

Lab Exercise: Visual Exploration of Shape Features (Self-Contained)

Course: Data Visualization & Exploratory Analysis

Duration: 2 hours (lab) + 10 minutes (presentation)

Abstract

This lab is self-contained: students receive the PDF only. The PDF contains a data-generation script, full runnable code snippets to produce each required figure (histograms, heatmaps, rankings, trends), and step-by-step instructions so no external notebooks or data are required.

Prerequisites

Python (3.8+), pandas, seaborn, matplotlib, plotly, scikit-learn, ipywidgets

What to expect in this handout

- A small synthetic dataset generator (Python) to create `feature_vectors/*.csv` locally.
- Ready-to-run plotting scripts for: histograms, interactive comparison, correlation heatmaps, pairplots, clustered heatmap, ranking charts (including diverging and stacked), and trend analyses.
- AI prompts (for instructors) used to create the synthetic data.
- Deliverables, tasks, and graded questions.

1 Quick setup (5 minutes)

1. Create and activate your Python environment, then install required packages:

```
1 python -m pip install --upgrade pip
2 pip install pandas seaborn matplotlib plotly scikit-learn ipywidgets
```

2. Save the data-generator script below as `generate_synthetic_feature_vectors.py` and run it from the project root:

```
1 python generate_synthetic_feature_vectors.py
```

3. This will create a `feature_vectors/` folder with 5 CSV files (club, dining, lamp, mug, table). All remaining code snippets assume that folder exists in the current working directory.

2 Data generation script (save as `generate_synthetic_feature_vectors.py`)

```
1 """
2 Creates synthetic feature CSVs for five categories: club, dining, lamp, mug, table.
3 Each CSV contains numeric features used in this lab and two metadata columns: source_file
, category.
```

```

4 """
5 import numpy as np
6 import pandas as pd
7 from pathlib import Path
8
9 np.random.seed(0)
10 out = Path('feature_vectors')
11 out.mkdir(exist_ok=True)
12
13 def make_cat(cat, n=200):
14     # base ranges altered per category to create separability
15     if cat == 'table':
16         sa = np.random.normal(500, 80, n).clip(50, 2000)
17         vol = np.random.normal(200, 40, n).clip(5, 1000)
18     elif cat == 'mug':
19         sa = np.random.normal(120, 30, n).clip(10, 800)
20         vol = np.random.normal(45, 15, n).clip(1, 400)
21     elif cat == 'lamp':
22         sa = np.random.normal(180, 60, n).clip(10, 1200)
23         vol = np.random.normal(60, 25, n).clip(1, 500)
24     elif cat == 'club':
25         sa = np.random.normal(90, 25, n).clip(10, 600)
26         vol = np.random.normal(35, 12, n).clip(1, 300)
27     else: # dining
28         sa = np.random.normal(300, 70, n).clip(20, 1500)
29         vol = np.random.normal(120, 30, n).clip(1, 800)
30
31 mean_curv = np.random.normal(25, 7, n).clip(0, 100) * (1 + (np.log1p(sa)/10 - 1))
32 median_curv = mean_curv * (0.9 + 0.2*np.random.rand(n))
33 silhouette_complexity = np.random.normal(10, 5, n).clip(0, 50)
34 skeleton_complexity = np.random.normal(8, 4, n).clip(0, 50)
35 aspect_ratio_y = np.random.normal(1.0 + (vol/500), 0.15, n).clip(0.2, 10)
36 hollow_ratio = np.random.beta(2,5, n) * (1 + sa/2000)
37 surface_to_volume_ratio = sa / (vol + 1e-6)
38
39 df = pd.DataFrame({
40     'mean_curvature': mean_curv,
41     'median_curvature': median_curv,
42     'surface_area': sa,
43     'volume': vol,
44     'silhouette_complexity': silhouette_complexity,
45     'skeleton_complexity': skeleton_complexity,
46     'aspect_ratio_y': aspect_ratio_y,
47     'hollow_ratio': hollow_ratio,
48     'surface_to_volume_ratio': surface_to_volume_ratio,
49 })
50 df['category'] = cat
51 df['source_file'] = [f"{cat}_{i:04d}.obj" for i in range(len(df))]
52 return df
53
54 cats = ['club', 'dining', 'lamp', 'mug', 'table']
55 for c in cats:
56     df = make_cat(c, n=200)
57     df.to_csv(out / f"{c}_features.csv", index=False)
58 print('Synthetic data written to', out)

