



# DATA ANALYSIS AND CLASSIFICATION

## TP1

ANA CAROLINA MORAIS N°2021222056

EDUARDO FERREIRA N°2021218018

# ÍNDICE

<b>01</b>	<b>Introdução</b>	
	Introdução .....	03
	Objetivo .....	03
	Base de dados .....	03
<b>02</b>	<b>Tarefa 1</b>	
	Importação de dados .....	04
	ICA .....	04
	Observações .....	05
<b>03</b>	<b>Tarefa 2 e 3</b>	
	Windowing .....	06
	Spectrum .....	06
	Observações .....	06
<b>04</b>	<b>Tarefa 4</b>	
	Filtros .....	07
	FIR .....	07
	IIR .....	07
	Observações .....	07
<b>05</b>	<b>Tarefa 6</b>	
	Wavelet .....	08
	CWT .....	08
	DWT .....	08

# ÍNDICE

<b>06</b>	<b>Tarefa 7</b>	
	Redução de dimensionalidade .....	09
	PCA .....	09
	MDS .....	10
<b>07</b>	<b>Tarefa 8</b>	
	Clustering .....	12
	Resultados .....	13
<b>08</b>	<b>Conclusão</b>	
	Conclusão .....	14
	Referências .....	15

# 01|Introdução

A detecção automática de emoções humanas através da atividade cerebral tem aplicações valiosas em diversas áreas, como ‘interfaces’ cérebro-computador, terapia e neurociência. Os sinais de eletroencefalograma (*EEG*) são uma ferramenta poderosa para capturar estados emocionais, mas apresentam desafios significativos devido à sua complexidade e suscetibilidade a ruído. Para superar essas dificuldades, torna-se essencial um processamento cuidadoso e a aplicação de técnicas avançadas de análise de dados. Neste trabalho, propomos implementar um *pipeline* de classificação de emoções com base em gravações brutas de *EEG*, recorrendo a métodos como a Transformada *Wavelet*, Redução de Dimensionalidade e *Clustering*.

## 01|Objetivo

O objetivo deste projeto é desenvolver e implementar um *pipeline* de análise de dados em *MATLAB* para classificar emoções humanas a partir de sinais *EEG*. Neste contexto, trabalhamos especificamente com os dados brutos do participante número 2, abrangendo três sessões experimentais distintas. As etapas incluem o carregamento e pré-processamento dos sinais, extração de características espectrais, redução de dimensionalidade e aplicação de algoritmos de clustering para distinguir os estados emocionais de felicidade, neutralidade e tristeza. A avaliação do desempenho do modelo será realizada através da construção de uma matriz de confusão e métricas adequadas. É, possível ver através da fig. 1, um esquema do que pretendemos implementar.



Fig.1. Pipeline

## 01|Base de dados

O *dataset* utilizado é o SEED (SJTU Emotion EEG Dataset), que contém registros de EEG de 15 participantes, capturados enquanto assistiam a 15 vídeos de filmes projetados para induzir três emoções: feliz, neutro e triste. Cada participante passou por três sessões experimentais em dias diferentes. Não esquecendo que o participante utilizado neste projeto é o n.º2.

Os dados estão armazenados em ficheiros *.cnt*, contendo os sinais EEG de 62 canais, adquiridos pelo sistema 10-20. Há também arquivos auxiliares com a sequência dos vídeos e marcações de tempo.

## 02| Tarefa 1 | Importação dos dados e Tratamento

Os sinais EEG foram carregados a partir dos ficheiros *.cnt* utilizando a função *loadcnt.m*. Em seguida, foram excluídos os canais não correspondentes a EEG (VEO, HEO, M1, M2), mantendo somente os 62 canais relevantes. Para otimizar o processamento, a taxa de amostragem foi reduzida por um fator de 4 usando a função *downsample*. A segmentação dos dados foi realizada com base nos tempos de início e fim de cada ‘trial’, ajustados para a nova taxa de amostragem.

Durante o pré-processamento, foram identificados e removidos alguns ‘trials’ que apresentavam picos anómalos. Essa decisão foi tomada para preservar a integridade dos dados, por ser preferível eliminar um ‘trial’ inteiro do que comprometer um canal específico. Em diferentes experiências, um canal pode estar comprometido num ‘trial’, mas ser válido noutra. Na fig. 2 podemos observar um ‘trial’ que eliminámos para a experiência 2\_1 o ‘trial’ n.º 4, onde é possível observar uma variação brusca de amplitude em um ou mais canais, indicando a presença de artefactos que poderiam afetar a análise posterior. Para a experiência 2\_2 eliminámos o ‘trial’ n.º10 e para a experiência 2\_3 eliminámos o ‘trial’ n.º13, todos apresentam bastante semelhança entre eles, daí serem excluídos.

Por fim, para visualização dos sinais, foram gerados gráficos dos 62 canais distribuídos num ‘grid’ de  $8 \times 8$  subplots, onde cada canal é plotado ao longo do tempo, permitindo uma análise detalhada da evolução dos sinais EEG em cada ‘trial’. Conforme a fig. 3 para o ‘trial’ n.º1, da experiência 2\_1.

