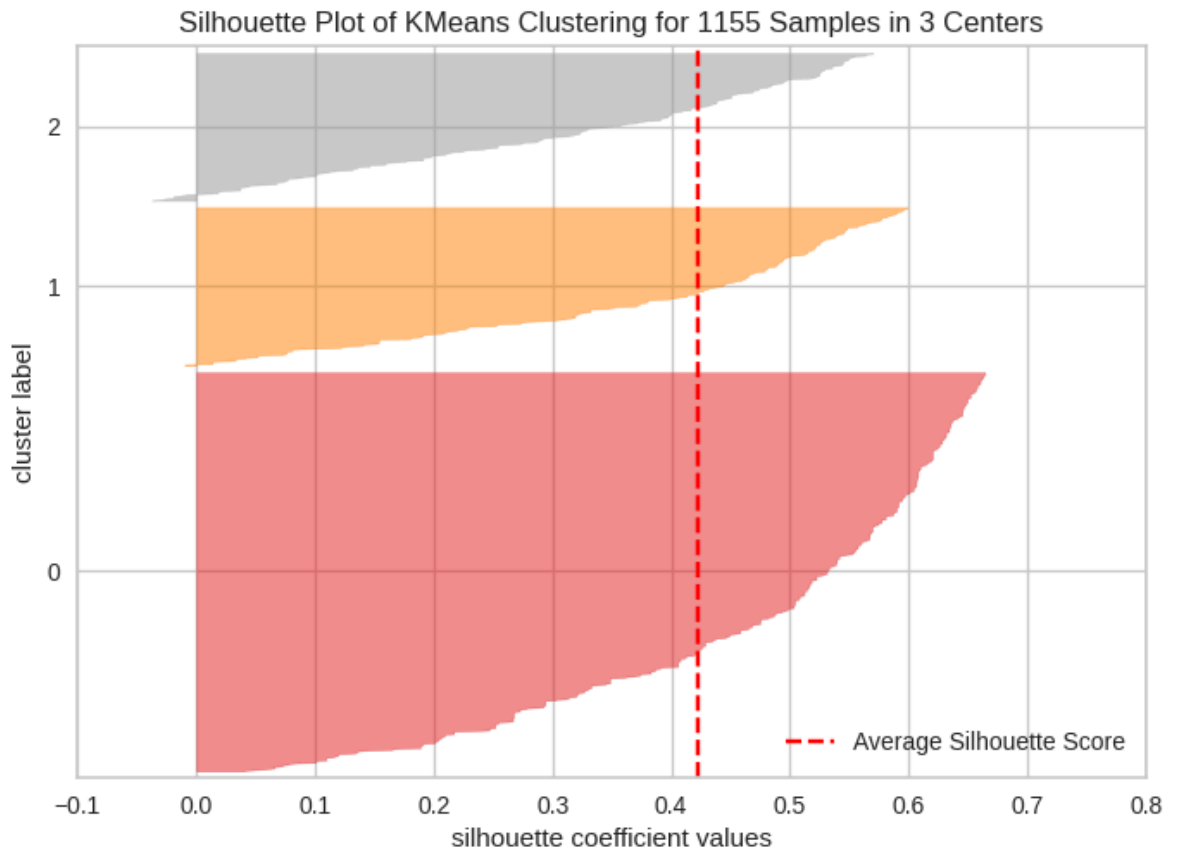
1 from sklearn.cluster import KMeans
2 from yellowbrick.cluster import KElbowVisualizer
3 from yellowbrick.cluster import SilhouetteVisualizer
4 n_clusters_kmeans = 3
5
6 kmeans_model = KMeans(n_clusters_kmeans, random_state=42)
7 visualizer = SilhouetteVisualizer(kmeans_model, timings=False)
8

```