```

How to run the analysis

Copy-paste these snippets into a .py file or run in Jupyter:

3 SECTION 1 — Histograms and basic stats

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 from pathlib import Path
5
6 root = Path('.')
7 fv_dir = root / 'feature_vectors'
8 csvs = sorted(fv_dir.glob('*_features.csv'))
9 dfs = [pd.read_csv(p).assign(category=p.stem.split('_')[0]) for p in csvs]
10 all_df = pd.concat(dfs, ignore_index=True)
11
12 def find_col(df, target):
13     t = target.lower()
14     for c in df.columns:
15         if t in c.lower().replace('_', ' '):
16             return c
17     return None
18
19 col = find_col(all_df, 'Surface-to-Volume Ratio')
20 vals = pd.to_numeric(all_df[all_df['category']=='mug'][col], errors='coerce').dropna()
21 plt.figure(figsize=(8,4))
22 sns.histplot(vals, bins=30, kde=True)
23 plt.title('Surface-to-Volume Ratio - mugs')
24 plt.show()
25 print('n =', len(vals))
26 print('mean =', vals.mean(), 'std =', vals.std())
```

4 SECTION 2 — Interactive comparison (overlay or violin/box)

```
1 # Overlay two categories
2 c1, c2 = 'mug', 'table'
3 f = 'surface_to_volume_ratio'
4 v1 = pd.to_numeric(all_df[all_df['category']==c1][f], errors='coerce').dropna()
5 v2 = pd.to_numeric(all_df[all_df['category']==c2][f], errors='coerce').dropna()
6 plt.figure(figsize=(8,4))
7 sns.histplot(v1, color='tab:blue', alpha=0.5, bins=30, label=c1)
8 sns.histplot(v2, color='tab:orange', alpha=0.5, bins=30, label=c2)
9 plt.legend()
10 plt.title(f'Overlay: {f} - {c1} vs {c2}')
11 plt.show()
12
13 # Normalized (z-score)
14 import numpy as np
15 combined = pd.concat([v1, v2])
16 mu, sigma = combined.mean(), combined.std()
17 v1n = (v1 - mu)/sigma
18 v2n = (v2 - mu)/sigma
19 plt.figure(figsize=(8,4))
20 sns.histplot(v1n, color='tab:blue', alpha=0.5, bins=30, label=c1)
21 sns.histplot(v2n, color='tab:orange', alpha=0.5, bins=30, label=c2)
22 plt.legend()
23 plt.title('Normalized overlay (z-score)')
24 plt.show()
```

5 SECTION 3 — Heatmaps and pairplot

```
1 import numpy as np
2 num_cols = all_df.select_dtypes(include='number').columns.tolist()
3 cols = num_cols[:20]
```

```

4 corr = all_df[cols].corr()
5 plt.figure(figsize=(10,8))
6 mask = np.triu(np.ones_like(corr, dtype=bool))
7 cmap = sns.diverging_palette(220, 10, as_cmap=True)
8 sns.heatmap(corr, mask=mask, cmap=cmap, annot=True, fmt=' .2f ')
9 plt.title('Correlation heatmap')
10 plt.show()
11
12 # Pairplot (choose -46 features)
13 top = all_df[cols].var().sort_values(ascending=False).index[:6].tolist()
14 sns.pairplot(all_df[top + ['category']].dropna(), hue='category', diag_kind='kde', corner
   =True)
15 plt.show()

```

6 SECTION 4—Rankings (bar charts, error bars, stacked, diverging, lollipop)

