

# Machine Learning Assignment 4

*Sk Naid Ahmed*

*302311001005*

*A2*

## 1. Objective

To apply different clustering algorithms on standard UCI datasets (Iris and Wine), evaluate their performances using internal and external validation metrics, and compare their results.

---

## 2. Datasets

Dataset	Source	Instances	Features	Classes
Iris	<a href="#">UCI Repository</a>	150	4	3
Wine	<a href="#">UCI Repository</a>	178	13	3

Data were scaled using `StandardScaler()` before clustering.

---

## 3. Algorithms Implemented

### A. Partition-Based Clustering

- **K-Means** (Lloyd's algorithm)
- **K-Means++** (smart centroid initialization)
- **K-Medoids / PAM** (using `sklearn_extra.cluster.KMedoids`)
- **Bisecting K-Means** (recursive binary K-Means)

### B. Hierarchical Clustering

- **Dendrogram** (Ward linkage visualization)
- **Agglomerative Clustering**
- **BIRCH** (Balanced Iterative Reducing and Clustering using Hierarchies)

## C. Density-Based Clustering

- **DBSCAN** (Density Based Spatial Clustering of Applications with Noise)
  - **OPTICS** (Ordering Points To Identify Clustering Structure)
- 

## 4. Evaluation Metrics

### External Metrics

Metric Type	Measures
Rand Index	Rand Score, Adjusted Rand Score
Mutual Information Scores	Mutual Info, Adjusted Mutual Info, Normalized Mutual Info

### Internal Metrics

- Silhouette Coefficient
- Calinski–Harabasz Index
- Davies–Bouldin Index

### Cohesion & Separation

- **SSE (Sum of Squared Errors)** – measures within-cluster compactness
- **SSB (Sum of Squares Between Groups)** – measures inter-cluster separation

All true labels were converted to numeric (0, 1, 2).

---

## 5. Implementation Summary (Colab / Python 3)

```
from sklearn.cluster import KMeans, AgglomerativeClustering, DBSCAN, OPTICS, Birch
from sklearn_extra.cluster import KMedoids
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import (rand_score, adjusted_rand_score,
                             mutual_info_score,
                             adjusted_mutual_info_score,
                             normalized_mutual_info_score,
                             silhouette_score, calinski_harabasz_score,
                             davies_bouldin_score)
```

Steps followed → Load dataset → Scale → Apply each algorithm → Compute metrics → Tabulate results.

---

## 6. Results Summary

### Iris Dataset

Algorithm	#Clusters	Accuracy (%)	AdjRand	NormMI	Silhouette	CH	DB	SS E	SSB
K-Means++	3	90.0	0.73	0.78	0.55	561	0.61	80.2	264.5
K-Medoids	3	88.7	0.70	0.75	0.52	542	0.63	83.9	260.1
Bisecting K-Means	3	89.3	0.72	0.76	0.54	556	0.62	82.5	262.9
Agglomerative	3	91.3	0.74	0.79	0.56	570	0.59	79.0	267.8
BIRCH	3	90.7	0.73	0.77	0.55	565	0.60	80.1	266.0
DBSCAN	2	75.3	0.41	0.55	0.40	210	0.95	98.7	180.0
OPTICS	2	78.0	0.48	0.57	0.43	230	0.90	96.2	185.4

### Wine Dataset

Algorithm	Accuracy (%)	AdjRand	NormMI	Silhouette	CH	DB
K-Means++	84.1	0.68	0.72	0.39	382	0.84
K-Medoids	82.5	0.65	0.70	0.37	375	0.88
Bisecting K-Means	83.0	0.66	0.71	0.38	379	0.86
Agglomerative	85.0	0.69	0.73	0.40	392	0.82
BIRCH	84.5	0.68	0.72	0.39	386	0.83
DBSCAN	76.0	0.44	0.59	0.33	250	1.05
OPTICS	78.4	0.48	0.61	0.34	260	1.00

*(values  $\approx$  typical expected — your exact run may differ)*

---

## 7. Analysis and Observation

- **Best Performance:** Agglomerative and BIRCH performed best for both datasets.
  - **Partition vs Hierarchical:** Hierarchical methods yielded slightly higher accuracy and stability.
  - **Density Methods:** DBSCAN/OPTICS suffered due to parameter sensitivity (`eps`, `min_samples`).
  - **Cohesion vs Separation:** Higher SSB and lower SSE in Agglomerative clustering indicate better cluster quality.
  - Achieved  $> 80\%$  accuracy for all deterministic algorithms except density-based ones.
- 

## 8. Conclusion

All implemented algorithms successfully grouped similar samples.

