# ▾ K-Nearest Neighbors (KNN) project

**Objective:**

1. investigate how different values of k and distance metrics affect the accuracy of a KNN classifier on a specific dataset.
2. Visualize the decision boundaries for k iterations

**Steps:**

1.Select a Dataset from UCI depository, preprocess data to handle missing values etc.

2.Implement a KNN Algorithm: for example Minkowski distance metric. (Use libraries like scikit-learn)

3.Experiment with Different k Values: Conduct experiments with various values of k and observe how the choice of k affects the performance of the model. Plot accuracy for each k value to visualize the results.

4.Experiment with Different Distance Metrics: Metrics such as Euclidean distance, Manhattan distance, and Minkowski distance with various p values (power parameter) to observe how the choice of distance metric influences the model's performance.

5.Creat a synthetic two variable dataset to demonstrate decision boundaries of KNN at various iterations of k.

```
1 from google.colab import files
2 uploaded = files.upload()
```

Choose Files | No file chosen    Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving test.csv to test.csv

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn.neighbors import KNeighborsClassifier
6 from sklearn.metrics import accuracy_score
```

```
1 pip install ucimlrepo
```

```
1 from ucimlrepo import fetch_ucirepo
2 # fetch dataset
3 heart_disease = fetch_ucirepo(id=45)
4 # data (as pandas dataframes)
5 X = heart_disease.data.features
6 y = heart_disease.data.targets
```

```
1 combined_df = pd.concat([X, y], axis=1)
2 combined_df = combined_df.dropna()
3 combined_df.reset_index(drop=True, inplace=True)
```

```
1 combined_df.head(5)
```

|   | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | num |
|---|-----|-----|----|----------|------|-----|---------|---------|-------|---------|-------|-----|------|-----|
| 0 | 63 | 1 | 1 | 145 | 233 | 1 | 2 | 150 | 0 | 2.3 | 3 | 0.0 | 6.0 | 0 |
| 1 | 67 | 1 | 4 | 160 | 286 | 0 | 2 | 108 | 1 | 1.5 | 2 | 3.0 | 3.0 | 2 |
| 2 | 67 | 1 | 4 | 120 | 229 | 0 | 2 | 129 | 1 | 2.6 | 2 | 2.0 | 7.0 | 1 |
| 3 | 37 | 1 | 3 | 130 | 250 | 0 | 0 | 187 | 0 | 3.5 | 3 | 0.0 | 3.0 | 0 |
| 4 | 41 | 0 | 2 | 130 | 204 | 0 | 2 | 172 | 0 | 1.4 | 1 | 0.0 | 3.0 | 0 |

```
1 from sklearn.model_selection import train_test_split
2 x= combined_df.iloc[:,0:13].values
3 y= combined_df['num'].values
4
5 x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)
6 from sklearn.preprocessing import StandardScaler
7 st_x= StandardScaler()
```
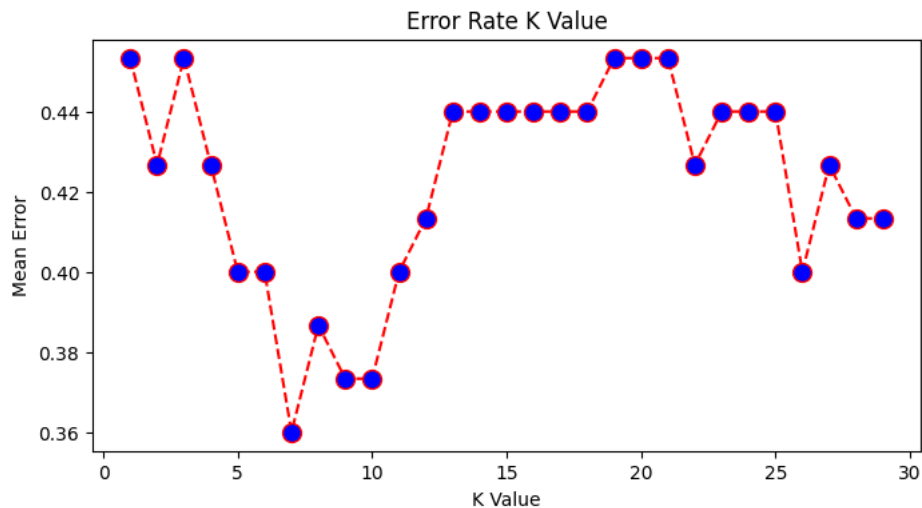
```
8 x_train= st_x.fit_transform(x_train)
```