```

9 visualizer.fit(X_scaled)          # Fit the data to the visualizer
10 visualizer.show()

```

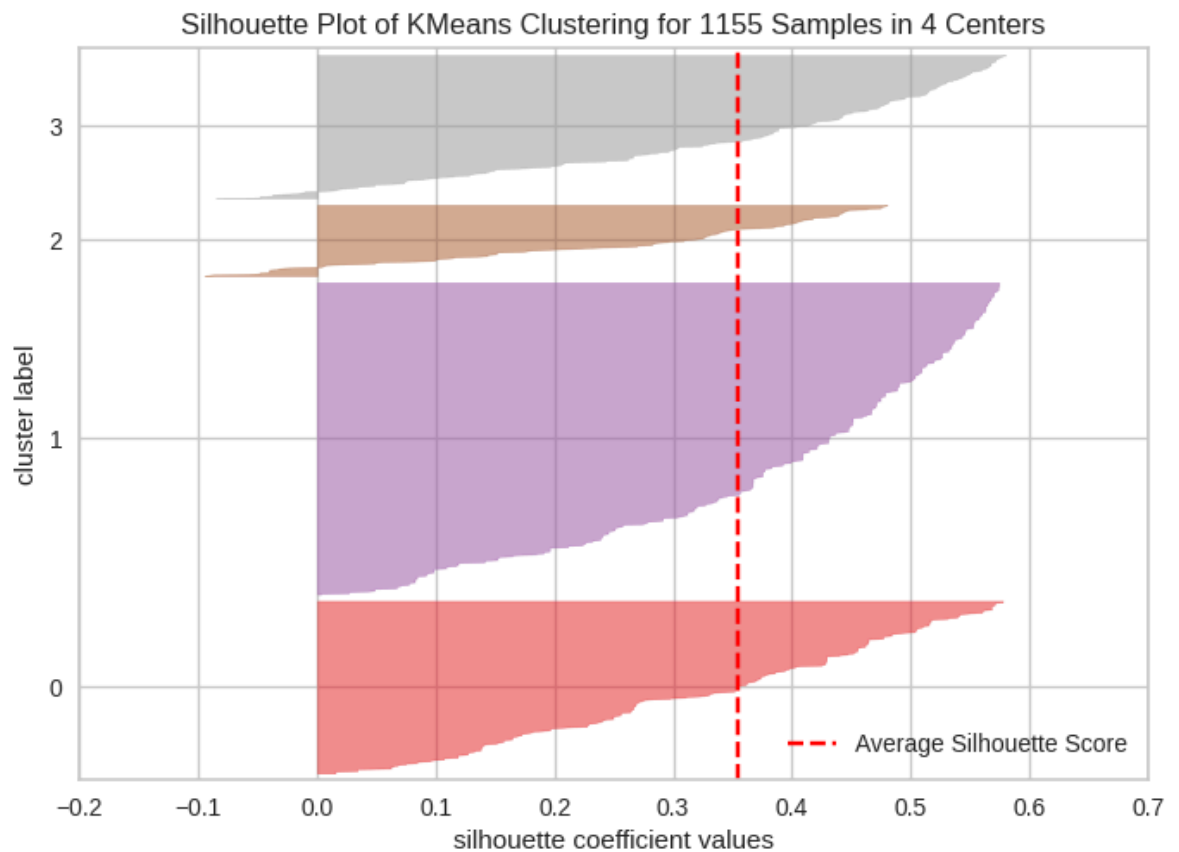


<Axes: title={'center': 'Silhouette Plot of KMeans Clustering for 1155

```

1 from sklearn.cluster import KMeans
2 from yellowbrick.cluster import KElbowVisualizer
3 from yellowbrick.cluster import SilhouetteVisualizer
4 n_clusters_kmeans = 4
5
6 kmeans_model = KMeans(n_clusters_kmeans, random_state=42)
7 visualizer = SilhouetteVisualizer(kmeans_model, timings=False)
8
9 visualizer.fit(X_scaled)          # Fit the data to the visualizer
10 visualizer.show()

```

<Axes: title={'center': 'Silhouette Plot of KMeans Clustering for 1155

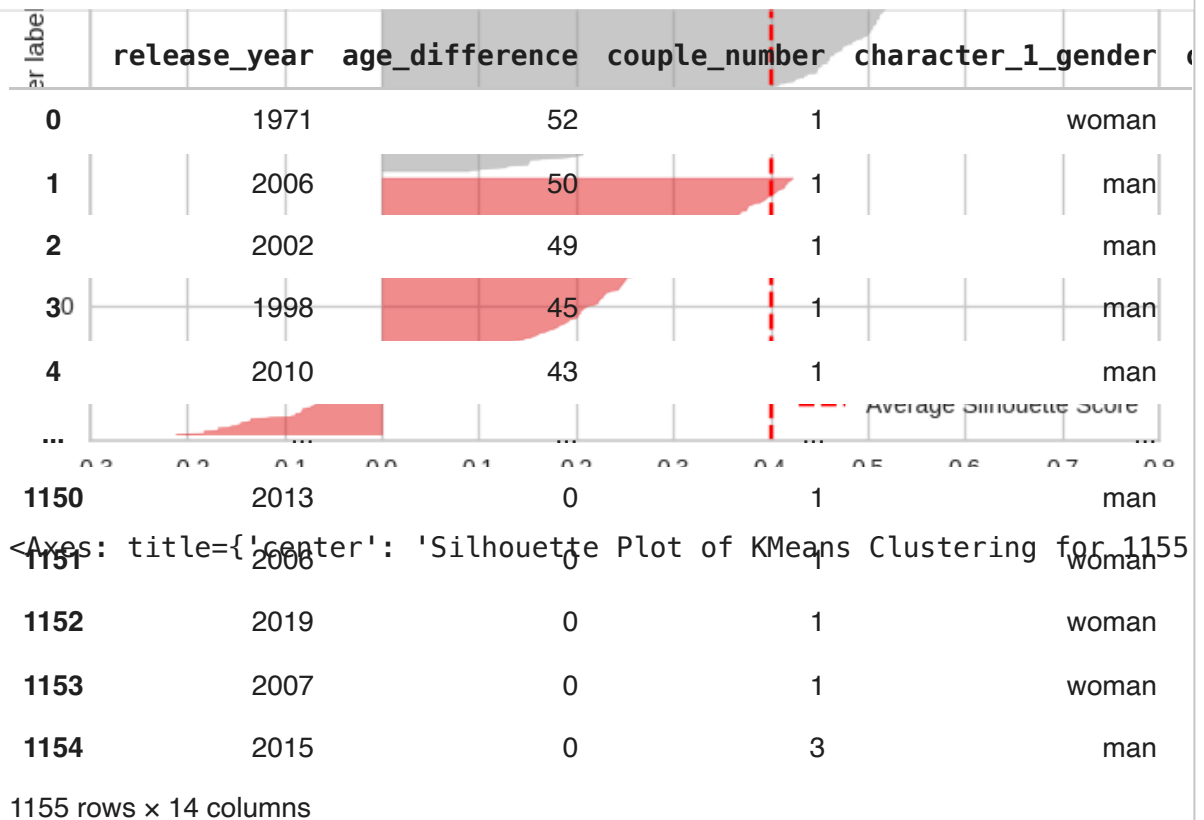
```