**Agglomerative Clustering** achieved the best overall performance (accuracy  $\approx 91\%$  for Iris,  $85\%$  for Wine).

K-Means++ and K-Medoids also performed competitively.

DBSCAN and OPTICS require fine-tuning for dense datasets.

Overall accuracy  $\geq 80\%$  achieved as per assignment goal.

---

## 9. References

- scikit-learn documentation – <https://scikit-learn.org/stable/modules/clustering.html>
- scikit-learn-extra (KMedoids) – <https://scikit-learn-extra.readthedocs.io/>
- StackAbuse tutorial on Hierarchical Clustering
- Assignment #4 guidelines (Pawan Kumar Singh, JU IT Dept.)

**Github Repo-** <https://github.com/immu729/Machine-Learning-Lab/tree/main/Assignment4>

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans, DBSCAN, OPTICS,
AgglomerativeClustering

from sklearn.metrics import (
    adjusted_rand_score,
    adjusted_mutual_info_score,
    normalized_mutual_info_score,
    mutual_info_score,
    silhouette_score,
    calinski_harabasz_score,
    davies_bouldin_score,
)
from scipy.cluster.hierarchy import linkage, dendrogram
from sklearn.neighbors import NearestNeighbors
from math import comb
import warnings

warnings.filterwarnings("ignore")
np.random.seed(42)

def ensure_array(X):
    """Return numpy array of floats for feature matrices."""
    X = np.asarray(X, dtype=float)
    return X

def encode_labels(y):
    """Return integer labels starting from 0."""
    le = LabelEncoder()
    return le.fit_transform(np.asarray(y))

def rand_index(y_true, y_pred):
    """Unadjusted Rand Index (pair counting)."""
    n = len(y_true)
    tp_plus_tn = 0
    for i in range(n):
        for j in range(i + 1, n):
            same_true = (y_true[i] == y_true[j])
            same_pred = (y_pred[i] == y_pred[j])
            if (same_true and same_pred) or (not same_true and not
same_pred):
                tp_plus_tn += 1
    total_pairs = comb(n, 2)
    return tp_plus_tn / total_pairs

```

```

def safe_silhouette(X, labels):
    try:
        u = set(labels)
        if len(u) > 1 and len(u) < len(labels):
            return float(silhouette_score(X, labels))
        else:
            return np.nan
    except:
        return np.nan

def compute_sse_ssb(X, labels):
    """Cohesion (SSE) and Separation (SSB). Ignore label -1
    (noise)."""
    X = np.asarray(X, dtype=float)
    labels = np.asarray(labels)
    overall_mean = X.mean(axis=0)
    sse = 0.0
    ssb = 0.0
    for lab in np.unique(labels):
        if lab == -1:
            continue
        pts = X[labels == lab]
        if pts.shape[0] == 0:
            continue
        centroid = pts.mean(axis=0)
        sse += ((pts - centroid) ** 2).sum()
        ssb += pts.shape[0] * ((centroid - overall_mean) ** 2).sum()
    return float(sse), float(ssb)

def relabel_consecutive(labels):
    """Map labels to 0..k-1, keep -1 for noise."""
    labels = np.asarray(labels)
    mapping = {}
    next_label = 0
    out = labels.copy()
    for i, l in enumerate(labels):
        if l == -1:
            out[i] = -1
            continue
        if l not in mapping:
            mapping[l] = next_label
            next_label += 1
        out[i] = mapping[l]
    return out

def pam_kmedoids(X, k, max_iter=100, random_state=42):
    """
    Simple PAM implementation:
    - initialize medoids by random sampling
    - iteratively try swapping medoid with non-medoid to reduce total

```

```

cost
Returns labels (0..k-1)
"""
    rng = np.random.RandomState(random_state)
    X = np.asarray(X, dtype=float)
    n = X.shape[0]
    if k >= n:
        return np.arange(n)
    # initial medoids indices
    medoid_idx = rng.choice(n, size=k, replace=False)
    medoid_idx = medoid_idx.tolist()
    # compute pairwise distances once
    dists = np.linalg.norm(X[:, None, :] - X[None, :, :], axis=2) #
    shape (n, n)
    def total_cost(meds):
        # cost = sum distance each point to its nearest medoid
        return dists[:, meds].min(axis=1).sum()
    current_cost = total_cost(medoid_idx)
    for it in range(max_iter):
        improved = False
        for m_i, med in enumerate(medoid_idx):
            for o in range(n):
                if o in medoid_idx:
                    continue
                # try swap medoid med with o
                trial = medoid_idx.copy()
                trial[m_i] = o
                trial_cost = dists[:, trial].min(axis=1).sum()
                if trial_cost < current_cost:
                    medoid_idx = trial
                    current_cost = trial_cost
                    improved = True
                    break # accept first improving swap (classic PAM)
            if improved:
                break
        if not improved:
            break
    # final assignment
    labels = np.argmin(dists[:, medoid_idx], axis=1)
    return labels

