

# Trabalho 02: Aprendizado Não Supervisionado

Universidade Estadual de Campinas (UNICAMP), Instituto de Computação (IC)

Prof. Jacques Wainer, 2021s2

Aluno: Maurício Pereira Lopes - RA: 225242

```
In [1]: 1 # carregamento de bibliotecas
2 import pandas as pd
3 from sklearn.cluster import KMeans
4 from sklearn.mixture import GaussianMixture
5 from sklearn import metrics
6 import numpy as np
7 import matplotlib.pyplot as plt
8
9 # configurações
10 random_state = 42
```

**Leia este dataset que é um dataset de 1000 dados com 13 atributos. Os dados não precisam ser normalizados.**

Fonte do arquivo de dados: <https://www.ic.unicamp.br/~wainer/cursos/2s2021/433/ex2-data.csv>  
(<https://www.ic.unicamp.br/~wainer/cursos/2s2021/433/ex2-data.csv>)

```
In [3]: 1 # Carregar dados como um Pandas DataFrame e converter para um numpy array
2 url = "https://www.ic.unicamp.br/~wainer/cursos/2s2021/433/ex2-data.csv"
3 df = pd.read_csv(url, header = None, sep = r"\s+")
4 data = df.to_numpy()
5 data
```

```
Out[3]: array([[ -10.63,  -3.91,  27.69, ...,  5.32,   1.35,  -4.56],
 [ 12.56, -19.5 ,   4.39, ..., -3.41,  -1.85,   2.55],
 [   4.19, -12.3 , -22.25, ...,   5.5 ,  11.08,  -2.19],
 ...,
 [   9.61, -21.71,  10.19, ...,  -2.   ,   2.94,   7.31],
 [-14.32,  -4.46,  -4.53, ..., -3.52,  -1.13,   5.46],
 [ 17.4 ,  11.46,  14.03, ..., -3.86,  -4.54,   3.02]])
```

## 1- K-Means

- Rode o K-Means com k de 2 a 15

```
In [4]: 1 # Definir uma lista com os valores de k indo de 2 a 15
2 k_list = list(range(2, 16, 1))
3 kmeans = {}
4 for k in k_list:
5     kmeans[k] = KMeans(n_clusters = k, random_state = random_state).fit(data)
```

- Use silhueta e pelo menos alguma outra medida interna de qualidade

```
In [5]: 1 # Silhouette Coefficient
2 silhouette = {}
3 for k in k_list:
4     labels = kmeans[k].labels_
5     silhouette[k] = metrics.silhouette_score(data, labels, metric = 'euclidean')
```

```
In [6]: 1 # Calinski-Harabasz Index
2 calinski = {}
3 for k in k_list:
4     labels = kmeans[k].labels_
5     calinski[k] = metrics.calinski_harabasz_score(data, labels)
```

```
In [7]: 1 # Davies-Bouldin Index
2 davies = {}
3 for k in k_list:
4     labels = kmeans[k].labels_
5     davies[k] = metrics.davies_bouldin_score(data, labels)
```

```
In [8]: 1 # Sum of Squared Errors (SSE)
2 sse = {}
3 for k in k_list:
4     # Inertia: Sum of distances of samples to their closest cluster center
5     sse[k] = kmeans[k].inertia_
```

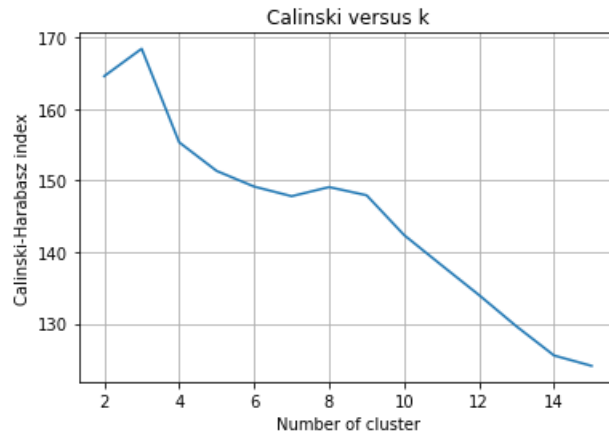
- Mostre o gráfico dessas medidas versus o k

```
In [9]: 1 # Silhouette versus k
2 '''
3 Mais próximo de 1 indica clusterização mais densa.
4 Valores perto de zero indicam sobreposição entre clusters.
5 '''
6 plt.plot(list(silhouette.keys()), list(silhouette.values()), label = 'linear')
7 plt.grid()
8 plt.title("Silhouette versus k")
9 plt.xlabel("Number of cluster")
10 plt.ylabel("Silhouette Coefficient")
11 plt.show()
```