1 from sklearn.cluster import KMeans
2 from yellowbrick.cluster import KElbowVisualizer
3 from yellowbrick.cluster import SilhouetteVisualizer
4 n_clusters_kmeans = 2
5
6 kmeans_model = KMeans(n_clusters_kmeans, random_state=4
7 visualizer = SilhouetteVisualizer(kmeans_model, timings
8
9 visualizer.fit(X_scaled)           # Fit the data to the v
10 visualizer.show()
```

Silhouette Plot of KMeans Clustering for 1155 Samples in 2 Centers



```
1 n_clusters = 3
2 data_hiera['cluster_labels'] = fcluster(Z_average, n_cl
3 data_hiera
```



Next steps:

[Generate code with data_hiera](#)[New interactive sheet](#)

```
1 data_hiera['cluster_labels'].value_counts()
```

count	
cluster_labels	
1	1133
2	16
3	6

dtype: int64

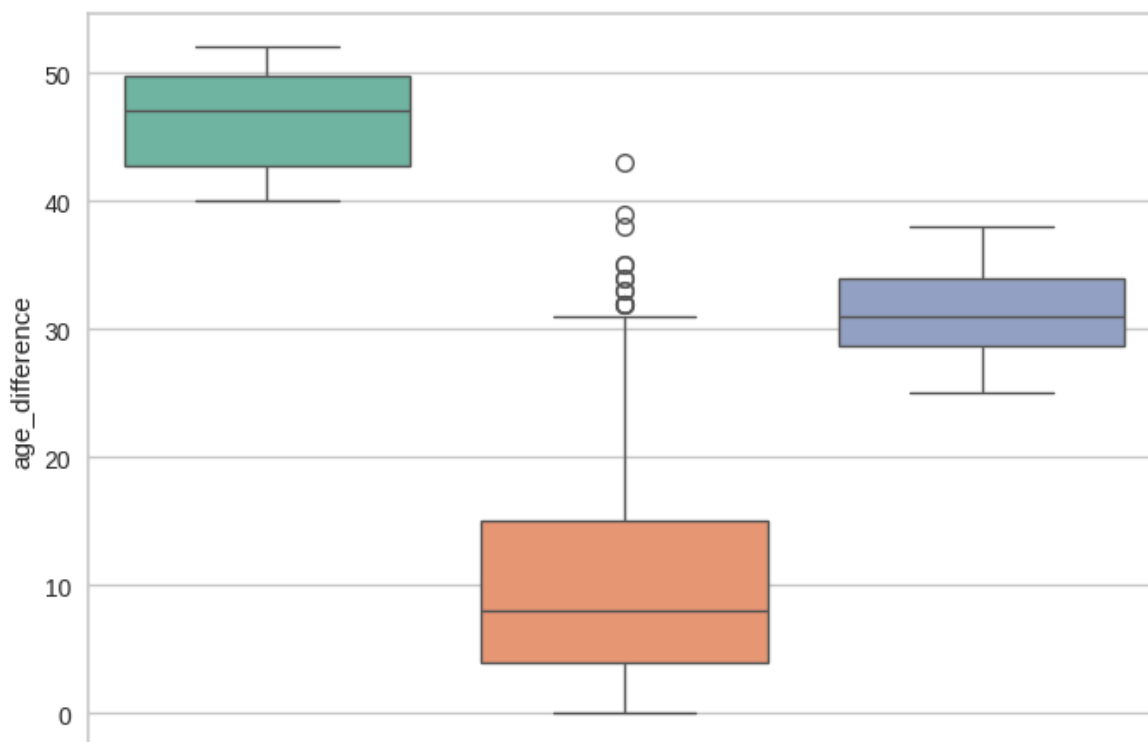
```
1 data_hiera.groupby('cluster_labels')['age_difference'].describe()
```

	count	mean	std	min	25%	50%	75%	max	
cluster_labels									
1	1133.0	9.943513	7.777061	0.0	4.00	8.0	15.00	43.0	
2	16.0	31.000000	4.016632	25.0	28.75	31.0	34.00	38.0	
3	6.0	46.333333	4.760952	40.0	42.75	47.0	49.75	52.0	