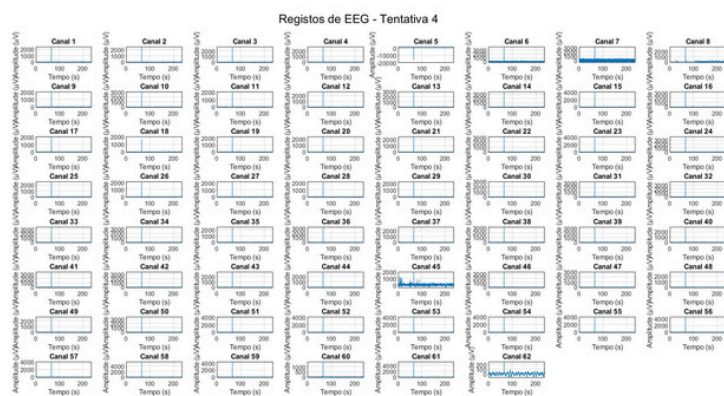


Fig.2. Representação de um ‘trial’ que foi eliminado



Fig.3. Representação dos 62 canais para uma dada tentativa

## 02| Tarefa 1 | ICA

Para minimizar artefactos e melhorar a qualidade do sinal EEG, aplicámos a técnica de Análise de Componentes Independentes (ICA). O ICA é um método amplamente utilizado para separar sinais misturados em componentes independentes, permitindo identificar e remover fontes de ruído, como piscadelas oculares, movimentos musculares e interferências elétricas.

A nossa implementação utilizou a função *apply\_ICA*, onde o algoritmo *FastICA* foi aplicado individualmente a cada ‘trial’. Como estratégia de otimização, testámos duas abordagens distintas de extração de componentes: *kurtosis* e *negentropy*. A escolha da melhor decomposição para cada ‘trial’ foi feita com base na variância das componentes obtidas, seleccionando o método que apresentou maior dispersão média das componentes independentes, pois isso indica uma melhor separação dos sinais de origem, o que neste caso o melhor método foi o *negentropy*.

Após a extração das componentes independentes, cada ‘trial’ foi inspecionado visualmente para identificar componentes indesejadas associadas a artefactos. As componentes consideradas não fisiológicas foram definidas como zero antes de reconstruir o sinal EEG limpo. Devido à complexidade do sinal e à persistência de artefactos, foi necessário realizar duas iterações do ICA para melhorar a separação. No entanto, mesmo após várias iterações, não foi possível identificar todas as componentes indesejadas em alguns casos.

A reconstrução do EEG limpo foi realizada utilizando a transformação inversa do ICA:

$$EEG\_Reconstruído = T * Zica + m$$

onde T representa a matriz de mistura, Zica as componentes independentes (já filtradas), e m, a média original do sinal.

## 02 | Tarefa 1 | Observações

Após a aplicação do ICA, observa-se uma redução significativa de artefactos nos sinais de EEG, especialmente aqueles relacionados a piscadelas e atividade muscular. A amplitude das ondas diminuiu em vários canais, indicando a remoção de componentes de alta amplitude, e o ruído de alta frequência foi atenuado, resultando num sinal mais limpo. No entanto, em certos canais, a remoção pode ter sido excessiva, resultando na possível perda de informações neurais relevantes. Mesmo após duas iterações de ICA, alguns artefactos ainda persistem, o que sugere que o ICA, por si só, pode não ser suficiente para uma separação completa das fontes de ruído. Para lidar com esses artefactos remanescentes, uma abordagem complementar pode envolver o uso de técnicas de filtragem, o que iremos abordar na próxima tarefa. Assim, em vez de depender exclusivamente da seleção e remoção manual de componentes, a combinação do ICA com outras técnicas pode oferecer um equilíbrio mais eficaz entre a limpeza do sinal e a preservação das informações neurais relevantes.

### 03|Tarefa 2 & 3 | Windowing | Spectrum

Para evitar artefactos espectrais causados pela truncagem dos sinais, aplicaram-se diferentes funções de janela a cada segmento reconstruído do EEG. As funções de janela escolhidas foram *Hamming*, *Blackman* e *Hann* e *Triangular*, tendo como objetivo suavizar as extremidades dos sinais e minimizar a dispersão espectral. A multiplicação de cada segmento pela respetiva função de janela foi realizada antes da análise no domínio da frequência.

A comparação do espectro de amplitude obtido para cada janela foi realizada através da Transformada Rápida de Fourier (FFT). Como esperado, a maioria da energia do sinal EEG concentrou-se nas frequências mais baixas, até cerca de 1 - 100 Hz, correspondendo às bandas de frequência mais relevantes na atividade cerebral. Observou-se que todas as funções de janela produziram espectros semelhantes, sugerindo que a escolha da janela não introduziu diferenças significativas na análise espectral para este tipo de sinal. (Fig. 4 - Comparação dos Espectros para as diferentes janelas).