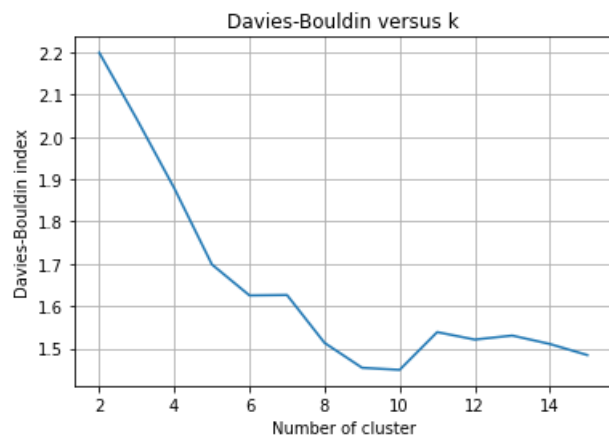
In [10]:

```
1 # Calinski-Harabasz versus k
2 '''
3 O valor é alto quando os clusters são densos e
4 bem separados.
5 '''
6 plt.plot(list(calinski.keys()), list(calinski.values()), label = 'linear')
7 plt.grid()
8 plt.title("Calinski versus k")
9 plt.xlabel("Number of cluster")
10 plt.ylabel("Calinski-Harabasz index")
11 plt.show()
```

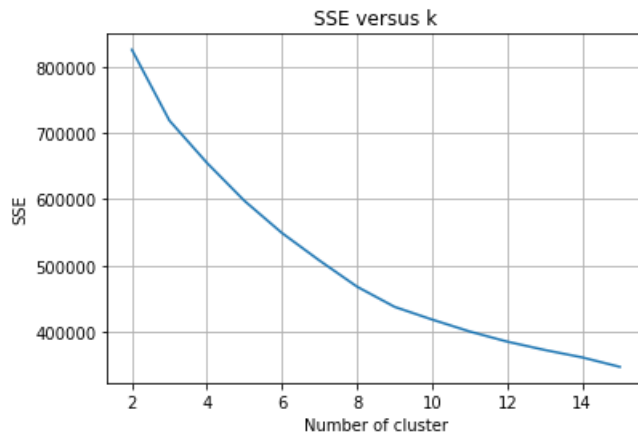


In [11]:

```
1 # Davies-Bouldin versus k
2 '''
3 Zero é o menor valor possível.
4 Valores próximos de zero indicam melhor clusterização.
5 '''
6 plt.figure()
7 plt.plot(list(davies.keys()), list(davies.values()))
8 plt.grid()
9 plt.title("Davies-Bouldin versus k")
10 plt.xlabel("Number of cluster")
11 plt.ylabel("Davies-Bouldin index")
12 plt.show()
```



```
In [12]: 1 # Sum of Squared Errors (SSE) versus k
2 '''
3 Estamos buscando por um cotovelo (elbow) no gráfico,
4 a partir do qual, não haveria ganho significativo na
5 redução da soma dos erros quadrados ao se aumentar o k.
6 '''
7 plt.figure()
8 plt.plot(list(sse.keys()), list(sse.values()))
9 plt.grid()
10 plt.title("SSE versus k")
11 plt.xlabel("Number of cluster")
12 plt.ylabel("SSE")
13 plt.show()
```



## 2 - Escolha um k

Analisando os gráficos das métricas acima, um k razoável poderia ser do 8 ao 10. Nesta feixa vemos o seguinte comportamento nos gráficos:

- Maior coeficiente de silhueta
- Valor intermediário no índice de Calinski-Harabasz
- Valor bem próximo ao mínimo do índice de Davies-Bouldin
- SSE encaminhando para a estabilidade