```
1 sns.boxplot(data=data_hiera, x='cluster_labels', y="age_difference")
```

/tmp/ipython-input-375687973.py:1: FutureWarning:
Passing `palette` without assigning `hue` is deprecated and will be removed in a future version. Assigning `hue` to the `cluster_labels` variable will silence this warning.

```
sns.boxplot(data=data_hiera, x='cluster_labels', y="age_difference",  
<Axes: xlabel='cluster_labels', ylabel='age_difference'>
```

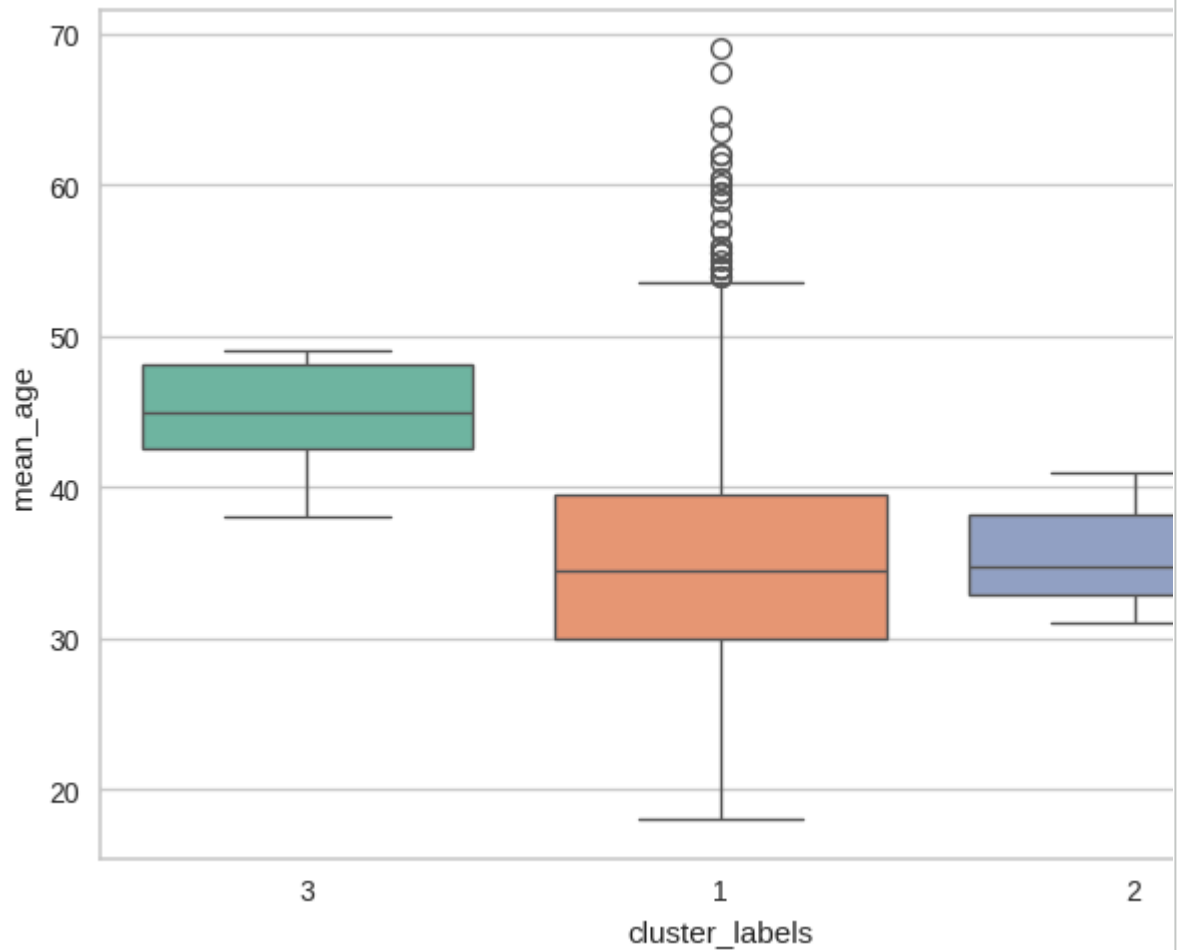


