

## Trabalho Prático 2 - Machine Learning

O objetivo deste trabalho é implementar e avaliar os algoritmos Naive Bayes, KNN e Árvores de Decisão aplicados a diferentes problemas de classificação e regressão. Adicionalmente o algoritmo K-Médias deve ser utilizado para clustering.

### 1. (Classificação - Base de Dados Iris) :

#### 1.1. Conforme visto nas primeiras aulas, busque e carregue o **Iris Dataset**:

- 4 Features: Sepal Length, Sepal Width, Petal Length, Petal Width
- 3 Classes: Versicolor, Setosa e Virginica



#### 1.2. Com objetivo de visualizar os dados, plote os gráficos a seguir:

- As 6 combinações em 2D de dois dos quatro atributos do Iris dataset.
- As 4 combinações em 3D de três dos quatro atributos do Iris dataset.

#### 1.3. **Implemente** o algoritmo Gaussian Naive Bayes e aplique à base de dados.

- Plote a superfície de decisão obtida pelo algoritmo.

#### 1.4. **Implemente** o algoritmo KNN e aplique à base de dados.

- Encontre o melhor valor para o parâmetro  $K(1...m)$  e melhor métrica de distância (pelo menos 3 diferentes).
- Plote a superfície de decisão obtida pelo algoritmo.

#### 1.5. **Utilize** o algoritmo Arvore de Decisão para classificar a base de dados.

- Diferentemente dos itens 1.3 e 1.4, você não precisará implementar o algoritmo, utilize a implementação DecisionTreeClassifier do pacote sklearn.
- Entenda o funcionamento do algoritmo DecisionTreeClassifier e busque pelos melhores parâmetros (melhor métrica de pureza e melhor critério de parada).
- Plote a fronteira de decisão gerada por este método bem como a árvore final gerada (Utilize sklearn.tree).

#### 1.6. Compare a eficácia dos classificadores com relação à acurácia, utilizando a divisão dos dados (Treinamento, Validação e Teste). Execute cada experimento 10 vezes e compare as médias dos resultados.

#### 1.7. Retire os rótulos da base de dados e **Implemente** o algoritmo K-Médias para encontrar agrupamentos nos dados. Utilizando $K=3$ , teste diferentes métricas de distância (pelo menos 3) e avalie os clusters obtidos.

- Plote os clusters e também os centróides finais obtidos.
- Compare visualmente os clusters com as classes reais.

2. (**Classificação - Base de Dados Bi-dimensional**) Com o objetivo de fixar o conteúdo, você deverá implementar os algoritmos Gaussian Naive Bayes, KNN e Árvore de decisão e aplicar à base de dados para predição de aprovação de um estudante com base nos resultados de 2 avaliações realizadas por ele. A base de dados(arquivo ex2data1.txt utilizada no TP1) contem dados históricos referentes a avaliações passadas, onde as colunas da base são: Avaliação 1, Avaliação 2 e resultado(aprovado ou reprovado).

Após implementar os 3 algoritmos, mostre a superfície de decisão gerada por eles e compare os resultados(Acurácia média e tempo computacional para o treinamento e o teste) utilizando validação cruzada.

3. (**Classificação - Base Câncer de Mama**) O objetivo agora é implementar os 3 algoritmos para resolução do problema de detecção de pacientes com câncer de mama. Os algoritmos deverão ser comparados ao final, onde você deverá apontar os prós e contras dos algoritmos com relação ao desempenho obtido. Não se esqueça de **normalizar** e **separar** os dados(treinamento, validação e teste) para efetuar uma avaliação correta.

A base de dados a ser utilizada nesta questão é conhecida como **Wisconsin Diagnostic Breast Cancer - WDBC** e foi produzida pelo Dr. William H. Wolberg, pesquisador do departamento de Cirurgia Geral da Universidade de Wisconsin. Ela está disponível publicamente no UCI Machine Learning Repository (<<https://archive.ics.uci.edu/static/public/17/breast+cancer+wisconsin+diagnostic.zip>>). Esta base de dados é composta por 569 pacientes (357 saudáveis e 212 com câncer) e 32 campos (dos quais 30 atributos são úteis) no total, sendo que o primeiro pode ser descartado, por se tratar do identificador do paciente e o último campo é o rótulo da classe (0 - Saudável, 1 - com Câncer).

4. (**Regressão - Base de Dados Food Truck**) Implemente os algoritmos KNN e Árvore de decisão para resolver o problema de Regressão de prever o lucro para uma empresa de food truck. Este é o mesmo problema utilizado no TP1. Você deverá alterar o KNN para que ela possa resolver problemas de regressão e, quanto ao algoritmo Árvore de decisão, utilize a implementação do pacote sklearn adequada à tarefa de regressão. O arquivo ex1data1.txt contém os dados que deverão ser utilizados. Cada linha do arquivo é uma cidade; a primeira coluna contem a população e a segunda o lucro, da cidade correspondente. Após aplicar os algoritmos e obter os resultados para os dados de teste, compare os algoritmos com relação ao erro médio quadrático obtido e plote a curva gerada por ambos.

**O que deve ser entregue:** Deve ser entregue um relatório contendo o código e um relatório contendo a análise de cada algoritmo aplicado aos problemas. Um notebook contendo relatório e códigos também pode ser enviado.