Estou optando por um k = 9.

```
In [13]: 1 # atribuindo o valor escolhido para k
2 k = 9
3 # obtendo os labels de cada dado conforme a clusterização com o k escolhido
4 labels_kmeans = kmeans[k].labels_
```

## 3 - GMM

Usando o k escolhido no item anterior, rode o GMM

- Com gaussianas esféricas (matrizes de covariancia sao a matriz identidade vezes uma constante)

```
In [14]: 1 labels_gmm_esf = GaussianMixture(n_components = k, covariance_type = "spherical",
2                                     random_state = random_state).fit_predict(data)
```

- Com gaussianas diagonais (matrizes de covariancia sao matrizes diagonais)

```
In [15]: 1 labels_gmm_diag = GaussianMixture(n_components = k, covariance_type = "diag",
2                                       random_state = random_state).fit_predict(data)
```

- Com gaussianas sem restrição (as matrizes de covariancia são livres)

```
In [16]: 1 labels_gmm_full = GaussianMixture(n_components = k, covariance_type = "full",
2                                           random_state = random_state).fit_predict(data)
```

#### 4 - Medidas externas para comparar duas clusterizações

- Use pelo menos 2 medidas externas para comparar a solução do GMM sem restrição (full) com as outras 2 (esférica e diagonal).

```
In [17]: 1 # Rand index
2 rand_full_esf = metrics.rand_score(labels_gmm_full, labels_gmm_esf)
3 rand_full_diag = metrics.rand_score(labels_gmm_full, labels_gmm_diag)
4 print("Rand index entre\nGMM Full e GMM Esférico:", rand_full_esf)
5 print("GMM Full e GMM Diagonal:", rand_full_diag)
```

```
Rand index entre
GMM Full e GMM Esférico: 0.9323103103103103
GMM Full e GMM Diagonal: 0.9464984984984985
```

```
In [18]: 1 # Mutual Information based scores
2 mutual_full_esf = metrics.adjusted_mutual_info_score(labels_gmm_full, labels_gmm_esf)
3 mutual_full_diag = metrics.adjusted_mutual_info_score(labels_gmm_full, labels_gmm_diag)
4 print("Mutual information score entre\nGMM Full e GMM Esférico:", mutual_full_esf)
5 print("GMM Full e GMM Diagonal:", mutual_full_diag)
```

```
Mutual information score entre
GMM Full e GMM Esférico: 0.7873659750308852
GMM Full e GMM Diagonal: 0.7966804927436857
```

```
In [19]: 1 # Fowlkes-Mallows scores
2 fms_full_esf = metrics.fowlkes_mallows_score(labels_gmm_full, labels_gmm_esf)
3 fms_full_diag = metrics.fowlkes_mallows_score(labels_gmm_full, labels_gmm_diag)
4 print("Fowlkes-Mallows score entre\nGMM Full e GMM Esférico:", fms_full_esf)
5 print("GMM Full e GMM Diagonal:", fms_full_diag)
```

```
Fowlkes-Mallows score entre
GMM Full e GMM Esférico: 0.748992634365838
GMM Full e GMM Diagonal: 0.7769945956528681
```

Como uma comparação adicional, não solicitada no enunciado, foi calculado o Rand index entre GMM Full e KMeans com o k escolhido. Pode-se ver que, pelo Rand Index, o KMeans chegou a um resultado muito próximo ao GMM\_Full.

```
In [20]: 1 # Rand index
2 rand_full_kmeans = metrics.rand_score(labels_gmm_full, labels_kmeans)
3 print("Rand index entre\nGMM Full e KMeans:", rand_full_kmeans)
```

```
Rand index entre
GMM Full e KMeans: 0.9461161161161161
```

Em todas as métricas externas calculadas acima, os melhores resultados foram entre GMM\_Full e GMM\_Diagonal. Isso indica que estes dois modelos geraram clueterizações mais parecidas. Se os labels obtidos com o GMM\_Full fossem considerados como Ground Truth, o GMM\_Diag teria sido o modelo com melhor resultado de clusterização.

```
In [ ]: 1
```