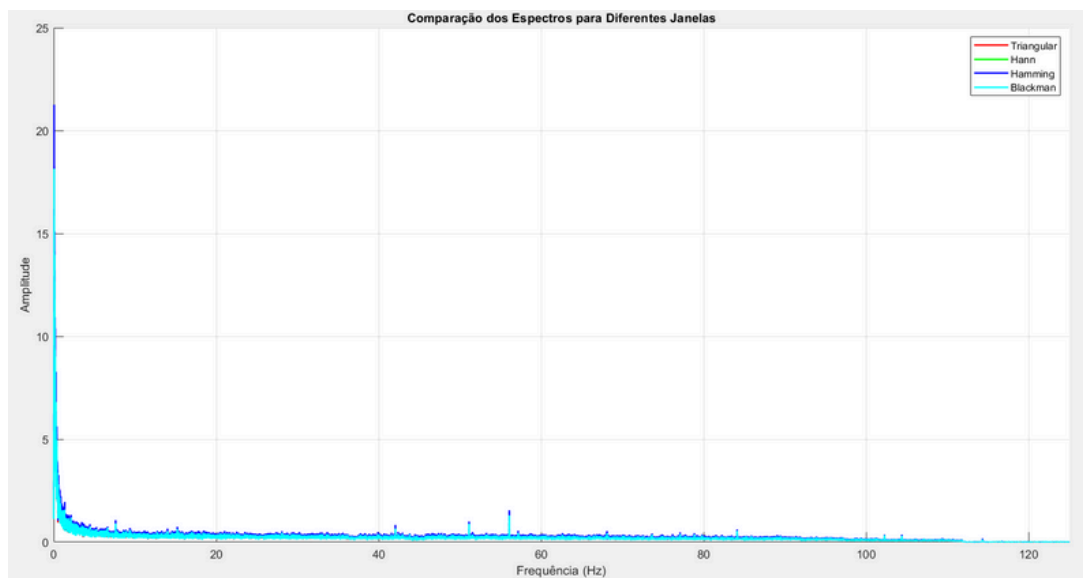


Fig.4. Comparação dos Espectros para Diferentes Janelas

A janela *Hamming* demonstrou um compromisso entre resolução espectral e atenuação dos lóbulos secundários, enquanto a janela *Blackman* apresentou a melhor supressão dos lóbulos secundários, embora à custa de uma menor resolução devido ao alargamento do lóbulo principal. A janela *Hann* revelou-se uma solução intermédia, com melhor atenuação dos lóbulos secundários em comparação com a *Hamming*, mas sem alcançar a eficácia da *Blackman* na redução da dispersão espectral. Já a janela *Triangular*, apesar de ser conceitualmente mais simples, demonstrou um desempenho próximo ao da *Hann*, mas com uma atenuação ligeiramente inferior dos lóbulos secundários.

Apesar das diferenças teóricas entre as janelas, os resultados práticos não evidenciaram variações expressivas, indicando que o impacto da função de janela foi reduzido para este conjunto de dados. No entanto, entre as opções testadas, a janela *Hann* mostrou-se a mais equilibrada, proporcionando um bom compromisso entre resolução espectral e minimização da dispersão, tornando-se, assim, a escolha mais recomendável para a análise espectral destes sinais EEG.

## 04|Tarefa 4 | Filtros

Na etapa de filtragem, foram testadas abordagens baseadas em filtros de Resposta Finita ao Impulso (FIR) e Resposta Infinita ao Impulso (IIR) para remover componentes espectrais fora da faixa de interesse (1 Hz - 100 Hz). No caso dos filtros FIR, foram exploradas duas configurações: um filtro passa-banda (Band-Pass) e uma combinação de filtros passa-baixa (Low-Pass) e passa-alta (High-Pass). O filtro passa-banda mostrou-se eficaz na remoção das frequências indesejadas, mas a abordagem baseada na aplicação sequencial de filtros passa-baixa e passa-alta permitiu um controle mais preciso sobre as transições de frequência, reduzindo oscilações na resposta espectral.

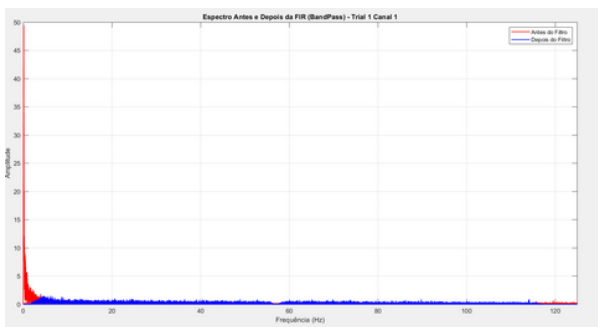


Fig. 5 - FIR - Aplicação do filtro Bandpass

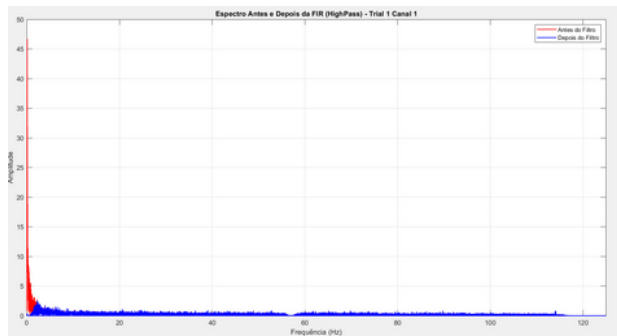


Fig. 6 - FIR -Aplicação do filtro passa-baixo seguido de um passa-alto

Relativamente aos filtros IIR, optou-se por utilizar filtros Butterworth, tanto na configuração passa-baixa quanto passa-alta. Esses filtros foram escolhidos devido à sua resposta suave e ausência de ondulações na banda passante, garantindo uma transição gradual entre as frequências eliminadas e preservadas. No entanto, apesar de atingirem a atenuação desejada com uma ordem significativamente menor do que os FIR, os filtros IIR introduziram distorção de fase, um fator crítico na análise de sinais EEG, onde a preservação da integridade temporal das oscilações é essencial para uma correta interpretação dos dados.

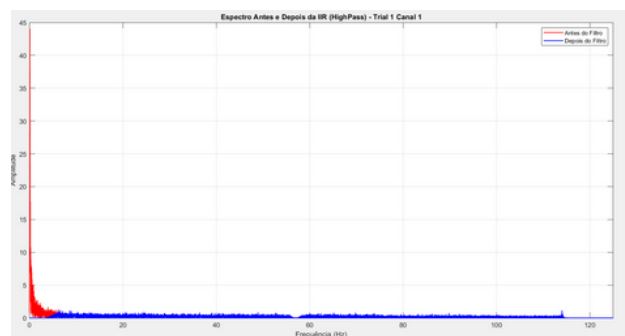


Fig. 7 - IIR -Aplicação do filtro passa-baixo seguido de um passa-alto

Diante dessa análise, a escolha final recaiu sobre os filtros FIR na configuração passa-baixa e passa-alta, pois, embora exigissem uma ordem mais elevada e, consequentemente, maior custo computacional, garantiram estabilidade e uma resposta de fase linear. Esse compromisso entre eficiência espectral e preservação da forma de onda assegurou que os sinais EEG processados mantivessem a sua fidelidade estrutural, sem distorções que pudessem comprometer etapas posteriores da análise.



