

Instruções:

O trabalho será realizado em grupos já previamente estabelecidos

Ao realizar a entrega da avaliação, siga o padrão abaixo para o nome do arquivo

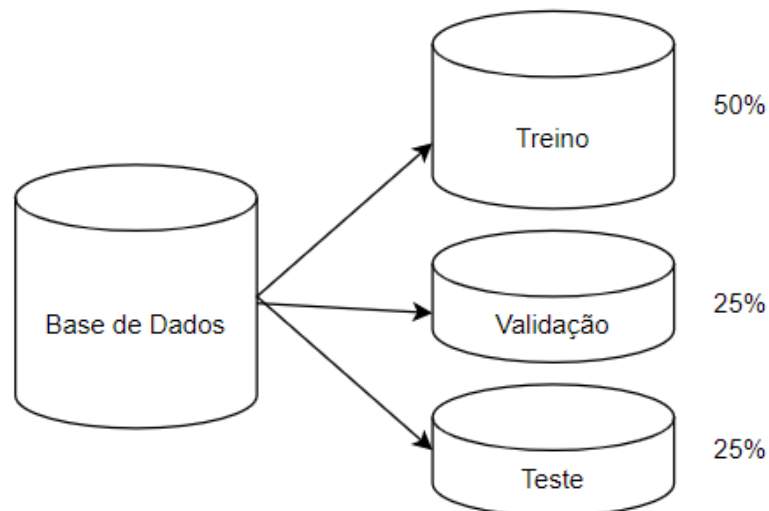
NomeSobrenome1-NomeSobrenome2-NomeSobrenome3.zip

Qualquer dúvida, entre em contato via email:
angeloorsatto@fag.edu.br

O objetivo deste trabalho consiste em comparar o comportamento, em termos de acurácia, de classificadores baseados em diferentes conceitos sobre uma mesma base de dados.

Para que o processo tenha base para análise, deverão ser executadas 10 repetições. Os valores a serem comparados deverão ser os valores médios das 10 execuções.

O primeiro passo consistirá na divisão da base original em três subconjuntos mutuamente exclusivos: treino, teste e validação. A instância que for designada para um conjunto não deve aparecer nos outros. O conjunto de treino deverá possuir 50% do tamanho do conjunto original. Já as bases de validação e teste, terão 25% da dimensão.



IMPORTANTE 1: no momento de separar a base original nos três conjuntos (treino, teste e validação), deve-se manter as proporções originais das classes. Por exemplo, se um conjunto possui 200 instâncias da classe A e 100 da classe B, o conjunto de treino terá 100 instâncias da classe A e 50 da classe B.

IMPORTANTE 2: a escolha das instâncias que formarão cada um dos conjuntos deve ser totalmente aleatória.

IMPORTANTE 3: lembre-se de sempre “bagunçar” os conjuntos de dados antes de fazer as divisões e de realizar o treinamento. A adoção de aleatoriedade adiciona robustez ao processo.

Depois de formados os conjuntos, o passo seguinte será o treinamento dos modelos de classificação. Deverão ser implementadas as estratégias dos K Vizinhos mais próximos (KNN), Árvore de Decisão (AD), Máquina de Vetor de Suporte (SVM) e Multilayer Perceptron (MLP).

Para se determinar quais os melhores parâmetros dos métodos de classificação, deve-se adotar o conjunto de validação. Por exemplo, digamos que estamos treinando um KNN e queremos decidir qual o melhor K a ser empregado. Deve-se treinar o classificador com o conjunto de treino e então variar o valor de K e analisar a acurácia do classificador sobre o conjunto de validação. O valor de K que obtiver a maior acurácia será utilizado no momento de classificar o conjunto de teste.