```
1 sns.boxplot(data=data_hiera, x='cluster_labels', y="mean_age", pa
```

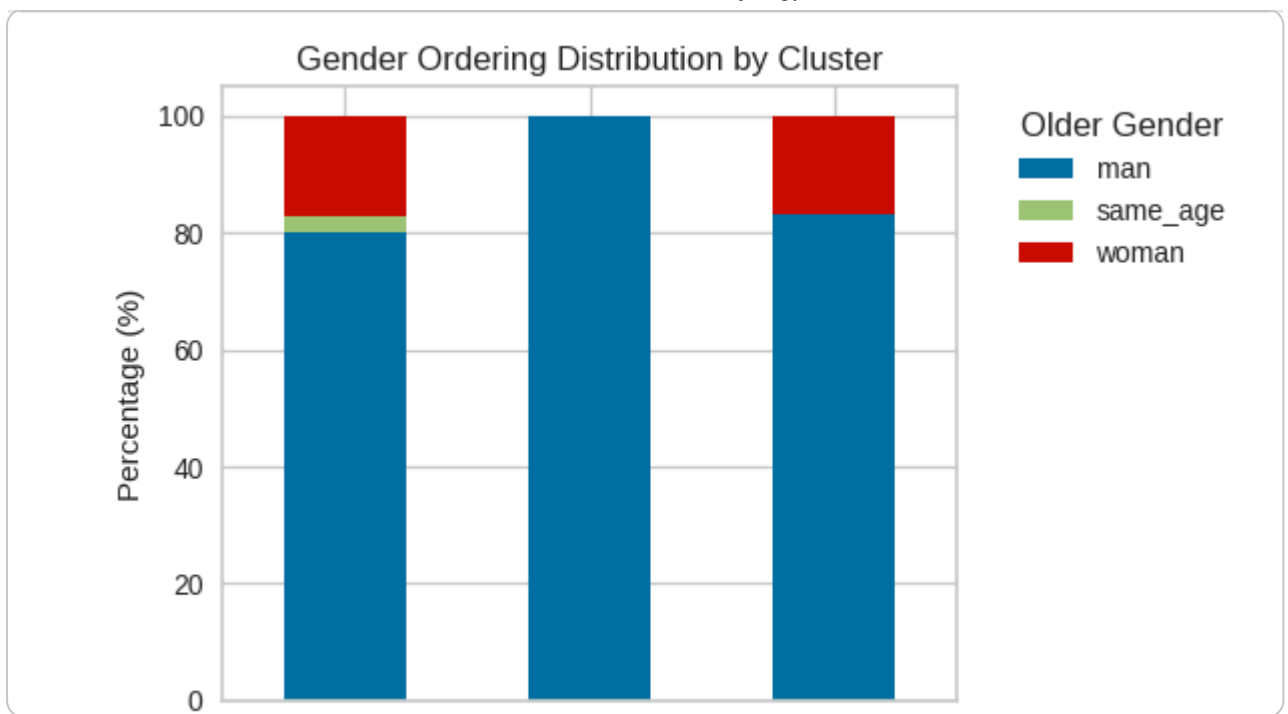
/tmp/ipython-input-924421717.py:1: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be re

```
sns.boxplot(data=data_hiera, x='cluster_labels', y="mean_age", palet
<Axes: xlabel='cluster_labels', ylabel='mean_age'>
```



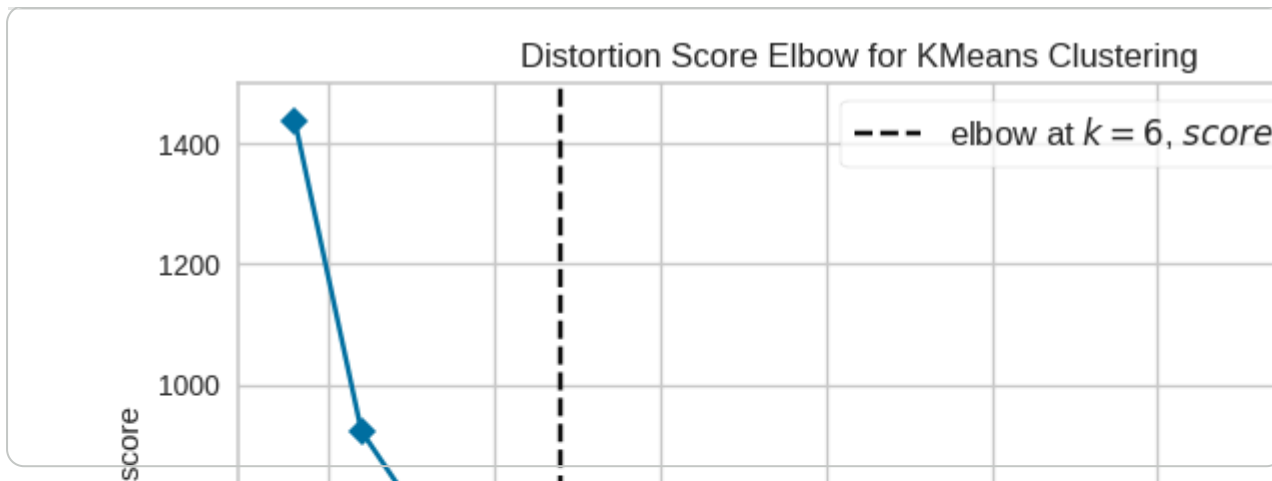
```
1 import matplotlib.pyplot as plt
2
3 gender_pct = (
4     pd.crosstab(data_hiera['cluster_labels'], data_hiera['older_g
5     * 100
6 )
7
8 gender_pct.plot(
9     kind='bar',
10    stacked=True,
11    figsize=(6, 4)
12 )
13
14 plt.ylabel('Percentage (%)')
15 plt.xlabel('Cluster')
16 plt.title('Gender Ordering Distribution by Cluster')
17 plt.legend(title='Older Gender', bbox_to_anchor=(1.05, 1), loc='u
18 plt.tight_layout()
19 plt.show()
```



✓ K-Mean