## 05|Tarefa 6 | Wavelet

Numa fase inicial, procurou-se identificar os canais que poderiam ter maior relevância para a análise. Para reduzir a carga de trabalho na avaliação dos resultados das transformadas de wavelet, selecionaram-se, com base no artigo [1], os canais considerados críticos na interpretação da resposta emocional. Estes canais estão mais associados às frequências beta e gama, tendo sido escolhidos os canais T7 e T8, uma vez que integram o perfil de elétrodos mais simples mencionados no referido artigo.

Considerando um dos ‘trials’ do Participante 2, que daqui em diante será referido como ‘Trial’ 14, observa-se que, ao comparar os espectros dos dois canais, presentes na Figura 8, a magnitude do canal T8 (índice 32) é significativamente inferior à do canal T7 (índice 24). Este resultado pode indicar uma maior atividade na região de T7 em resposta ao estímulo apresentado no ‘Trial’ 14.

No que diz respeito à distribuição de energia, verifica-se que ambos os canais apresentam baixa energia nas frequências mais reduzidas, o que era esperado. Isto ocorre porque as bandas delta e teta, associadas a essas frequências, não são consideradas críticas para a resposta emocional, conforme indicado no artigo previamente referenciado.

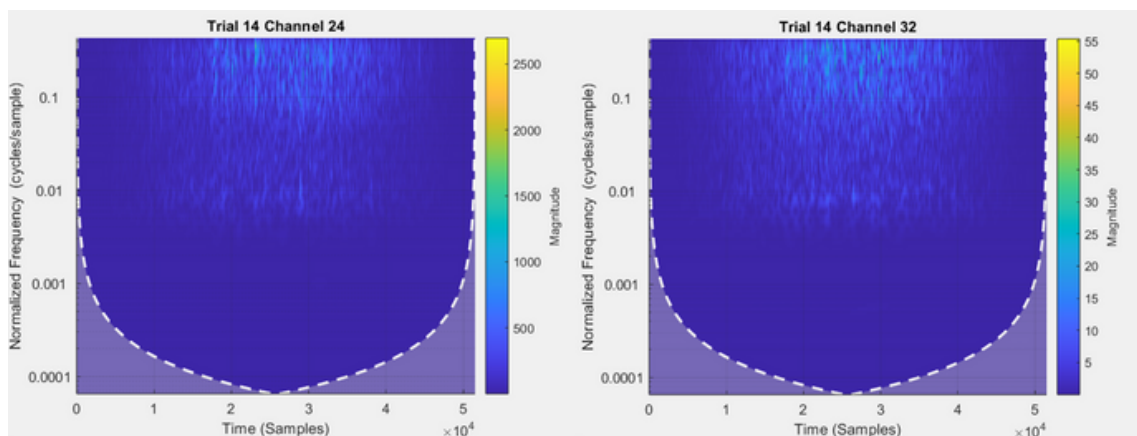


Fig. 8 - Comparação de espectros para diferentes canais

Na análise através da Transformada Wavelet Discreta (DWT), o primeiro passo foi determinar o número adequado de níveis de decomposição. Esta escolha é essencial para garantir que as bandas de frequência de interesse sejam devidamente separadas e analisadas. Tendo em conta que a frequência de amostragem do sinal é de 250 Hz, a frequência máxima representada será de 125 Hz, conforme estabelecido pelo teorema de Nyquist. Em cada nível de decomposição da DWT, a frequência é sucessivamente dividida por 2, resultando em sub-bandas de espectro progressivamente mais estreitas.

Para uma análise eficaz das diferentes bandas de frequência do EEG, é necessário que a decomposição atinja uma resolução suficiente para isolar as bandas delta (1-4 Hz) e teta (4-8 Hz). Considerando esse critério, definiu-se que cinco níveis de decomposição seriam apropriados, por permitirem uma segmentação adequada dessas bandas, facilitando a extração de características relevantes para a análise da resposta emocional.

É relevante salientar que, considerando os intervalos resultantes das sucessivas divisões, para determinar a potência de cada banda, é necessário, no caso específico da banda Gama (30-100 Hz), somar ambas as subdivisões de maior frequência, uma vez que ambas englobam esta faixa de frequência.

## 06|Tarefa 7 | Redução de Dimensionalidade

A redução de dimensionalidade é uma técnica essencial na análise de dados, especialmente quando se trabalha com conjuntos de dados de alta dimensão. O objetivo principal é diminuir o número de variáveis, preservando a maior quantidade possível da variabilidade dos dados originais. Essa abordagem não só reduz o custo computacional, como também facilita a visualização dos dados e pode auxiliar na remoção de ruído.

- **Análise de Componentes Principais (PCA)**

O PCA é um dos métodos mais utilizados para a redução de dimensionalidade. Ele transforma um conjunto de variáveis correlacionadas num novo conjunto de variáveis não correlacionadas, denominadas componentes principais, ordenadas por ordem de variância explicada. O objetivo é projetar os dados num espaço de menor dimensão enquanto se preserva a maior quantidade possível da informação dos dados originais.

No projeto, o PCA foi aplicado à matriz *all\_power*, que representa *features* extraídas de 62 canais ao longo de três ‘trials’, resultando num total de 42 observações. Diferentes técnicas de normalização foram testadas, sendo que a normalização “*medianiqr*” apresentou os melhores resultados por equilibrar a escala das variáveis e reduzir o impacto de outliers. Para a decomposição, foi utilizado o algoritmo *SVD*, e outros métodos, como *EIG* e *ALS*, foram testados, produzindo resultados equivalentes.