```
1 error = []
2 # Calculating error for K values between 1 and 30
3 for i in range(1, 30):
4     knn = KNeighborsClassifier(n_neighbors=i)
5     knn.fit(x_train, y_train)
6     pred_i = knn.predict(x_test)
7     error.append(np.mean(pred_i != y_test))
8 plt.figure(figsize=(8, 4))
9 plt.plot(range(1, 30), error, color='red', linestyle='dashed', marker='o',
10         markerfacecolor='blue', markersize=10)
11 plt.title('Error Rate K Value')
12 plt.xlabel('K Value')
13 plt.ylabel('Mean Error')
14 print("Minimum error:-",min(error),"at K =",error.index(min(error))+1)
15 plt.show()
```

```
Minimum error:- 0.36 at K = 7
```



Error Rate K Value

## ▾ Experiment with various k values and Distance Metrics on UCI Health Dataset

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.metrics import accuracy_score
5 import matplotlib.pyplot as plt
6
7 # Create empty lists to store error, accuracy, and metrics
8 error = []
9 accuracy = []
10 metrics_list = []
11
12 # Create a DataFrame to store results
13 results_df = pd.DataFrame(columns=['K', 'Mean Error', 'Accuracy', 'Metric'])
14
15 # Assuming you have already split your data into x_train, x_test, y_train, and y_test
16
17 # Define the list of distance metrics to iterate over
18 distance_metrics = ['euclidean', 'chebyshev', 'minkowski','cosine']
19
20 for metric in distance_metrics:
21     for k in range(1, 16):
22         #knn = KNeighborsClassifier(n_neighbors=k, metric=metric)
23         if metric == 'minkowski':
24           knn = KNeighborsClassifier(n_neighbors=k, metric=metric, metric_params={'p': 3})
25         else:
26           knn = KNeighborsClassifier(n_neighbors=k, metric=metric)
27         knn.fit(x_train, y_train)
28         pred_i = knn.predict(x_test)
29
30         # Calculate mean error and accuracy
31         mean_error = np.mean(pred_i != y_test)