```
1 X_scaled_k=X_scaled_df.drop(columns=['cluster_labels'])
```

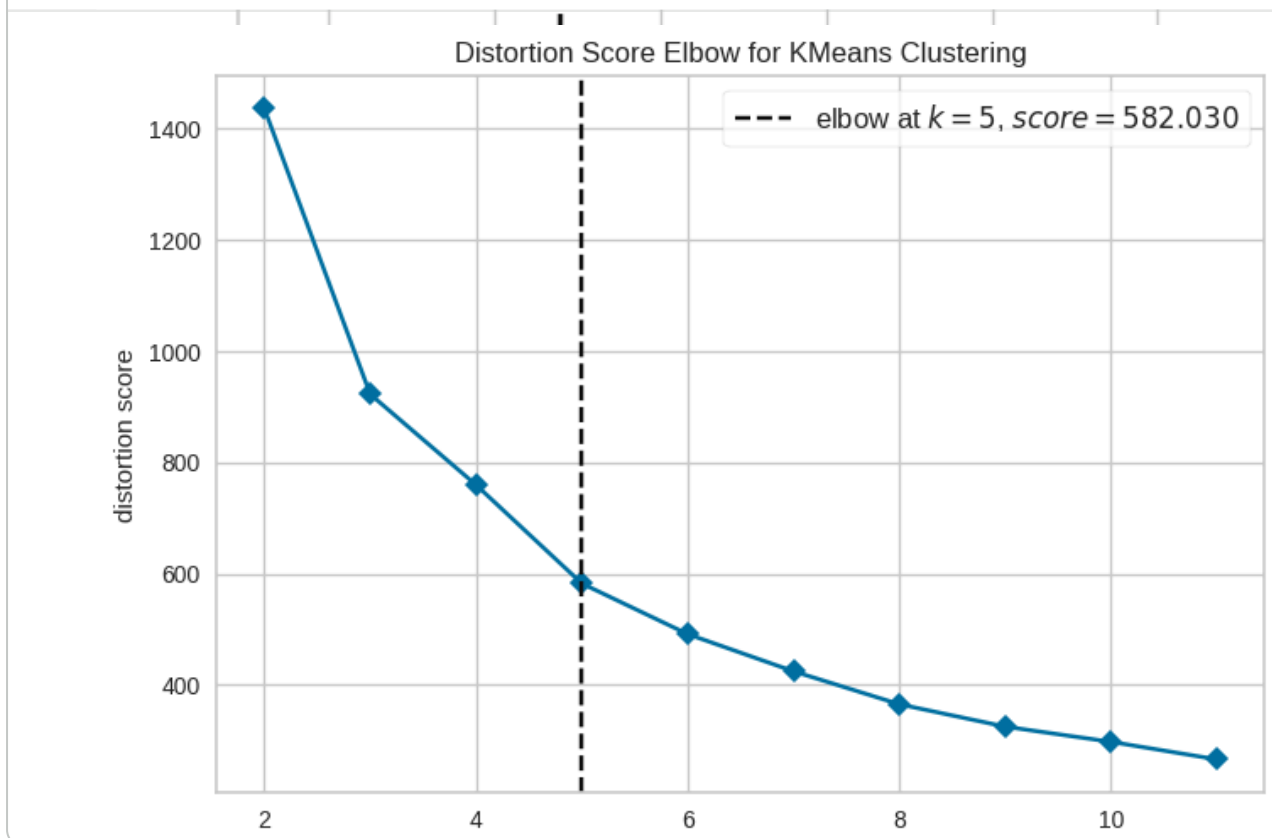
```
1 from sklearn.cluster import KMeans
2 from yellowbrick.cluster import KElbowVisualizer
3
4 clustering = KMeans(random_state=42)
5 visualizer = KElbowVisualizer(
6     clustering, k=(2, 20), metric='distortion', timings
7 )
8
9 visualizer.fit(X_scaled_k)          # Fit the data to the
10 visualizer.show()
```