Os resultados do PCA foram analisados por via três tipos de visualizações.

1. **Variância Explicada Cumulativa** – Indica quantos componentes são necessários para atingir um determinado limiar de variância preservada. Utilizamos um critério de 95% de variância explicada para selecionar o número ótimo de componentes.
2. **Valores Próprios das Componentes** – Avaliam a relevância de cada componente principal.
3. **Projeção dos Dados** – Representação dos dados no espaço das primeiras duas ou três componentes principais, permitindo uma visualização mais intuitiva da estrutura dos dados.

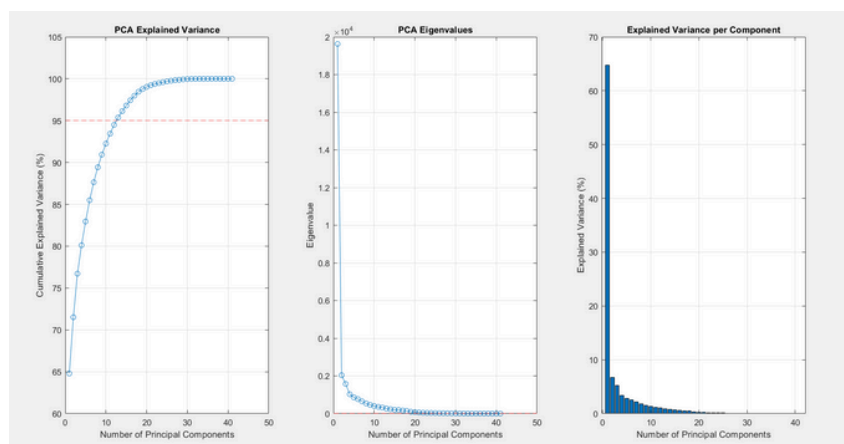


Fig. 9 - Análise da Variância Explicada pelo PCA e os valores próprios (eigenvalues) associados a cada componente principal

## 06|Tarefa 7 | Redução de Dimensionalidade

- **Comparação entre PCA com 2 dimensões e PCA com 3 Dimensões**

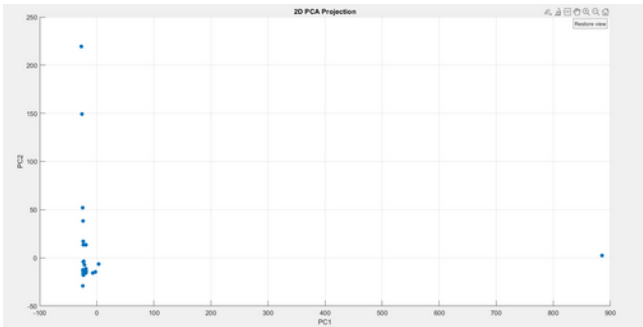


Fig. 10 - Projeção dos dados em 2D

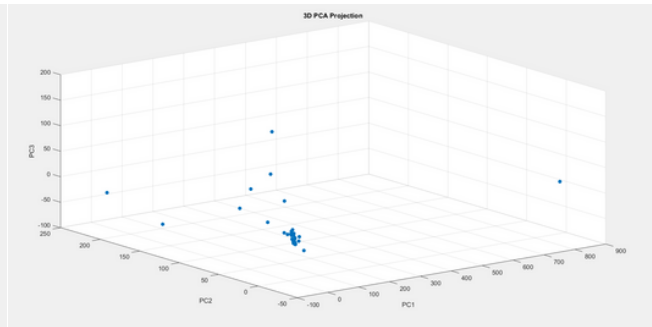


Fig. 11 - Projeção dos dados em 3D

A projeção bidimensional dos dados (Fig. 9), considerando as duas primeiras componentes principais (PC1 e PC2), possibilitou uma visualização mais clara da estrutura dos dados, permitindo observar agrupamentos e identificar tendências de variação. Contudo, algumas observações ainda estavam sobrepostas, o que indicou que a inclusão de uma terceira dimensão poderia contribuir para uma melhor separação entre os dados. A variação explicada pelas duas primeiras componentes principais foi considerável, mas não atingiu a marca dos 95% de variância explicada, sugerindo que componentes adicionais seriam necessárias para capturar toda a informação presente no conjunto de dados.

Ao incorporar uma terceira componente principal na projeção tridimensional (PC1, PC2 e PC3 - Fig. 10), foi possível evidenciar mais claramente algumas observações que estavam sobrepostas na projeção 2D. A distribuição dos dados tornou-se mais distinta, indicando que a terceira componente ajudou a capturar informações adicionais e relevantes.

A análise dos gráficos sugere que a projeção tridimensional proporciona uma separação mais eficaz entre algumas observações, indicando que ao menos três componentes principais são necessárias para capturar adequadamente a estrutura do conjunto de dados. Caso a prioridade seja a simplicidade e a facilidade de interpretação visual, a projeção 2D pode ser suficiente para a compreensão dos dados.

- **Análise Multidimensional Escalonada (MDS)**

O MDS é uma técnica de redução de dimensionalidade que visa representar dados de alta dimensionalidade num espaço de menor dimensão, preservando as distâncias ou semelhanças entre as observações. O *'stress'* é uma métrica fundamental para avaliar o quão bem essa representação foi alcançada, sendo valores mais baixos indicativos de uma melhor correspondência entre as distâncias originais e as representadas. Ao contrário do PCA, que foca na maximização da variância explicada, o MDS busca manter a estrutura de proximidade dos dados originais, tornando-o útil para identificar padrões e agrupamentos complexos em dados multidimensionais.