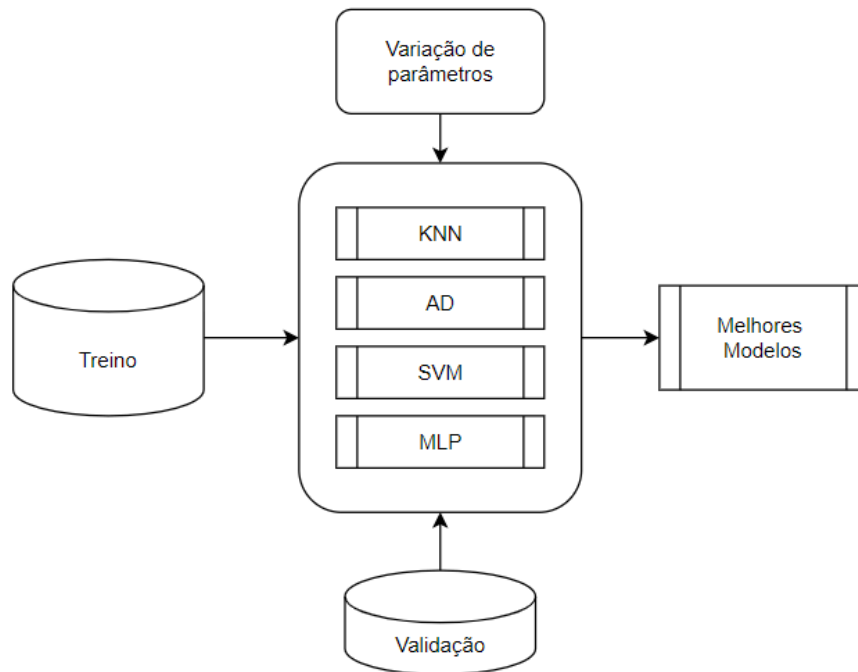
Os parâmetros que deverão ser definidos para cada classificador:

KNN: Valor de K e pesos dos votos (uniforme e distâncias)

AD: Com poda ou sem poda e critério para mensurar a qualidade de uma divisão (critério de Gini e entropia)

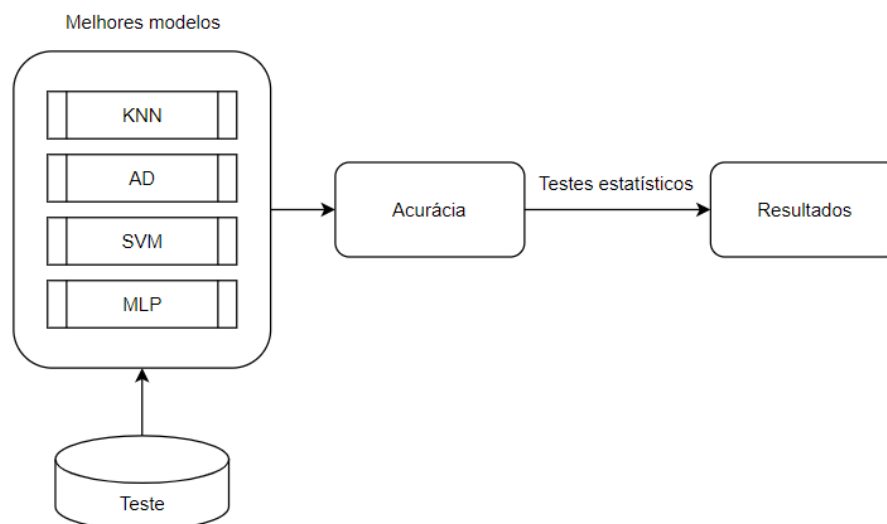
SVM: Valor do erro (C), Tipo de Kernel (Polinomial ou Radial)

MLP: Número de épocas de treino, taxa de aprendizagem, número de camadas escondidas, função de ativação (identidade, logística sigmóide, tangente hiperbólica e ReLu)



Definidos os melhores parâmetros para cada classificador, o passo seguinte será avaliar os seus desempenhos sobre o conjunto de teste. Nesta etapa deverá ser guardada a acurácia de cada classificador ao longo das 10 execuções.

Por fim, a última etapa consiste na comparação das acurácias dos métodos para descobrir qual deles obteve o melhor desempenho. Para tanto, deve-se realizar dois testes estatísticos. O primeiro teste será o método Kruskal-Wallis, que servirá para detectar se há diferença entre o desempenho dos algoritmos (independente de qual foi melhor ou pior). O segundo teste estatístico, será o teste de Mann-Whitney, a fim de comparar, dois a dois, os classificadores com o objetivo de avaliar se eles têm desempenhos significativamente diferentes e quem é o melhor.



Como fazer:

A linguagem e tecnologias utilizadas são de escolha de cada grupo

Não é necessário implementar os métodos de classificação. Neste caso, pode-se e é indicado, que sejam utilizadas implementações prontas dos métodos, ficando a carga da dupla apenas a implementação do framework e análise dos parâmetros e resultados. Da mesma forma, para o carregamento, aleatorização, divisão e sorteio dos conjuntos de treino, teste e validação podem ser utilizadas funções próprias das linguagens.