```

1 from sklearn.cluster import KMeans
2 from yellowbrick.cluster import KElbowVisualizer
3
4 clustering = KMeans(random_state=42)
5 visualizer = KElbowVisualizer(
6     clustering, k=(2, 12), metric='distortion', timings=False, ra
7 )
8
9 visualizer.fit(X_scaled_k)          # Fit the data to the visualize
10 visualizer.show()

```



```
1 n_clusters_kmeans = 4
```

```

1 data_k=age_gaps_copy.copy()
2 data_k=data_k.drop(columns=['character_1_gender', 'character_2_

```

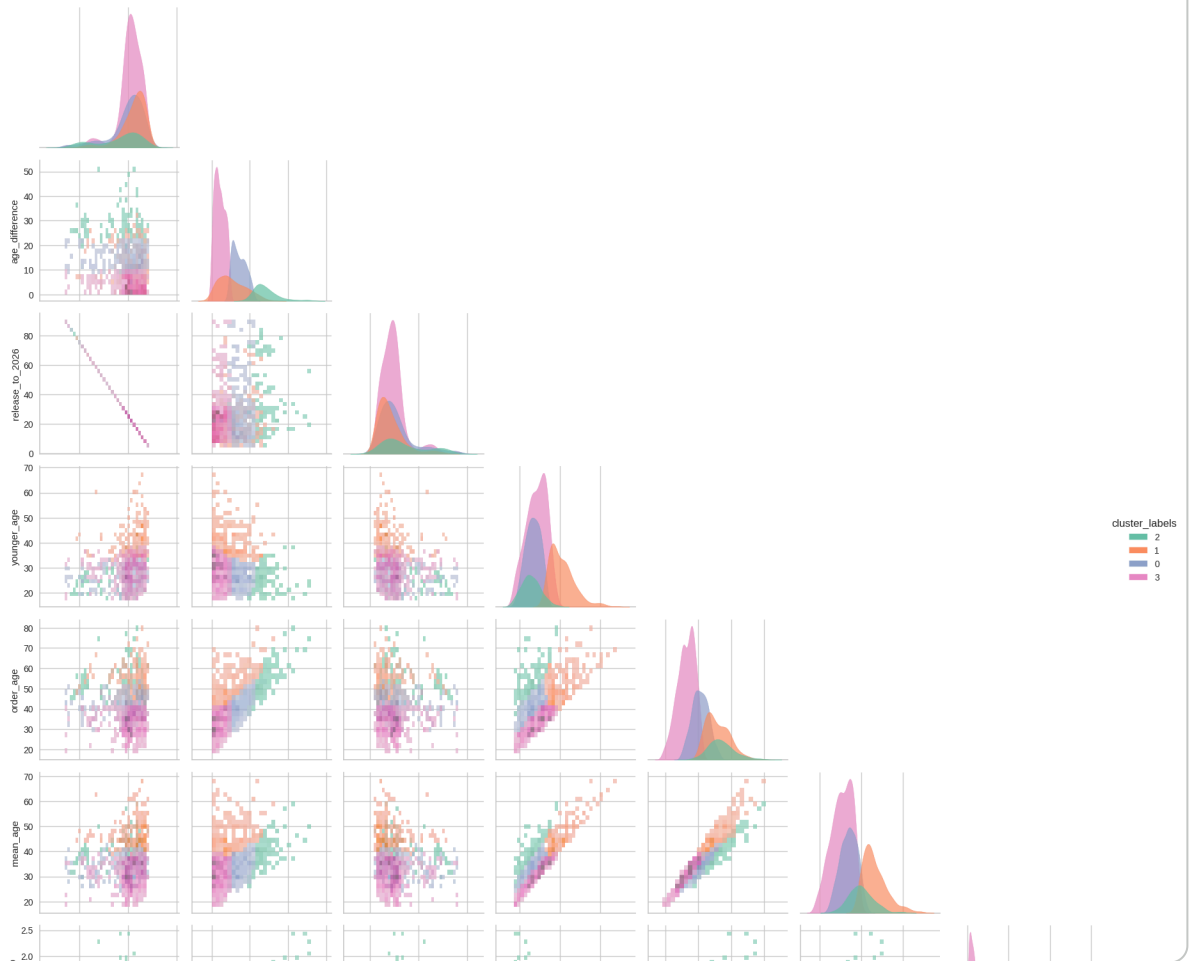
```
1 from sklearn.cluster import KMeans
2
3 # Initialize KMeans model with n_clusters and random_state=42
4 kmeans = KMeans(n_clusters=n_clusters_kmeans, random_state=42) #
5
6 # Fit the model and predict cluster labels, excluding the existing
7 data_k['cluster_labels'] = kmeans.fit_predict(X_scaled_df).astype
8
9 # Display the value counts of the new cluster labels
10 data_k['cluster_labels'].value_counts().sort_index()
```

	count
cluster_labels	
0	279
1	238
2	112
3	526

dtype: int64