No nosso projeto, abordamos o problema do MDS aplicando tanto o método métrico quanto o não métrico, com o objetivo de avaliar como cada abordagem influencia a representação dos dados em duas (2D) e três dimensões (3D). Além disso, testamos diferentes métodos de normalização, como *zscore*, *norm*, *center* e *range*, para verificar como cada técnica afeta a distribuição e a interpretação dos dados. Após normalizar os dados conforme o método escolhido, calculamos as distâncias entre as observações usando métricas como a *euclidiana* e a *cityblock*. A seguir, projetamos os dados nas duas e três dimensões, avaliando a qualidade das projeções por meio dos valores de ‘stress’ gerados para cada combinação de método de normalização, métrica de distância e critério (métrico ou não métrico).

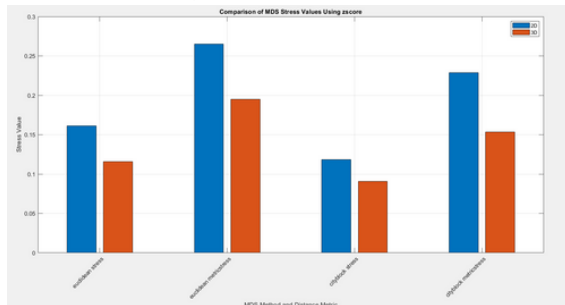


Fig. 12 - Comparação MDS usando Z-score

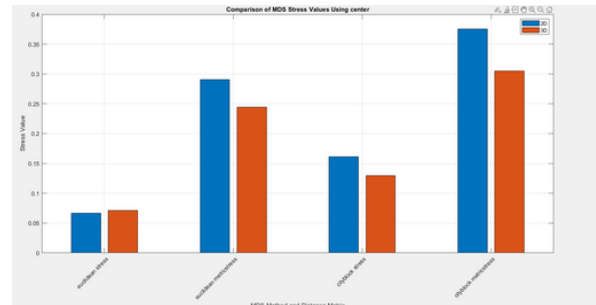


Fig. 13 - Comparação MDS usando Center

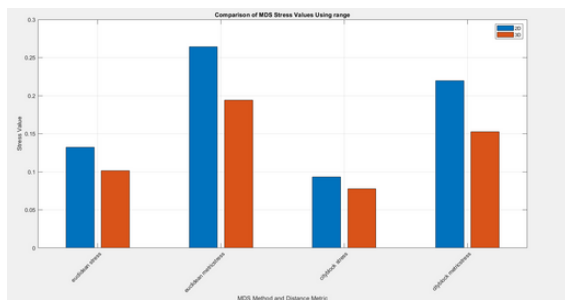


Fig. 14 - Comparação MDS usando Range

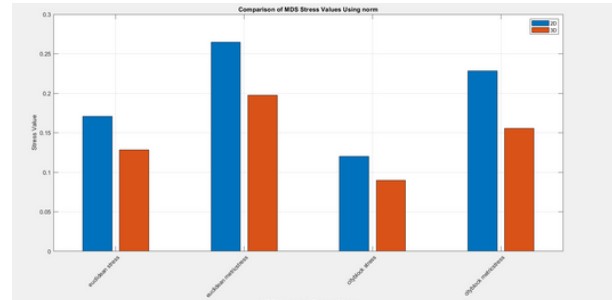


Fig. 15 - Comparação MDS usando Norm

Os resultados mostraram diferenças significativas entre os métodos métrico e não métrico, bem como entre as métricas de distância utilizadas.

- **MDS Métrico vs. Não Métrico:** O método não métrico apresentou consistentemente valores de ‘stress’ mais baixos do que o método métrico em quase todas as combinações de normalização e métrica de distância. Isto sugere que o MDS não métrico é mais eficaz em preservar a estrutura de proximidade dos dados quando esta não segue uma escala linear.
- **Métricas de Distância:** Entre as métricas de distância testadas, a distância euclidiana tendeu a produzir valores de ‘stress’ mais baixos em comparação com a distância cityblock. No entanto, houve casos em que a distância cityblock apresentou valores de ‘stress’ mais baixos, especialmente com determinadas normalizações. Este comportamento indica que a escolha da métrica de distância ideal pode depender do tipo de normalização aplicada, visto que a cityblock tende a enfatizar diferenças locais, enquanto a euclidiana capta melhor as relações globais.

- **Normalizações:**

- *Range* e *Norm* foram as normalizações que melhor se adaptaram ao MDS, produzindo valores de ‘stress’ consistentemente baixos, especialmente quando combinadas com o MDS não métrico.
- *Z-score* também apresentou um bom desempenho, embora ligeiramente inferior às duas anteriores.
- *Center* foi a normalização com o pior desempenho, resultando em valores de ‘stress’ mais altos, independentemente da métrica de distância ou do método de MDS utilizado.

Os resultados indicam que o MDS não métrico é superior ao métrico para preservar as relações de distância em dados complexos, especialmente quando combinado com as normalizações *Range* e *Norm*. Embora a métrica euclidiana tenda a apresentar melhores resultados em geral, a *cityblock* mostrou-se competitiva em certas situações específicas, sugerindo a importância de considerar ambas as métricas dependendo do contexto do problema. Por outro lado, a normalização *Center* mostrou-se inadequada, possivelmente devido à sua incapacidade de lidar com distribuições de dados que não estejam centradas em zero. Assim, conclui-se que para a próxima tarefa, escolhemos a normalização através do método *Norm*, o método não métrico de MDS, “Stress”, e a distância *cityblock*.

## 07|Tarefa 8 | Clustering

O *clustering* é uma técnica de aprendizagem não supervisionada cujo objetivo é agrupar dados com características semelhantes. No contexto deste estudo, aplicamos métodos de clustering para distinguir três estados emocionais: feliz, neutro e triste. A eficácia dos algoritmos de *clustering* foi avaliada com base na matriz de confusão e em métricas como *accuracy*, sensibilidade e especificidade.

