

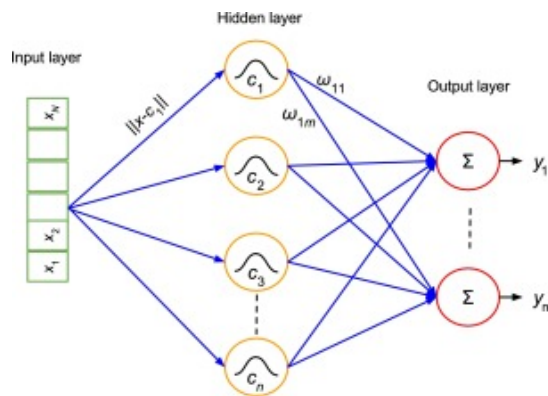
Exercício 7

Arthur Felipe Reis Souza

April 29, 2024

1 Introduction

As redes neurais do tipo RBF (Radial Basis Function) são redes neurais com uma camada de entrada, uma camada intermediária e uma camada de saída. A camada intermediária desse tipo de rede aplica N (N é o número de neurônios da camada intermediária) funções radiais sobre os dados de entrada, gerando uma saída H que alimentará a camada de saída. Uma função radial é uma função que mensura a distância entre os pontos e um valor arbitrário. A distância pode ser entendida como uma medida de dissimilaridade entre dois ou mais valores. Há várias métricas de distâncias que podem ser utilizadas, sendo elas: Euclidean distance, Manhattan distance, Minowski distance (forma generalizada das duas anteriores). Portanto, a ideia principal das redes do tipo RBF é: pegar as entradas da rede e aplicar sobre essa entrada um conjunto de N funções radiais (função gaussiana é comumente utilizada). A saída de cada função radial será combinada linearmente com os parâmetros da camada de saída da rede, que poderá aplicar uma função de ativação sigmoidal ou identidade sobre essa combinação.

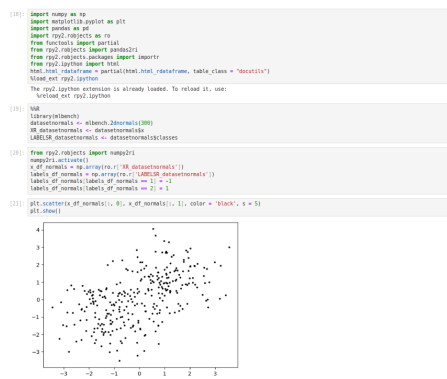


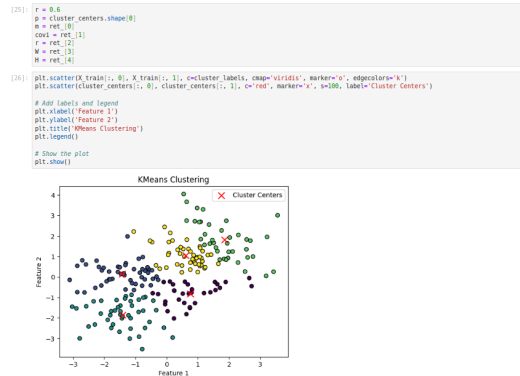
Faz-se mister, portanto, encontrar a localização dos centros, ou do valor central dessas funções radiais que irão aplicar, sobre os dados de entrada, uma

A primeira parte do exercícios 7 consiste em utilizar, nos dados obtidos no último Exercício ELM, o algoritmo KMeans para agrupar os mesmos. Após aplicar o algoritmo KMeans, usando $K = 5$, obtemos os seguintes resultados abaixo :