O que deverá ser entregue:

1 - Relatório detalhado contendo:

- descrição dos dados utilizados (o que a base significa, atributos, classes, número de classes, etc.);
- os parâmetros utilizados para calibração dos modelos;
- tecnologias utilizadas;
- acurácias obtidas por cada classificação ao longo das execuções, como no exemplo abaixo:

Repetição	KNN	SVM	AD	MLP
1	Acc	Acc	Acc	Acc
2	Acc	Acc	Acc	Acc
...
10	Acc	Acc	Acc	Acc
	Média / DP	Média / DP	Média / DP	Média / DP

- comparação individual entre cada modelo;
- análises estatísticas indicando se há diferença significativa entre os desempenhos dos classificadores;
- conclusões do trabalho

Todos os integrantes do grupo devem ser identificados no relatório.

O relatório deve ser claro, explicando cada processo e o que foi realizado durante cada etapa do trabalho.

O relatório deve ser entregue em PDF no formato de publicações e artigos da SBC ([Templates para Artigos e Capítulos de Livros da SBC](#))

Obs.: excluir a marcação de capítulo (chapter) contida no arquivo do link acima.

2 - Arquivos utilizados para a implementação deste trabalho:

Atenção: não compartilhem link dos arquivos de código

Todos os arquivos devem ser entregues compactados em um único arquivo, seguindo o padrão de nomenclatura: NomeSobrenome1-NomeSobrenome2-NomeSobrenome3-NomeSobrenome4-NomeSobrenome5.zip

3 - Apresentação:

A apresentação é de formato livre.

Cada grupo terá 20 minutos para apresentação e arguição do trabalho realizado, focando na descrição do problema, nos desempenhos obtidos e no resultado da análise estatística.

Obs.: não é necessário entregar o arquivo da apresentação, caso o grupo opte por elaborar outros arquivos para a apresentação.

A avaliação da apresentação será feita de forma individual.

4 - Entrega:

O trabalho deverá ser submetido na tarefa no Classroom até dia 05/11/2024

As apresentações serão nos dias 06/11/2024, 07/11/2024 e 13/11/2024, caso seja necessário a última data.

GRUPOS

Grupo 1:

- **Integrantes:** Juliano Notari, Vitor José Dietrich, Gustavo Miguel
- **Base de dados:** Ionosphere
- **Link:** <https://archive.ics.uci.edu/dataset/52/ionosphere>

Grupo 2:

- **Integrantes:** Guilherme Mattge, Igor Dall Forno, Henrique Nicolli
- **Base de dados:** Maternal Health Risk
- **Link:** <https://archive.ics.uci.edu/dataset/863/maternal+health+risk>

Grupo 3:

- **Integrantes:** Antonio Alencar, Bruno Romao, Carlos Rosa, Lucas Henrique e Richardson Korp
- **Base de dados:** AIDS Clinical Trials Group Study 175
- **Link:** <https://archive.ics.uci.edu/dataset/890/aids+clinical+trials+group+study+175>

Grupo 4:

- **Integrantes:** Daniel Henrique, Hudson Storti, Luiz Moll, Bruno Jobim e Kaique da Cruz
- **Base de dados:** Letter Recognition
- **Link:** <https://archive.ics.uci.edu/dataset/59/letter+recognition>

Grupo 5:

- **Integrantes:** Daniel Urbaneki, Everton Moscaleski, Lourenço Devequi, Wellington Augusto
- **Base de dados:** Glass Identification
- **Link:** <https://archive.ics.uci.edu/dataset/42/glass+identification>

Grupo 6:

- **Integrantes:** Saul Rian, Yago Kenji, Luiz Felipe Pronsati
- **Base de dados:** Spambase
- **Link:** <https://archive.ics.uci.edu/dataset/94/spambase>

Grupo 7:

- **Integrantes:** Alan Fernando Lenz, George Paiva Thome Speltz, Pedro Felipe Galvan
- **Base de dados:** Rice (Cammeo and Osmancik)
- **Link:** <https://archive.ics.uci.edu/dataset/545/rice+cammeo+and+osmancik>

Grupo 8:

- **Integrantes:** Caio Augusto, Gabriel Augusto, Gabriel Tomiazzi, João Pedro Rebello, Pedro Augusto
- **Base de dados:** Predict Students' Dropout and Academic Success
- **Link:** <https://archive.ics.uci.edu/dataset/697/predict+students+dropout+and+academic+success>

Grupo 9:

- **Integrantes:** Gabriel Ceola, Arthur Baron, Guilherme da Cruz, Karol Dambrosio
- **Base de dados:** Wine
- **Link:** <https://archive.ics.uci.edu/dataset/109/wine>

Grupo 10:

- **Integrantes:** Eduardo Araujo, Rafael Pagliari, Igor Wolf, Eduardo Arvelino, Gustavo Tomadon
- **Base de dados:** Dry Bean
- **Link:** <https://archive.ics.uci.edu/dataset/602/dry+bean+dataset>