Para garantir que os ‘clusters’ fossem minimamente coerentes e representativos, realizámos previamente uma etapa de remoção de *outliers*. Para tal, aplicámos um método baseado na média e no desvio padrão das variáveis, eliminando valores que se desviavam significativamente do padrão esperado. A remoção de outliers foi essencial para reduzir o impacto de observações extremas e permitir que os algoritmos de *clustering* captassem com maior precisão as estruturas subjacentes nos dados. Após esta etapa, testámos quatro métodos de clustering amplamente utilizados: K-Means, K-Medoids, Clustering Hierárquico e Fuzzy C-Means.

- **K-Means:** método baseado na minimização da variância intra-cluster, onde os dados são atribuídos ao ‘cluster’ cujo centroide está mais próximo. Apesar da sua eficiência computacional, é sensível a outliers.
- **K-Medoids:** semelhante ao K-Means, mas utiliza elementos reais dos dados como centros dos ‘clusters’, tornando-o mais robusto a outliers.

- **Clustering Hierárquico (Agglomerative):** método que constrói uma hierarquia de ‘clusters’ com base nas distâncias entre observações, permitindo uma representação estruturada dos dados.
- **Fuzzy C-Means:** abordagem baseada na teoria dos conjuntos difusos, onde uma mesma observação pode pertencer a múltiplos ‘clusters’ com diferentes graus de pertinência, proporcionando uma classificação mais flexível dos dados.

Para avaliar adequadamente o desempenho dos algoritmos, implementamos um método de realinhamento de ‘clusters’ que associa os ‘clusters’ descobertos às classes verdadeiras. O algoritmo prioriza ‘clusters’ maiores e maximiza a sobreposição entre ‘clusters’ previstos e classes reais. Após o realinhamento, calculamos as seguintes métricas:

- **Accuracy:** Proporção de instâncias corretamente classificadas
- **Sensitivity (Sensibilidade):** Média das taxas de verdadeiros positivos para cada classe
- **Specificity (Especificidade):** Média das taxas de verdadeiros negativos para cada classe