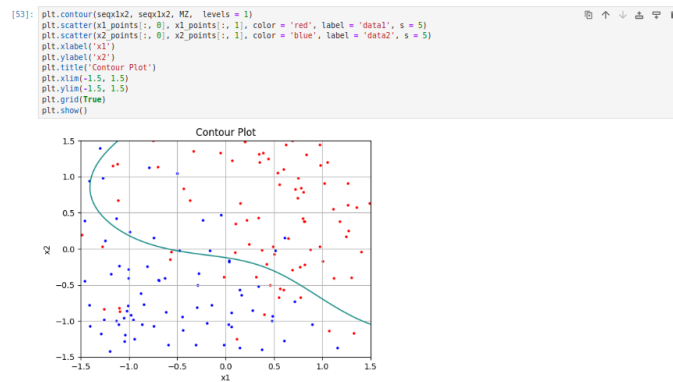
1.1 2dnormals

Dados para $K = 5$, dataset normals :





A superfície de separação para $K = 5$ está registrada abaixo :



Aumentando muito o valor de K teoricamente levaria o modelo ao overfitting. Porém, esse conjunto de dados é facilmente separável, e não foi possível observar um overfitting visualmente, porém com um K muito grande o modelo erra mais nos dados de teste. O valor de K que levará o modelo a ter uma melhor performance é obtido com o algoritmo de Grid Search Cross-Validation, e está registrado abaixo :

```
H shape : (210, 1)
Haug shape : (210, 2)
ylin shape : (210, 1)
M shape : (2, 1)
H shape : (210, 2)
Haug shape : (210, 3)
ylin shape : (210, 1)
M shape : (2, 1)
H shape : (210, 3)
Haug shape : (210, 4)
ylin shape : (210, 1)
M shape : (4, 1)
H shape : (210, 4)
Haug shape : (210, 5)
ylin shape : (210, 1)
M shape : (5, 1)
H shape : (210, 5)
Haug shape : (210, 6)
```

[illegible]

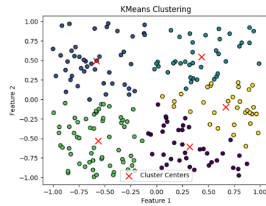
```
[21]: from sklearn.model_selection import train_test_split
      X_train, X_test, y_train, y_test = train_test_split(x, df_xor, labels=df_xor, random_state = 0, train_size = 0.7)

[22]: from sklearn.cluster import KMeans
      kmeans = KMeans(n_clusters = 5, random_state = 0, n_init = 'auto').fit(X_train)
      cluster_labels = kmeans.labels_
      cluster_centers = kmeans.cluster_centers_
      p = cluster_centers.shape[0]

[23]: plt.scatter(X_train[:, 0], X_train[:, 1], c=cluster_labels, cmap='viridis', markers='o', edgecolors='k')
      plt.scatter(cluster_centers[:, 0], cluster_centers[:, 1], c='red', markers='x', s=100, label='Cluster Centers')

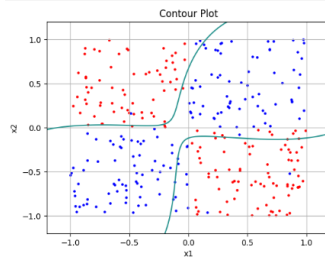
# Add labels and legend
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.title('KMeans Clustering')
plt.legend()

# Show the plot
plt.show()
```



A superfície de separação para $K = 5$ está registrada abaixo :

```
[198]: plt.contour(seqx1x2, seqx2x2, MZ, levels = 1)
      plt.scatter(x1_points[:, 0], x1_points[:, 1], color = 'red', label = 'data1', s = 5)
      plt.scatter(x2_points[:, 0], x2_points[:, 1], color = 'blue', label = 'data2', s = 5)
      plt.xlabel('x1')
      plt.ylabel('x2')
      plt.title('Contour Plot')
      plt.xlim(-1.2, 1.2)
      plt.ylim(-1.2, 1.2)
      plt.grid(True)
      plt.show()
```