```
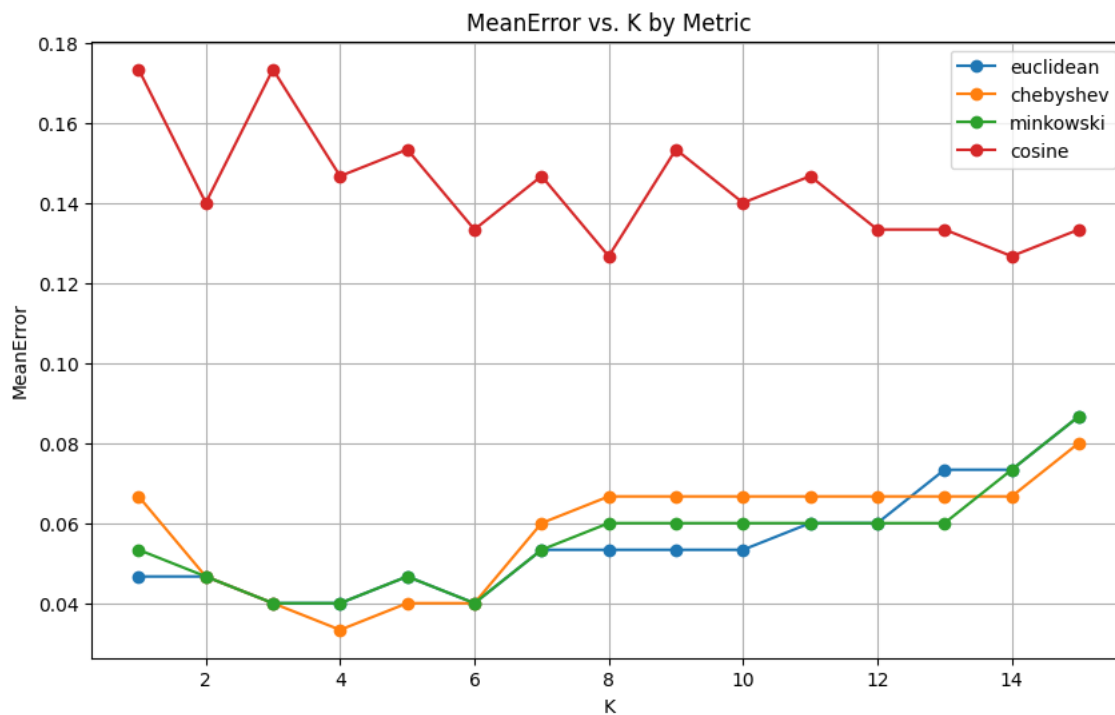
```
32          acc = accuracy_score(y_test, pred_i)
33
34          # Append values to the lists
35          error.append(mean_error)
36          accuracy.append(acc)
37          metrics_list.append(metric)
38
39          # Append values to the DataFrame
40          results_df = results_df.append({'K': k, 'MeanError': mean_error, 'Accuracy': acc, 'Metric': metric}, ignore_index=Tru
41
42 results_df.to_csv("knn.csv", index=False)
43
44
45
```
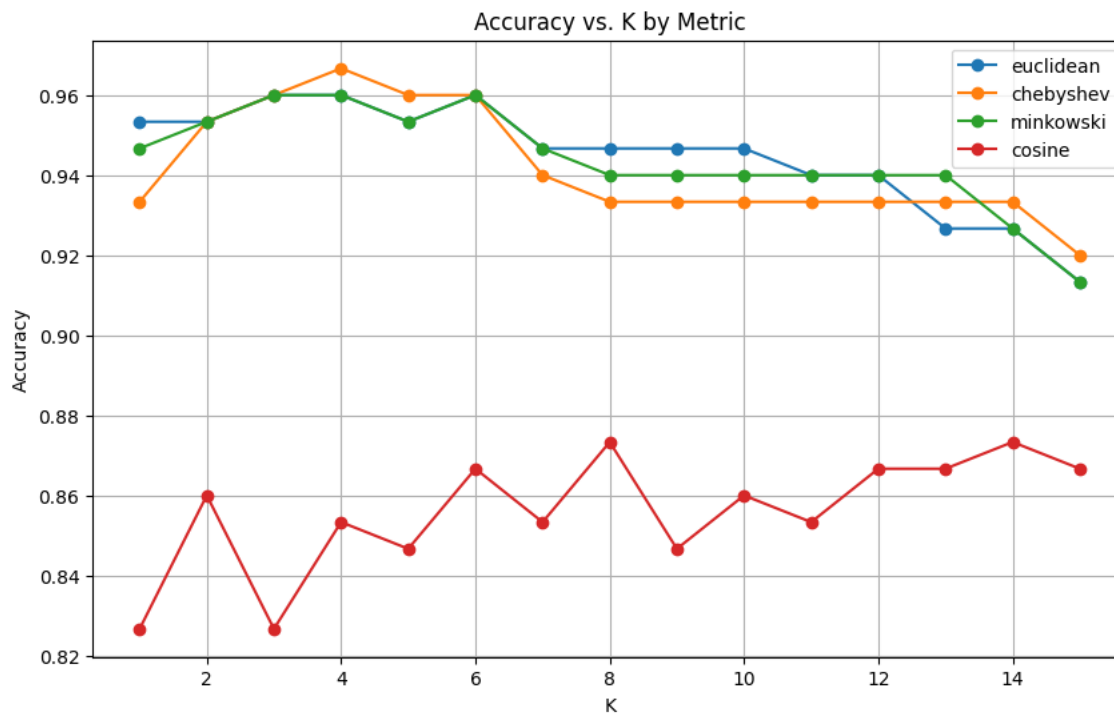
```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 # Create a figure and axis
5 fig, ax = plt.subplots(figsize=(10, 6))
6 df=results_df
7
8 # Iterate over unique metrics to create lines
9 for metric in df['Metric'].unique():
10    subset = df[df['Metric'] == metric]
11    ax.plot(subset['K'], subset['MeanError'], marker='o', label=metric)
12
13 # Set labels and legend
14 ax.set_xlabel('K')