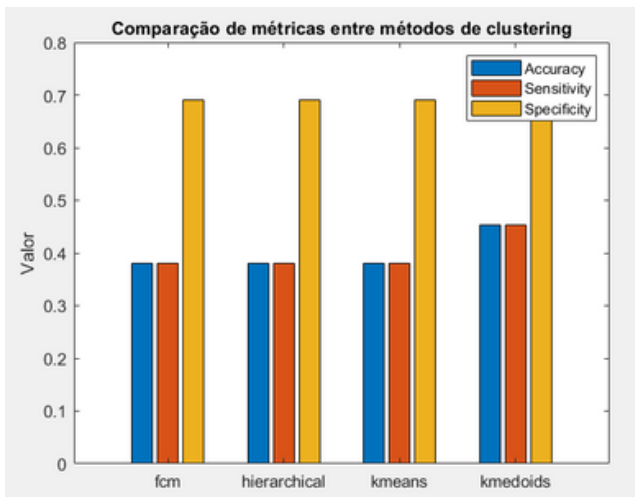


Fig. 16 - Comparação de métricas para diferentes métodos de Clustering - PCA 2D

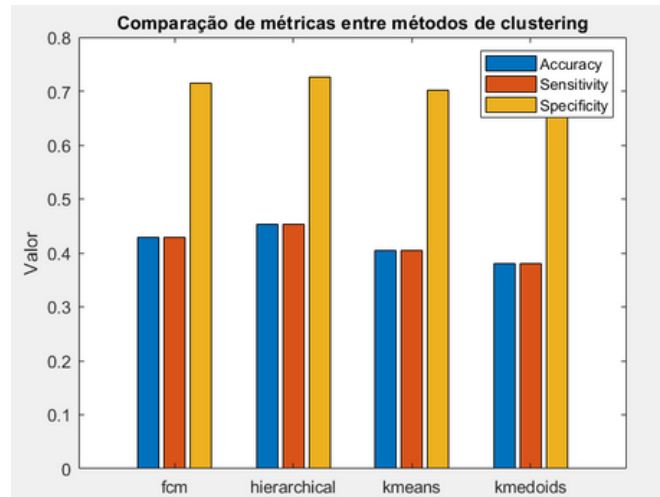


Fig. 17 - Comparação de métricas para diferentes métodos de Clustering - PCA 3D

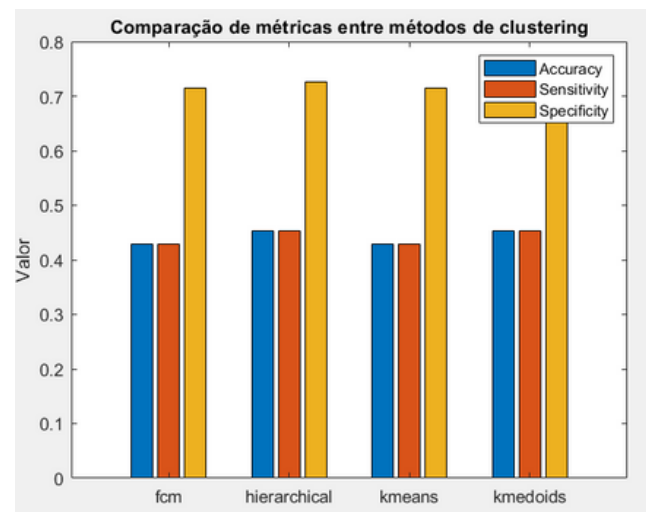


Fig. 18 - Comparação de métricas para diferentes métodos de Clustering - MDS 2D

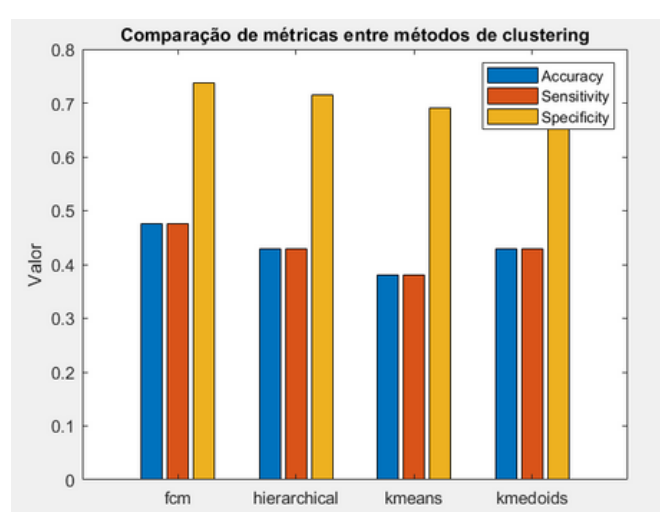


Fig. 19 - Comparação de métricas para diferentes métodos de Clustering - MDS 3D

Os resultados indicam que a escolha da redução de dimensionalidade e do método de *clustering* teve um impacto significativo no desempenho. Surpreendentemente, o Fuzzy C-Means com MDS em 3D obteve o melhor desempenho geral (*accuracy*: 0,48). O K-Medoids e o Clustering Hierárquico demonstraram consistência em diferentes configurações, com bom desempenho em PCA 2D, MDS 2D e, no caso do Hierárquico, também em PCA 3D. O K-Means não se destacou em nenhuma configuração específica, possivelmente devido à sua conhecida sensibilidade a outliers, apesar da tentativa de remoção.

A técnica MDS tendeu a proporcionar melhores resultados que o PCA, sugerindo que a preservação das distâncias entre os pontos é mais relevante para esta tarefa específica.

Notavelmente, todos os métodos apresentaram especificidade (0,69-0,74) consideravelmente mais alta que *accuracy* e sensibilidade (0,38-0,48), indicando boa capacidade de identificar verdadeiros negativos.

Em resumo, embora os valores de *accuracy* para todos os métodos sejam moderados, a combinação de MDS em 3D com Fuzzy C-Means revelou-se a mais eficaz para diferenciar os três estados emocionais, seguida pelas combinações de K-Medoids e Clustering Hierárquico com MDS 2D. Estes resultados destacam a importância de testar diferentes algoritmos e técnicas de redução de dimensionalidade para otimizar a detecção de padrões em dados emocionais.

## 08| Conclusão

Ao analisarmos os resultados do *clustering*, constatamos que a divisão dos pontos não reflete realistamente os valores reais, conforme evidenciado pelos baixos valores de precisão (*accuracy*). Esta limitação pode ser atribuída a diversos fatores. No entanto, considerando que foram extensivamente testadas várias opções em todas as etapas do procedimento — nomeadamente o tipo de janela aplicada, a *wavelet* selecionada, o método de normalização e o algoritmo de redução de dimensionalidade —, é provável que a reduzida precisão esteja relacionada com a fase de pré-processamento.

Em particular, existe a possibilidade de ter sido removida uma quantidade excessiva de informação, especialmente na aplicação do ICA. Nesta etapa, a separação e remoção de componentes ruidosas dos dados relevantes não foi realizada idealmente, o que pode ter resultado, mais uma vez, na eliminação excessiva de informação, bem como na falha do algoritmo em identificar determinados padrões ruidosos.

Com o objetivo de avaliar a eficácia do ICA por nós implementado, procedemos à aplicação de todos os passos do processo, excluindo especificamente o ICA, e analisámos os resultados obtidos. Observámos diferenças significativas, sobretudo após a aplicação do PCA.

## 08| Conclusão

Neste cenário, os pontos resultantes apresentam-se maioritariamente sobrepostos em regiões específicas, conduzindo, inevitavelmente, a baixas precisões.

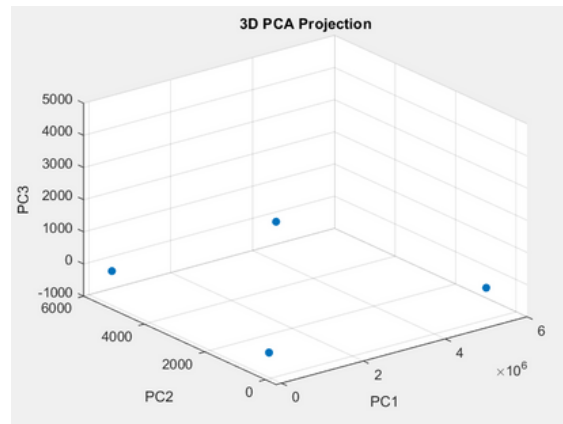


Fig. 20 - Comparação de métricas para diferentes métodos de Clustering - MDS 2D

Na aplicação do MDS, verificou-se uma distribuição dos pontos mais próxima do esperado. No entanto, as precisões obtidas foram ainda inferiores.

Com base nestes resultados, conclui-se que, apesar de a aplicação do ICA não ter sido executada de forma absolutamente ideal, esta etapa revela-se essencial para a remoção do ruído do sinal a analisar.

Em conclusão, embora os resultados não sejam particularmente encorajadores, foi possível aplicar e testar o procedimento de forma sistemática, maximizando, dentro do possível, a sua aproximação à realidade.

## | Referências

[1] Zheng, Wei-Long & Lu, Bao-Liang. (2015). Investigating Critical Frequency Bands and Channels for EEG-Based Emotion Recognition with Deep Neural Networks. IEEE Transactions on Autonomous Mental Development. 7. 1-1. 10.1109/TAMD.2015.2431497.