Aumentando muito o valor de K teoricamente levaria o modelo ao overfitting. O valor ideal de K é obtido com o metodo Grid Search CV e está registrado abaixo :

```
[201]: from sklearn.metrics import accuracy_score
      lst_acc_test = list()
      lst_param = list()

      for i in range(1, 200, 1):
          kmeans = KMeans(n_clusters = i, random_state = 0, n_init = 'auto').fit(X_train)
          cluster_labels = kmeans.labels_
          cluster_centers = kmeans.cluster_centers_
          ret_ = train_BRF_train(BRF(x_train, y_train, i, f))

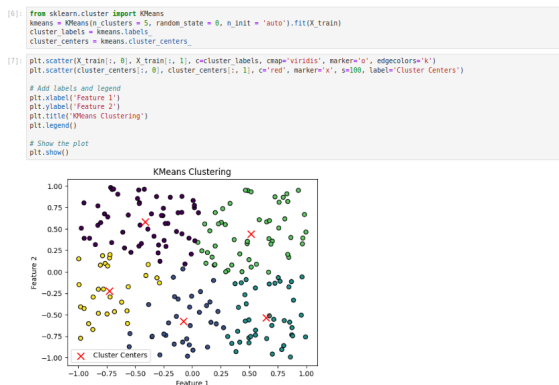
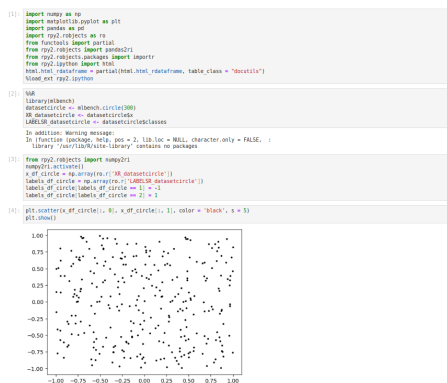
          yhat_test = train_BRF.y_BRF(X_test, ret_)
          yhat_test(yhat_test == -0.5) = 1
          yhat_test(yhat_test == -0.5) = 1
          yhat_train = train_BRF.y_BRF(X_train, ret_)
          yhat_train(yhat_train == -0.5) = 1
          yhat_train(yhat_train == -0.5) = 1
          lst_acc_test.append(accuracy_score(y_test, yhat_test))
          lst_param.append(i)
          lst_acc_test = np.array(lst_acc_test)
          lst_param = np.array(lst_param)

      H shape : (210, 1)
      H shape : (210, 2)
      yin shape : (210,)
      W shape : (2,1)
      H shape : (210, 2)
      H shape : (210, 3)
      yin shape : (210,)
      W shape : (3,1)
      H shape : (210, 3)
      H shape : (210, 4)
      yin shape : (210,)
      W shape : (4,1)
      H shape : (210, 4)
      H shape : (210, 5)
      yin shape : (210,)
      W shape : (5,1)
      H shape : (210, 5)
      H shape : (210, 6)

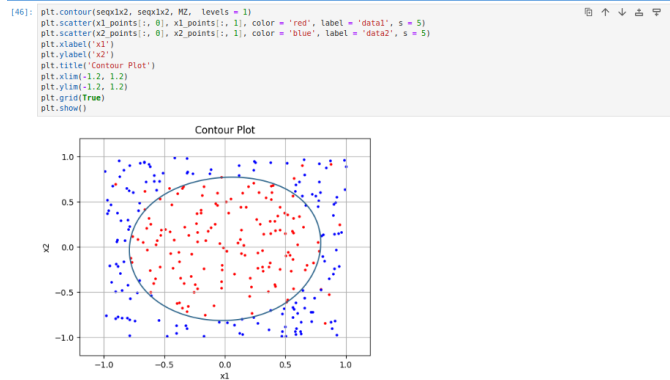
[202]: idx = np.argmax(lst_acc_test)
      print('The value of K that give us the best performance in test data is : {lst_param[idx]}')
      The value of K that give us the best performance in test data is : 44
```

1.3 circle

Dados para $K = 5$, dataset circle :



A superfície de separação para $K = 5$ está registrada abaixo :