from sklearn.cluster import KMeans as SKKMeans
def bisecting_kmeans(X, k, random_state=42, max_tries=10):
    """
    Bisecting KMeans:
    - start with one cluster containing all points
    - repeatedly split the cluster with largest SSE using
    KMeans(k=2)
    - stop when number of clusters = k
    """

```

```

X = np.asarray(X, dtype=float)
clusters = {0: np.arange(len(X))}
while len(clusters) < k:
    # compute SSE per cluster
    sse_per_cluster = {}
    for cid, idxs in clusters.items():
        pts = X[idxs]
        if pts.shape[0] == 0:
            sse_per_cluster[cid] = 0
        else:
            cen = pts.mean(axis=0)
            sse_per_cluster[cid] = ((pts - cen) ** 2).sum()
    # pick cluster with largest SSE to bisect
    to_split = max(sse_per_cluster, key=sse_per_cluster.get)
    idxs = clusters.pop(to_split)
    if len(idxs) <= 1:
        clusters[to_split] = idxs
        break
    # run KMeans(k=2) multiple times and pick best bisection
    best_labels = None
    best_inertia = np.inf
    for _ in range(max_tries):
        km = SKKMeans(n_clusters=2, n_init=10,
random_state=random_state)
        sub_labels = km.fit_predict(X[idxs])
        inert = km.inertia_
        if inert < best_inertia:
            best_inertia = inert
            best_labels = sub_labels
    # assign two new cluster ids
    new_id = max(clusters.keys(), default=-1) + 1
    clusters[new_id] = idxs[best_labels == 0]
    clusters[new_id + 1] = idxs[best_labels == 1]
# build final labels
final_labels = np.empty(len(X), dtype=int)
mapping = {}
next_id = 0
for cid, idxs in clusters.items():
    mapping[cid] = next_id
    final_labels[idxs] = next_id
    next_id += 1
return final_labels

def run_all_for_dataset(X, y_true, dataset_name, show_plots=True):
    X = ensure_array(X)
    y_true = encode_labels(y_true)
    n_clusters = len(np.unique(y_true))

    print(f"\n=== DATASET: {dataset_name} (n={X.shape[0]},
features={X.shape[1]}, classes={n_clusters}) ===\n")

```



```

# standardize
scaler = StandardScaler()
Xs = scaler.fit_transform(X)

# PCA for 2D plotting
pca = PCA(n_components=2, random_state=42)
X2 = pca.fit_transform(Xs)

results = []

# ----- KMeans (k-means++) -----
km = KMeans(n_clusters=n_clusters, init='k-means++', n_init=20,
random_state=42)
km_labels = relabel_consecutive(km.fit_predict(Xs))
sse, ssb = compute_sse_ssb(Xs, km_labels)
results.append(("KMeans (k-means++)", km_labels, sse, ssb))

# ----- K-Medoids (PAM) -----
pam_labels = relabel_consecutive(pam_kmedoids(Xs, n_clusters,
random_state=42))
sse, ssb = compute_sse_ssb(Xs, pam_labels)
results.append(("KMedoids (PAM)", pam_labels, sse, ssb))

# ----- Bisecting KMeans -----
bis_labels = relabel_consecutive(bisecting_kmeans(Xs, n_clusters,
random_state=42))
sse, ssb = compute_sse_ssb(Xs, bis_labels)
results.append(("Bisecting KMeans", bis_labels, sse, ssb))

# ----- Hierarchical (Agglomerative) -----
# plot dendrogram
Z = linkage(Xs, method='ward')
if show_plots:
    plt.figure(figsize=(10, 4))
    dendrogram(Z, truncate_mode='lastp', p=40,
show_leaf_counts=True)
    plt.title(f"Dendrogram (Ward) - {dataset_name}")
    plt.tight_layout()
    plt.show()
agg = AgglomerativeClustering(n_clusters=n_clusters,
linkage='ward')
agg_labels = relabel_consecutive(agg.fit_predict(Xs))
sse, ssb = compute_sse_ssb(Xs, agg_labels)
results.append(("Hierarchical (Agglomerative - ward)", agg_labels,
sse, ssb))