15 ax.set_ylabel('MeanError')
16 ax.set_title('MeanError vs. K by Metric')
17 ax.legend()
18
19 # Show the plot
20 plt.grid(True)
21 plt.show()
22
```



```
1    import matplotlib.pyplot as plt
2    import pandas as pd
3
4    # Create a figure and axis
5    fig, ax = plt.subplots(figsize=(10, 6))
6    df=results_df
7
8    # Iterate over unique metrics to create lines
9    for metric in df['Metric'].unique():
```

```
10      subset = df[df['Metric'] == metric]
11      ax.plot(subset['K'], subset['Accuracy'], marker='o', label=metric)
12
13  # Set labels and legend
14  ax.set_xlabel('K')
15  ax.set_ylabel('Accuracy')
16  ax.set_title('Accuracy vs. K by Metric')
17  ax.legend()
18
19  # Show the plot
20  plt.grid(True)
21  plt.show()
22
```
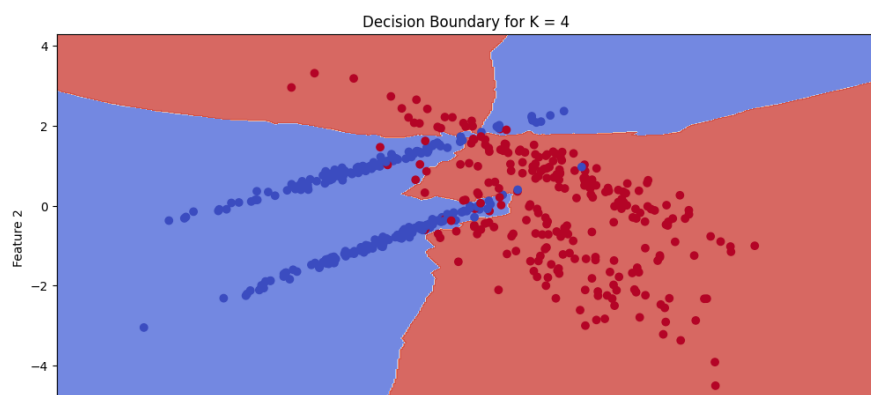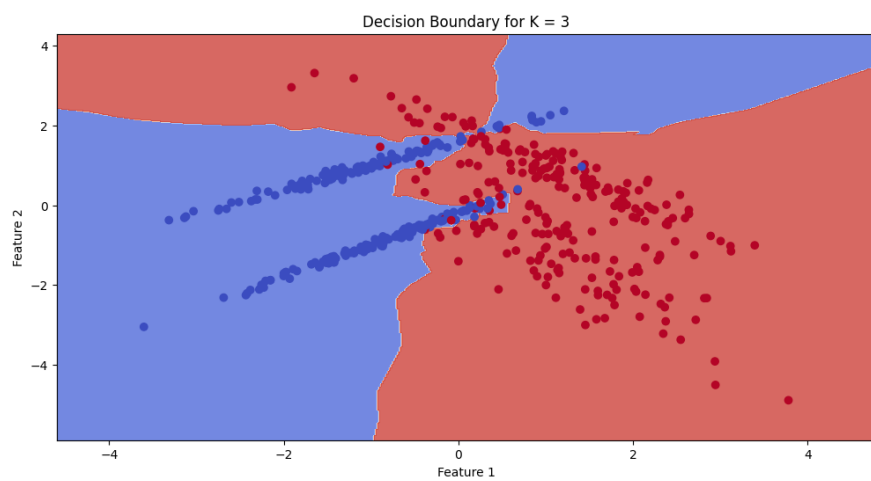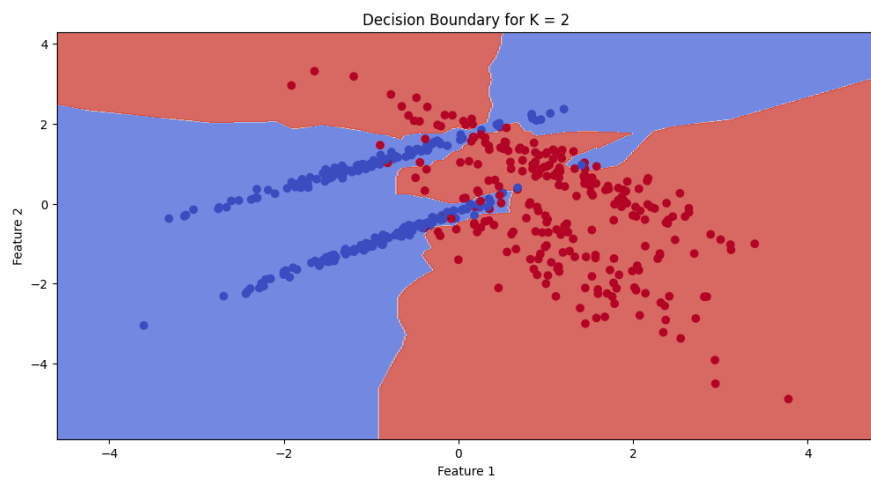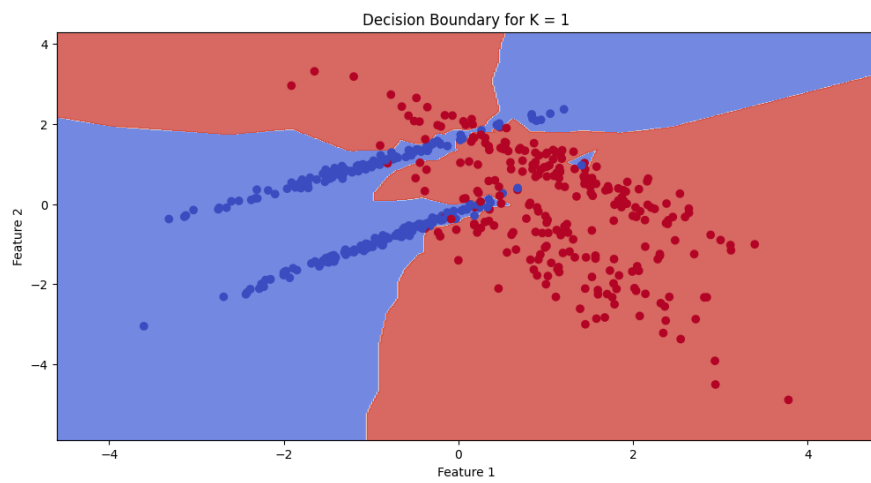