Aumentando muito o valor de K teoricamente levaria o modelo ao overfitting. O valor ideal de K é obtido com o metodo Grid Search CV e está registrado abaixo :

```
[56]: from sklearn.metrics import accuracy_score
list_acc_test = list()
list_param = list()

for i in range(1, 200, 1):
    kmeans = KMeans(n_clusters = 1, random_state = 0, n_init = 'auto').fit(X_train)
    cluster_labels = kmeans.labels_
    cluster_centers = kmeans.cluster_centers_
    ret = train_RBF_train_RBF(X_train, y_train, 1, r)

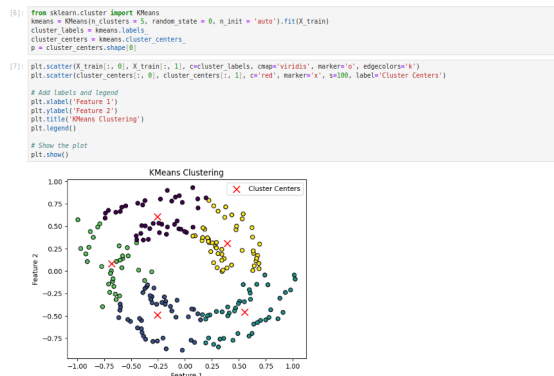
    yhat_test = train_RBF_y_RBF(X_test, ret)
    yhat_test[yhat_test <= -0.5] = -1
    yhat_test[yhat_test > -0.5] = 1
    yhat_train = train_RBF_y_RBF(X_train, ret)
    yhat_train[yhat_train <= -0.5] = -1
    yhat_train[yhat_train > -0.5] = 1
    list_acc_test.append(accuracy_score(y_test, yhat_test))
    list_param.append(i)
list_acc_test = np.array(list_acc_test)
list_param = np.array(list_param)

H shape : (210, 1)
Haup shape : (210, 2)
yin shape : (210,)
W shape : (2,)
H shape : (210, 2)
Haup shape : (210, 3)
yin shape : (210,)
W shape : (3,)
H shape : (210, 3)
Haup shape : (210, 4)
yin shape : (210,)
W shape : (4,)
H shape : (210, 4)
Haup shape : (210, 5)
yin shape : (210,)
W shape : (5,)
H shape : (210, 5)
Haup shape : (210, 6)

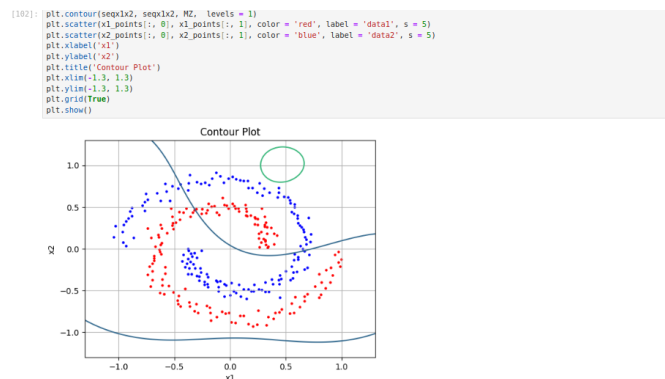
[57]: idx = np.argmax(list_acc_test)
print(f"The value of K that give us the best performance in test data is : {list_param[idx]}")
The value of K that give us the best performance in test data is : 82
```

1.4 spirals

Dados para K = 5, dataset spirals :



A superfície de separação para $K = 5$ está registrada abaixo :



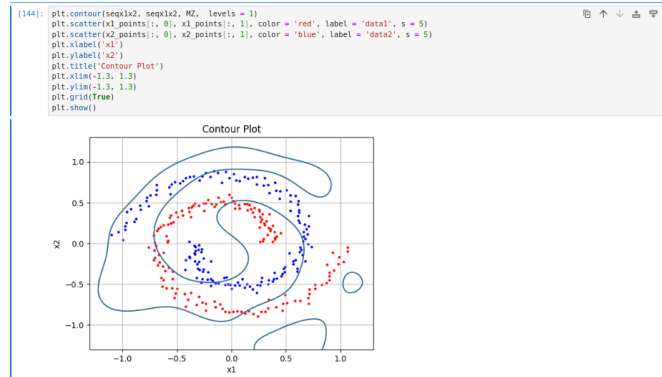
Aumentando muito o valor de K teoricamente levaria o modelo ao overfitting. Como pode ser observado para $K = 5$, o modelo não faz uma boa aproximação.