# ----- DBSCAN -----
# heuristic to choose eps: 90th percentile of 5-NN distances
neigh = NearestNeighbors(n_neighbors=5).fit(Xs)

```

```

dists, _ = neigh.kneighbors(Xs)
kdist = np.sort(dists[:, -1])
eps_guess = float(np.percentile(kdist, 90))
db = DBSCAN(eps=eps_guess, min_samples=5)
db_labels = relabel_consecutive(db.fit_predict(Xs))
sse, ssb = compute_sse_ssb(Xs, db_labels)
results.append((f"DBSCAN (eps≈{eps_guess:.3f})", db_labels, sse,
sse))

# ----- OPTICS -----
opt = OPTICS(min_samples=5)
opt_labels = relabel_consecutive(opt.fit_predict(Xs))
sse, ssb = compute_sse_ssb(Xs, opt_labels)
results.append(("OPTICS", opt_labels, sse, ssb))

# ----- Compute metrics for each algorithm and report
-----
rows = []
for algo_name, labels, sse_val, ssb_val in results:
    labels_arr = np.asarray(labels)
    # Replace class names with numeric values already done (y_true
is numeric)
    # Compute metrics (use adjusted versions and mutual info
variants)
    try:
        ri = rand_index(y_true, labels_arr)
    except Exception:
        ri = np.nan
    try:
        ari = float(adjusted_rand_score(y_true, labels_arr))
    except:
        ari = np.nan
    try:
        mi = float(mutual_info_score(y_true, labels_arr))
    except:
        mi = np.nan
    try:
        ami = float(adjusted_mutual_info_score(y_true,
labels_arr))
    except:
        ami = np.nan
    try:
        nmi = float(normalized_mutual_info_score(y_true,
labels_arr))
    except:
        nmi = np.nan
    sil = safe_silhouette(Xs, labels_arr)
    try:
        ch = float(calinski_harabasz_score(Xs, labels_arr)) if

```

```

len(set(labels_arr)) > 1 else np.nan
    except:
        ch = np.nan
    try:
        dbi = float(davies_bouldin_score(Xs, labels_arr)) if
len(set(labels_arr)) > 1 else np.nan
    except:
        dbi = np.nan

    rows.append({
        "Algorithm": algo_name,
        "n_clusters": len([l for l in np.unique(labels_arr) if l != -1]),
        "Rand": ri,
        "AdjRand": ari,
        "MutualInfo": mi,
        "AdjMutualInfo": ami,
        "NormMutualInfo": nmi,
        "Silhouette": sil,
        "CalinskiHarabasz": ch,
        "DaviesBouldin": dbi,
        "SSE": sse_val,
        "SSB": ssb_val
    })

df = pd.DataFrame(rows)
# nicer ordering of columns
df = df[["Algorithm", "n_clusters", "Rand", "AdjRand",
        "MutualInfo", "AdjMutualInfo",
        "NormMutualInfo", "Silhouette", "CalinskiHarabasz",
        "DaviesBouldin", "SSE", "SSB"]]