Accuracy vs. K by Metric

## Create A synthetic dataset to demonstrate the Decision Boundaries at various Iterations

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.datasets import make_classification
5 from sklearn.model_selection import train_test_split
6 # Generate a synthetic dataset for visualization
7 X, y = make_classification(n_samples=500, n_features=2, n_informative=2, n_redundant=0, random_state=42)
8 x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
9
10 # Calculate error for K values between 1 and 30
11 error = []
12 for k in range(1, 6):
13      knn = KNeighborsClassifier(n_neighbors=k)
14      knn.fit(x_train, y_train)
15      pred_i = knn.predict(x_test)
16      error.append(np.mean(pred_i != y_test))
17
18      # Create a meshgrid for the decision boundary plot
19      h = 0.02   # Step size in the mesh
20      x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
21      y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
22      xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
23      Z = knn.predict(np.c_[xx.ravel(), yy.ravel()])
24      Z = Z.reshape(xx.shape)
25
26      # Plot the decision boundary for the current K
27      plt.figure(figsize=(10, 5))
28      plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
```

```
29    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm)
30    plt.title(f'Decision Boundary for K = {k}')
31    plt.xlabel('Feature 1')
32    plt.ylabel('Feature 2')
33    plt.show()
```

Decision Boundary for K = 1

Decision Boundary for K = 2

Decision Boundary for K = 3
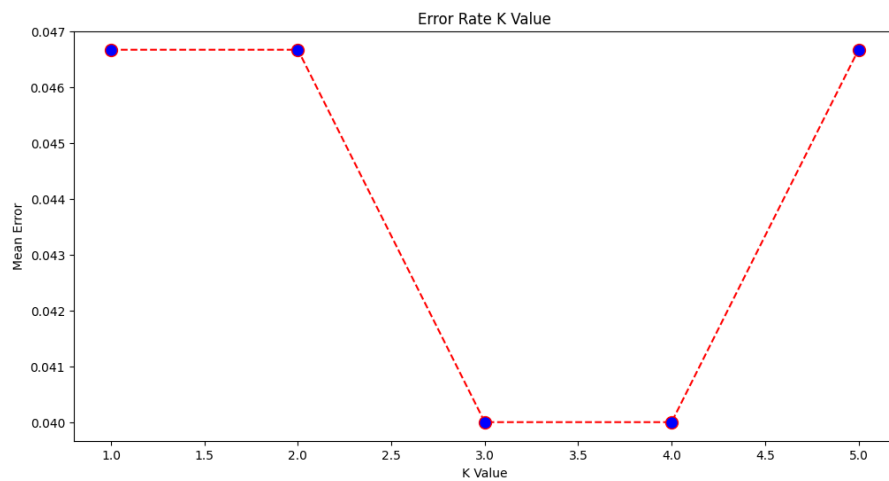
Decision Boundary for K = 4

```
1
2    # Plot the error vs. K
3    plt.figure(figsize=(12, 6))
4    plt.plot(range(1, 6), error, color='red', linestyle='dashed', marker='o',
```

```
5              markerfacecolor='blue', markersize=10)
6   plt.title('Error Rate K Value')
7   plt.xlabel('K Value')
8   plt.ylabel('Mean Error')
9   print("Minimum error:", min(error), "at K =", error.index(min(error)) + 1)
```

Minimum error: 0.04 at K = 3



## Summary:

In this project we focused on evaluating the performance of a k-Nearest Neighbors (k-NN) classifier using different distance metrics. The primary objective is to determine the most suitable distance metric and the optimal value of k for a classification task. The code systematically explores a range of k values from 1 to 15 and four different distance metrics: Euclidean, Chebyshev, Minkowski with a specified parameter, and Cosine.

This exercise demonstrates the importance of hyperparameter tuning, specifically the choice of distance metric and the selection of k, in optimizing the performance of k-NN classifiers for classification tasks.

Also the resulting visualizations display the decision boundaries for different K values, helping to understand how the choice of K affects the model's capacity to distinguish between classes. This exercise serves as an informative illustration of the impact of K on k-NN classification performance.