Para encontrar o melhor valor de K foi utilizado o algoritmo de Grid, que registra abaixo o melhor valor de K e também a superfície de separação para esse valor de K.

```
[113]: from sklearn.metrics import accuracy_score
      (list_acc_test = list())
      list_param = list()
      for i in range(1, 200, 1):
          kmeans = KMeans(n_clusters = 1, random_state = 0, n_init = 'auto').fit(X_train)
          cluster_labels = kmeans.labels_
          cluster_centers = kmeans.cluster_centers_
          ret = train_RBF_trainRBF(X_train, y_train, i, r)
          yhat_test = train_RBF_y_RBF(X_test, ret_)
          yhat_test(yhat_test <= -0.5) = -1
          yhat_test(yhat_test > -0.5) = 1
          yhat_train = train_RBF_y_RBF(X_train, ret_)
          yhat_train(yhat_train <= -0.5) = -1
          yhat_train(yhat_train > -0.5) = 1
          list_acc_test.append(accuracy_score(y_test, yhat_test))
          list_param.append(i)
      list_acc_test = np.array(list_acc_test)
      list_param = np.array(list_param)

      H shape : (210, 1)
      Havg shape : (210, 2)
      yin shape : (210,)
      W shape : (2,)
      H shape : (210, 2)
      Havg shape : (210, 3)
      yin shape : (210,)
      W shape : (3,)
      H shape : (210, 3)
      Havg shape : (210, 4)
      yin shape : (210,)
      W shape : (4,)
      H shape : (210, 4)
      Havg shape : (210, 5)
      yin shape : (210,)
      W shape : (5,)
      H shape : (210, 5)
      Havg shape : (210, 6)

[114]: idx = np.argmax(list_acc_test)
      print(f"The value of K that give us the best performance in test data is : {list_param[idx]}")
      The value of K that give us the best performance in test data is : 33
```



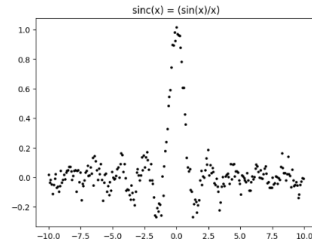
2 Sinc(x)

A figura abaixo retrata os dados gerados por uma função sinc(x) com 1 certo ruído gaussiano. Também mostra o resultado do clustering através do algoritmo KMeans :

```
[1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

[2]: x_sinc = np.arange(start = -10, stop = 10, step = 0.08)
y_sinc = np.sinc(x_sinc)
y_noise = y_sinc + np.random.normal(loc = 0, scale = 0.05, size = 250)

[3]: plt.scatter(x_sinc, y_noise, color = "black", s = 5)
plt.title('sinc(x) = [sin(x)/x]')
plt.show()
```



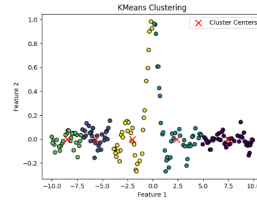
```
[4]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x_sinc, y_noise, random_state = 0, train_size = 0.7)

[5]: from sklearn.cluster import KMeans
X_train = X_train.reshape(-1, 1)
y_train = y_train.reshape(-1, 1)
kmeans = KMeans(n_clusters = 5, random_state = 0, n_init = 'auto').fit(X_train)
cluster_labels = kmeans.labels_
cluster_centers = kmeans.cluster_centers_

[6]: plt.scatter(X_train, y_train, c=cluster_labels, cmap='viridis', marker='o', edgecolor='k')
plt.scatter(cluster_centers[:, 0], np.zeros((cluster_centers.shape[0]),), c='red', marker='x', s=200, label='Cluster Centers')

# Add labels and legend
plt.xlabel('Feature 0')
plt.ylabel('Feature 1')
plt.title('KMeans Clustering')
plt.legend()

# Show the plot
plt.show()
```