```

1 rank_df = pd.DataFrame({'feature': cols, 'importance': np.linspace(1, len(cols), len(cols
   ))[::-1]})
2 # classic horizontal bar
3 top = rank_df.sort_values('importance', ascending=True).tail(15)
4 plt.figure(figsize=(8,6))
5 plt.barh(top['feature'], top['importance'], color='tab:blue')
6 plt.title('Top features – horizontal ranking (synthetic)')
7 plt.show()
8
9 # error bars using means/std across shapes
10 top_feats = rank_df['feature'].tolist()[:10]
11 means = all_df[top_feats].mean()
12 stds = all_df[top_feats].std()
13 order = means.sort_values(ascending=True).index
14 plt.figure(figsize=(8,6))
15 plt.barh(order, means[order], xerr=stds[order], color='tab:green')
16 plt.title('Means ± std (top features)')
17 plt.show()
18
19 # stacked bar: mean per category for a small set
20 feats = top_feats[:6]
21 pivot = all_df.groupby('category')[feats].mean()
22 pivot.plot(kind='bar', stacked=True)
23 plt.title('Stacked: mean per category')
24 plt.show()
25
26 # diverging: difference between two categories
27 c1, c2 = sorted(all_df['category'].unique())[:2]
28 mean1 = all_df[all_df['category']==c1][feats].mean()
29 mean2 = all_df[all_df['category']==c2][feats].mean()
30 diff = (mean1 - mean2).sort_values()
31 plt.figure(figsize=(8,4))
32 colors = ['tab:blue' if v>0 else 'tab:orange' for v in diff.values]
33 plt.barh(diff.index, diff.values, color=colors)
34 plt.axvline(0, color='k', linewidth=0.6)
35 plt.title(f'Diverging: {c1} - {c2}')
36 plt.show()
37
38 # lollipop
39 dfp = rank_df.sort_values('importance', ascending=True).tail(15)
40 plt.figure(figsize=(8,6))
41 plt.hlines(y=dfp['feature'], xmin=0, xmax=dfp['importance'], color='gray')
42 plt.plot(dfp['importance'], dfp['feature'], 'o', color='tab:red')
43 plt.title('Lollipop chart')
44 plt.show()

```

7 SECTION 5 — Trends (rolling mean, binned, quantiles)

```
1 from sklearn.decomposition import PCA
2 from sklearn.preprocessing import StandardScaler
3
4 # compute a proxy target (if no real score available)
5 num = all_df.select_dtypes(include='number').drop(columns=['source_file'], errors='ignore')
6 drop_ids = [c for c in num.columns if 'id' in c.lower() or c.lower().startswith('shape')]
7 num = num.drop(columns=drop_ids, errors='ignore')
8 X = StandardScaler().fit_transform(num.fillna(num.median()))
9 pc1 = PCA(n_components=1, random_state=0).fit_transform(X).ravel()
10 all_df['aesthetic_proxy'] = pc1
11
12 # rolling trend example
13 f = 'mean_curvature'
14 target = 'aesthetic_proxy'
15 df = all_df[[f, target]].dropna().sort_values(f)
16 df['y_roll'] = df[target].rolling(window=max(5, int(len(df)*0.12)), center=True).mean()
17 plt.figure(figsize=(8,4))
18 plt.scatter(df[f], df[target], s=10, alpha=0.3)
19 plt.plot(df[f], df['y_roll'], color='tab:red')
20 plt.title('Scatter + rolling mean')
21 plt.show()
```

8 Deliverables (submit PDF)

- A short summary (1 paragraph) of the lab and methods used.
- Three figures (a histogram, a heatmap/pairplot, and a ranking chart) with concise captions and interpretation (2-3 sentences each).
- Answer: "Which single feature best separates mugs from tables?" Provide a figure and justification.

9 Tips and grading

- Figures: 40%
- Written answers: 40%
- Clean runnable scripts/notebook: 15%
- Extra credit: 5%

10 Academic integrity

Work in pairs at most. Cite any AI or external code used.