

TDE3 – Busca Local

Neste projeto, você e seu grupo deverão implementar um **algoritmo genético (GA)** para o processo de seleção de features em um problema de classificação de dígitos manuscritos (de 0 a 9), considerando o conjunto de dados MNIST¹.

Cada exemplo desse conjunto de dados é composto por uma imagem de dimensões 28x28, conforme ilustrado na Figura 1. O dataset possui 60 mil exemplos de treinamento e 10 mil exemplos de teste, cada um com 784 features. Assim, o treinamento de um modelo com um subconjunto reduzido e adequado de features, será computacionalmente mais eficiente, terá ganhos de interpretabilidade e poderá ter resultados preditivos superiores.

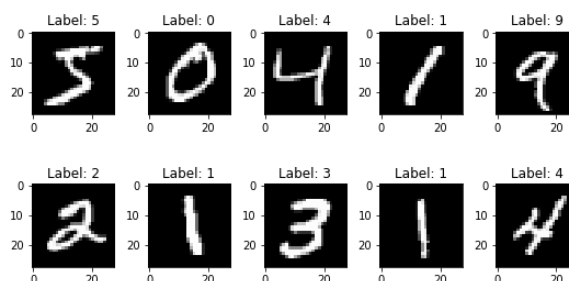


Figura 1. Exemplos do conjunto de dados MNIST.

Os resultados obtidos com a sua implementação de GA para a seleção de features deverão ser comparados com os resultados de outras implementações, também de sua responsabilidade. Em específico, deve-se considerar:

- Abordagem **wrapper**, podendo ser *backward selection* ou *forward selection*.
- **Baseline**, considerando um modelo com 100% das features disponíveis.

As três abordagens (GA, wrapper e baseline) devem ser comparadas em termos de:

- Acurácia no conjunto de teste;
- Porcentagem de features selecionadas (em relação aos dados originais);
- Tempo de treinamento;
- Tempo necessário para encontrar o subconjunto de features.

Para uma comparação justa de resultados, você deverá considerar o mesmo modelo de classificação (**Decision Tree**)² para todas as abordagens. Ao final do desenvolvimento do trabalho, você deverá ser capaz de gerar resultados para preencher a Tabela 1.

Tabela 1. Comparação entre as abordagens.

	GA	Wrapper	Baseline
Acurácia			
Porcentagem de features			
Tempo de treinamento			
Tempo para busca das features			

¹ Arquivos de **treino** e **teste** disponíveis em:

https://drive.google.com/file/d/1o5208m0lxrsiglPZOtGWoddWXpL2h_Eq

² Considere o algoritmo disponível na biblioteca *scikit-learn* e parâmetros padrão (`random_state=1`) -

<https://scikit-learn.org/1.5/modules/tree.html>

Pontos de atenção:

- A acurácia final dos modelos deve ser calculada com os dados de teste.
- A etapa de seleção de features deve ser realizada considerando somente dados de treinamento. Assim, o uso de dados de teste para o cálculo da função de fitness é incorreto. A porcentagem de dados utilizada nessa fase impacta diretamente no tempo de processamento e é uma decisão de projeto.
- A função de fitness pode ser uma combinação entre acurácia e porcentagem de features selecionada, penalizando soluções com muitas features.
- A representação do problema, a função de fitness, o método de seleção, operador de crossover, taxa de mutação, critério de parada do algoritmo, entre outras escolhas, são decisões de projeto que devem ser discutidas em grupo antes da implementação.

Critérios de avaliação:

- (40%) Capacidade de justificar as decisões de implementação (por exemplo, escolha de parâmetros como tamanho da população inicial, operador de crossover, taxa de mutação, critério de parada, etc);
- (30%) Corretude das implementações;
- (20%) Capacidade de comparar e interpretar os resultados obtidos;
- (10%) Clareza e organização do código.

Entregas:

- Código-fonte com todas as implementações (GA, wrapper e baseline);
- Tabela com os resultados obtidos;
- Vídeo de até **15 minutos** com a apresentação de todo o grupo. Nessa apresentação, o grupo deve focar em explicar cada uma das decisões de projeto na implementação do GA e como elas foram implementadas e testadas.

Aviso importante: A entrega do código por si só não garante a pontuação nos critérios de avaliação apresentados acima. Durante a avaliação/teste de autoria, cada equipe será questionada sobre as decisões de implementação e deverá ser capaz de justificá-las adequadamente. A não compreensão das etapas do trabalho ou a incapacidade de explicar escolhas fundamentais resultará em descontos na nota, independentemente de o código estar funcional ou não. Em casos extremos, o grupo todo ou estudantes do grupo poderão receber a nota mínima.