O resultado do erro para $K = 5$ neurônios está registrado abaixo :

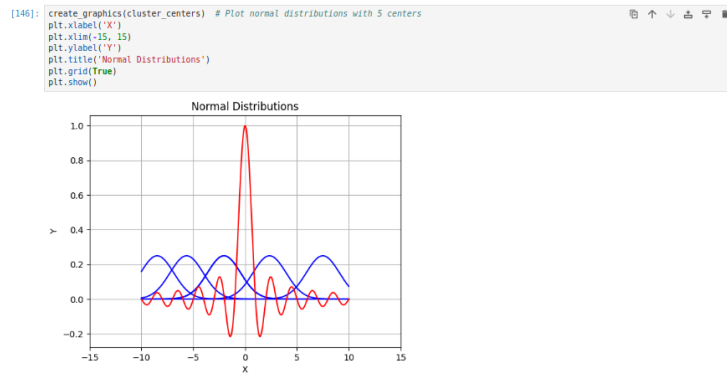
```
[9]: yhat_test = train_RBF.y_RBF(X_test, ret_)
yhat_train = train_RBF.y_RBF(X_train, ret_)
```

You're having the error tuple index out of range. So we will change the number of dimensions...

```
[10]: from sklearn.metrics import mean_squared_error
error_train = mean_squared_error(y_train, yhat_train)
error_test = mean_squared_error(y_test, yhat_test)
print(f"The train error is {error_train}")
print(f"The test error is {error_test}")
```

The train error is 0.0440884446602437
The test error is 0.05120952912089694

O plot abaixo mostra as funções gaussianas utilizadas na camada intermediária, e também a função $\text{sinc}(x)$ que gerou os dados.



É possível observar que, aumentando o número de funções radiais (neurônios da camada intermediária), o modelo tenderá a ter um overfitting. O valor ideal de K, usando o algoritmo de Grid Search foi de 24 neurônios radiais. A imagem abaixo mostra o algoritmo de Grid Search para encontrar o valor de K o qual o modelo terá a melhor performance.

```
[34]: from sklearn.metrics import accuracy_score
lst_error_test = list()
lst_param = list()
for i in range(1, 200, 1):
    kmeans = KMeans(n_clusters = i, random_state = 0, n_init = 'auto').fit(X_train)
    cluster_labels = kmeans.labels_
    cluster_centers = kmeans.cluster_centers_
    ret = train_RBF.trainRBF(X_train, y_train, i, 0.6)
    yhat_test = train_RBF.y_RBF(X_test, ret)

    error_test = mean_squared_error(y_test, yhat_test)
    lst_error_test.append(error_test)
    lst_param.append(i)
lst_error_test = np.array(lst_error_test)
lst_param = np.array(lst_param)

yin shape : (175, 1)
W shape : (5, 1)
You're having the error tuple index out of range. So we will change the number of dimensions...
H shape : (175, 5)
Haug shape : (175, 6)
yin shape : (175, 1)
W shape : (6, 1)
You're having the error tuple index out of range. So we will change the number of dimensions...
H shape : (175, 6)
Haug shape : (175, 7)
yin shape : (175, 1)
W shape : (7, 1)
You're having the error tuple index out of range. So we will change the number of dimensions...
H shape : (175, 7)
Haug shape : (175, 8)
yin shape : (175, 1)
W shape : (8, 1)
You're having the error tuple index out of range. So we will change the number of dimensions...
H shape : (175, 8)

[36]: idx = np.argmin(lst_error_test)
print(f'Using {lst_param[idx]} radial neurons in the hidden layer the model has the better performance.')
Using 24 radial neurons in the hidden layer the model has the better performance.
```