

Disciplina: Mineração de Dados

Professor: Murilo Vargas da Silva

Orientações para trabalho final da disciplina

Objetivo:

Estimular o aluno a aplicar os conhecimentos apresentados no decorrer da disciplina em problemas reais de mineração de dados, utilizando as técnicas de seleção, pré-processamento e transformação de dados, técnicas de visualização de dados, análise descritiva, análise de grupos, classificação e estimação/regressão.

Entregas:

- Relatório contendo a descrição da base de dados, detalhamentos das atividades executadas com descrição e imagens que permitam avaliar o resultado da atividade, data dia 23/06 via MOODLE;
- Slides utilizados na apresentação, na data da apresentação via MOODLE;
- Código fonte disponibilizado no GITHUB.

Atividades desejadas:

1 – Seleção e pré-processamento de dados:

- Escolha uma base de dados em <http://archive.ics.uci.edu/ml>;
- Avalie as características da base de dados: problema a ser investigado, número de amostras, número de atributos, tipos de atributos, possui valores ausentes?
- Utilizando a linguagem **Python**, junto com as bibliotecas **pandas** e **numpy**, crie um código que efetue uma limpeza de dados aplicando as técnicas apresentadas na aula de hoje.

2 – Normalização e redução de dados:

- Utilizando a base de dados escolhida;
- Primeiramente realize o procedimento de limpeza de dados Atividade 1;
- Utilize alguma técnica de normalização de dados;
- Utilize a técnica PCA e plot os dois principais componentes.

3 – Análise descritiva de dados - Visualização:

- Utilizando a base de dados escolhida;
- Primeiramente realize o procedimento de limpeza de dados Atividade 1;
- Realize a distribuição de frequência para algum(s) atributo(s) da base dados;
- Utilize alguma técnica de visualização para analisar os dados com base na distribuição de frequência. (Histograma, Gráfico de setores, dispersão, etc)

4 – Análise descritiva de dados - Medidas:

- Utilizando a base de dados escolhida;
- Primeiramente realize o procedimento de limpeza de dados Atividade 1;
- Calcular medidas de resumo apresentadas na aula de hoje:
 - Medidas de tendência central;
 - Medidas de dispersão;
 - Medidas de posição relativa;
 - Medidas de associação.

5 – Análise de grupos:

- Utilizando a base de dados escolhida;
- Primeiramente realize o procedimento de limpeza de dados Atividade 1;
- Utilizar o algoritmo KMeans
 - Utilizando apenas 2 principais componentes (PCA);
 - Rodar KMeans variando o número de grupos (parâmetro k)
 - Plotar os resultados dos agrupamentos para diferentes valores de K;
- Utilize medidas para avaliar a qualidade dos agrupamentos: coeficiente de forma, homogeneidade, etc

5 – Classificação - KNN:

- Utilizando a base de dados escolhida;
- Primeiramente realize o procedimento de limpeza de dados Atividade 1;
- Fazer os procedimentos de normalização que achar necessário;
- Utilizar o algoritmo K-NN
 - Fazer a divisão da base utilizando:
 - Holdout (Treinamento 70% e Teste 30%)
 - Cross-Validation (k=10)
 - Classificar com K-NN e calcular as seguintes métricas
 - Matrix de confusão
 - Acurácia
 - F1 Score

6 – Classificação - SVM:

- Utilizando a base de dados escolhida;
- Primeiramente realize o procedimento de limpeza de dados Atividade 1;
- Fazer os procedimentos de normalização que achar necessário;
- Utilizar o algoritmo SVM
 - Fazer a divisão da base utilizando:
 - Holdout (Treinamento 70% e Teste 30%)
 - Cross-Validation (k=10)
 - Classificar com SVM e calcular as seguintes métricas
 - Matrix de confusão
 - Acurácia
 - F1 Score