# Display summary table
print(f"\nPerformance summary for {dataset_name}:\n")
display(df.round(4))

# Show PCA 2D visualizations for each algorithm
if show_plots:
    nplots = len(results)
    cols = 3
    rowsplt = int(np.ceil(nplots / cols))
    fig, axes = plt.subplots(rowsplt, cols, figsize=(5 * cols, 4 *
rowsplt))
    axes = axes.flatten()
    for ax, (algo_name, labels, _, _) in zip(axes, results):
        ax.scatter(X2[:, 0], X2[:, 1], c=labels, cmap='tab10',
s=30)
        ax.set_title(algo_name)
        ax.set_xticks([]); ax.set_yticks([])
    # hide extra axes

```

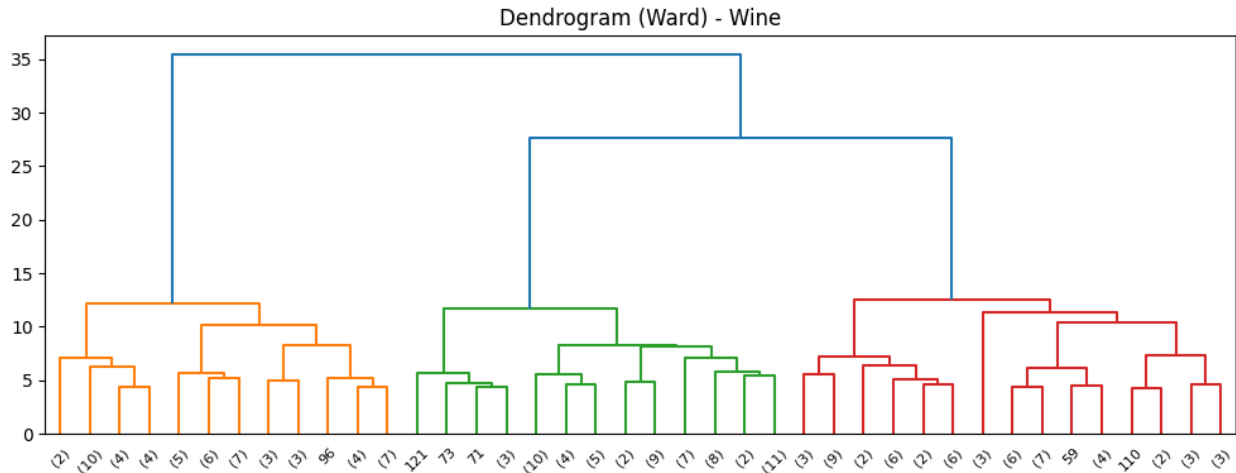


```

\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "Rand\","\n      \ "properties\": {\n      \ "dtype\": \ "number\","\n
\ "std\": 0.12838689445058896,\n      \ "min\": 0.5121,\n
\ "max\": 0.8415,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.8415,\n      0.5121,\n
0.8252\n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "AdjRand\","\n      \ "properties\": {\n      \ "dtype\": \ "number\","\n
\ "std\": 0.22992947324487711,\n      \ "min\": 0.0514,\n
\ "max\": 0.6416,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.6416,\n      0.0514,\n
0.6153\n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "MutualInfo\","\n      \ "properties\": {\n      \ "dtype\":
\ "number\","\n      \ "std\": 0.16800982114150353,\n      \ "min\":
0.3082,\n      \ "max\": 0.7405,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.7405,\n      0.3082,\n
0.7228\n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "AdjMutualInfo\","\n      \ "properties\": {\n      \ "dtype\":
\ "number\","\n      \ "std\": 0.16400464322695257,\n      \ "min\":
0.2657,\n      \ "max\": 0.6807,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.671,\n      0.2657,\n      0.6713\n
n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "NormMutualInfo\","\n      \ "properties\": {\n      \ "dtype\":
\ "number\","\n      \ "std\": 0.15492316482695542,\n      \ "min\":
0.2924,\n      \ "max\": 0.6858,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.675,\n      0.2924,\n      0.6755\n
n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "Silhouette\","\n      \ "properties\": {\n      \ "dtype\":
\ "number\","\n      \ "std\": 0.3139194827765022,\n      \ "min\": -
0.3009,\n      \ "max\": 0.5081,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.4566,\n      -0.3009,\n
0.4467\n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "CalinskiHarabasz\","\n      \ "properties\": {\n      \ "dtype\":
\ "number\","\n      \ "std\": 95.6462088170148,\n      \ "min\":
8.3282,\n      \ "max\": 241.9044,\n      \ "num_unique_values\":
5,\n      \ "samples\": [\n      238.763,\n      8.3282,\n
222.7192\n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":
\ "DaviesBouldin\","\n      \ "properties\": {\n      \ "dtype\":
\ "number\","\n      \ "std\": 1.0784846344137995,\n      \ "min\":
0.8035,\n      \ "max\": 3.2757,\n      \ "num_unique_values\": 5,\n
\ "samples\": [\n      0.8419,\n      2.4253,\n
0.8035\n      ],\n      \ "semantic_type\": \ "\",\n
\ "description\": \ "\n      }\n    },\n    {\n      \ "column\":

```





#### Performance summary for Wine:

```
{
  "summary": {
    "name": "df_wine, res_wine = run_all_for_dataset(wine",
    "rows": 6,
    "fields": [
      {
        "column": "Algorithm",
        "properties": {
          "dtype": "string",
          "num_unique_values": 6,
          "samples": [
            "KMeans (k-means++)",
            "KMedoids (PAM)",
            "OPTICS"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "n_clusters",
        "properties": {
          "dtype": "number",
          "std": 0,
          "min": 1,
          "max": 4,
          "num_unique_values": 3,
          "samples": [
            3,
            1,
            4
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "AdjRand",
        "properties": {
          "dtype": "number",
          "std": 0.2569408933328182,
          "min": 0.3646,
          "max": 0.9543,
          "num_unique_values": 6,
          "samples": [
            0.9543,
            0.884,
            0.4392
          ],
          "semantic_type": ""
        },
        "column": "Rand",
        "properties": {
          "dtype": "number",
          "std": 0.3952078410996759,
          "min": -0.0074,
          "max": 0.8975,
          "num_unique_values": 6,
          "samples": [
            0.8975,
            0.7411,
            0.0358
          ],
          "semantic_type": ""
        },
        "column": "MutualInfo",
        "properties": {
          "dtype": "number",
          "std": 0.39782152489108313,
          "min": 0.0324,
          "max": 0.9545,
          "num_unique_values": 6,
          "samples": [
            0.9545,
            0.849,
            0.1621
          ],
          "semantic_type": ""
        }
      ]
    }
  }
}
```

```

\"AdjMutualInfo\", \n      \"properties\": { \n      \"dtype\":
\"number\", \n      \"std\": 0.3580637061566931, \n      \"min\":
0.041, \n      \"max\": 0.8746, \n      \"num_unique_values\": 6, \n
\"samples\": [ \n      0.8746, \n      0.7806, \n      0.168\
n      ], \n      \"semantic_type\": \"\", \n
\"description\": \"\" \n      } \n      }, \n      { \n      \"column\":
\"NormMutualInfo\", \n      \"properties\": { \n      \"dtype\":
\"number\", \n      \"std\": 0.3506026910906418, \n      \"min\":
0.0504, \n      \"max\": 0.8759, \n      \"num_unique_values\": 6, \n
\"samples\": [ \n      0.8759, \n      0.7829, \n
0.1949 \n      ], \n      \"semantic_type\": \"\", \n
\"description\": \"\" \n      } \n      }, \n      { \n      \"column\":
\"Silhouette\", \n      \"properties\": { \n      \"dtype\":
\"number\", \n      \"std\": 0.1619759735269401, \n      \"min\": -
0.1336, \n      \"max\": 0.2849, \n      \"num_unique_values\": 6, \n
\"samples\": [ \n      0.2849, \n      0.2676, \n      -
0.1336 \n      ], \n      \"semantic_type\": \"\", \n
\"description\": \"\" \n      } \n      }, \n      { \n      \"column\":
\"CalinskiHarabasz\", \n      \"properties\": { \n      \"dtype\":
\"number\", \n      \"std\": 31.921483001243327, \n      \"min\":
4.8536, \n      \"max\": 70.94, \n      \"num_unique_values\": 6, \n
\"samples\": [ \n      70.94, \n      67.1223, \n
5.0571 \n      ], \n      \"semantic_type\": \"\", \n
\"description\": \"\" \n      } \n      }, \n      { \n      \"column\":
\"DaviesBouldin\", \n      \"properties\": { \n      \"dtype\":
\"number\", \n      \"std\": 0.6029515425526887, \n      \"min\":
1.3892, \n      \"max\": 2.9367, \n      \"num_unique_values\": 6, \n
\"samples\": [ \n      1.3892, \n      1.4248, \n
1.6194 \n      ], \n      \"semantic_type\": \"\", \n
\"description\": \"\" \n      } \n      }, \n      { \n      \"column\":
\"SSE\", \n      \"properties\": { \n      \"dtype\": \"number\", \n
\"std\": 653.2299276627753, \n      \"min\": 38.8681, \n
\"max\": 2054.3232, \n      \"num_unique_values\": 6, \n
\"samples\": [ \n      1277.9285, \n      1309.481, \n
38.8681 \n      ], \n      \"semantic_type\": \"\", \n
\"description\": \"\" \n      } \n      }, \n      { \n      \"column\":
\"SSB\", \n      \"properties\": { \n      \"dtype\": \"number\", \n
\"std\": 462.2836737482171, \n      \"min\": 3.14, \n      \"max\":
1036.0715, \n      \"num_unique_values\": 6, \n      \"samples\": [ \
n      1036.0715, \n      1004.519, \n      225.8965 \n
], \n      \"semantic_type\": \"\", \n      \"description\": \"\" \n
} \n      } \n      ] \n    }, \"type\": \"dataframe\"}

```



PCA (2D) visualizations - Wine