```
1 sns.pairplot(data_k, hue="cluster_labels", kind="hist", diag_kind:
2     plot_kws={'alpha': 0.7}, # Adjusts transparency of the scatter
3     diag_kws={'alpha': 0.7} # Adjusts transparency of the diagonal
4     )
```

<seaborn.axisgrid.PairGrid at 0x7ecbdb578620>



1 data_k

	release_year	age_difference	release_to_2026	younger_age	order_
0	1971	52	55	23	
1	2006	50	20	24	
2	2002	49	24	20	
3	1998	45	28	23	
4	2010	43	16	38	
...	
1150	2013	0	13	23	
1151	2006	0	20	42	
1152	2019	0	7	30	
1153	2007	0	19	21	
1154	2015	0	11	36	

1155 rows × 9 columns

Next steps:

[Generate code with data_k](#)

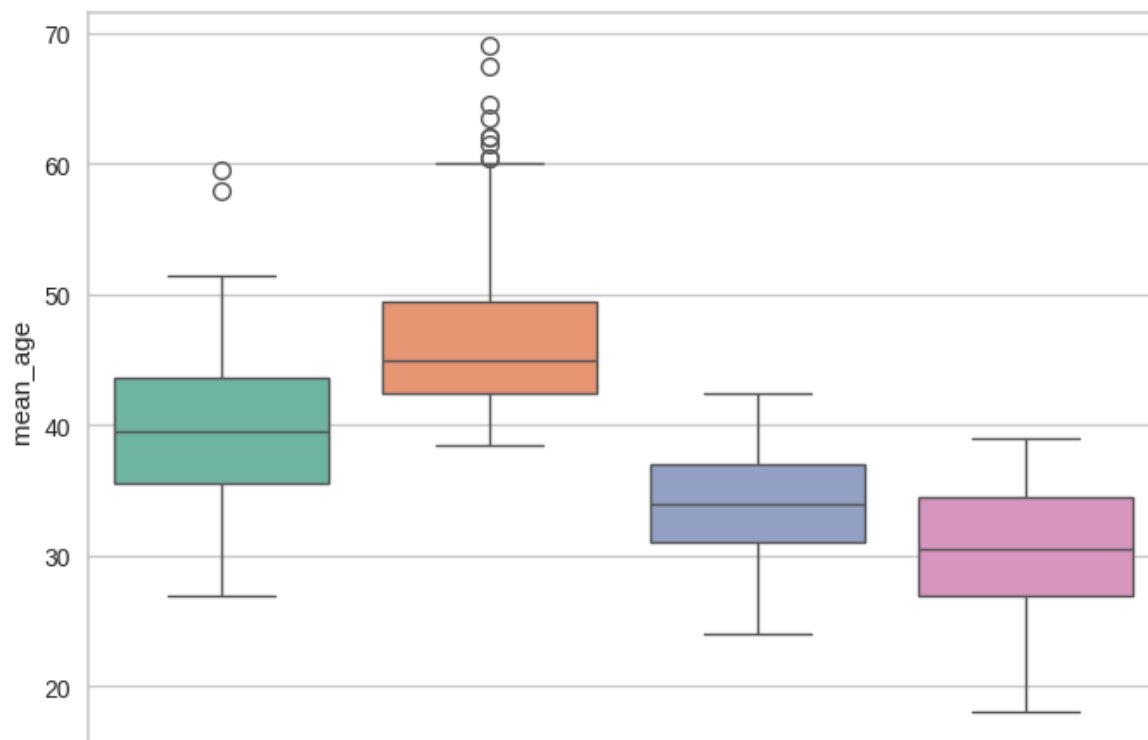
[New interactive sheet](#)


```
1 sns.boxplot(data=data_k, x='cluster_labels', y="mean_age", palette="
```

```
/tmp/ipython-input-1587560851.py:1: FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be re
```

```
sns.boxplot(data=data_k, x='cluster_labels', y="mean_age", palette="")
<Axes: xlabel='cluster_labels', ylabel='mean_age'>
```



```
1 sns.boxplot(data=data_k, x='cluster_labels', y="relative_age_gap"
```

/tmp/ipython-input-988000534.py:1: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be re

```
1 sns.boxplot(data=data_k, x='cluster_labels', y="order_age", palet
```

<Axes: xlabel='cluster_labels', ylabel='relative_age_gap'>

/tmp/ipython-input-21855133.py:1: FutureWarning:

2.5
Passing `palette` without assigning `hue` is deprecated and will be re

```
sns.boxplot(data=data_k, x='cluster_labels', y="order_age", palette=
```

<Axes: xlabel='cluster_labels', ylabel='order_age'>

