

TECHNICAL REPORT

Aluno: Vítor Barbosa da Silva

1. Introdução

O dataset de classificação contém dados relacionados à transações em cartões de crédito de diferentes usuários em diferentes lugares da Europa. Possui colunas que apresentam anonimamente os dados das transações, como horário, local, etc. A atividade em cima do dataset consiste em analisar os dados e treinar o sistema para classificar se a transação tem caráter fraudulento ou não.

Já o dataset de regressão contém dados relacionados à passagens aéreas compradas. Possui colunas que apontam os dados da passagem, como voo, linha aérea, cidade de origem e de destino, preço, etc. A atividade em cima do dataset consiste em analisar os dados e treinar o sistema para analisar o quanto os atributos, como classe, linha aérea e etc, podem contribuir para a variação no preço da passagem.

2. Observações

Devido ao grande volume de dados tive que realizar uma amostragem estratificada para ser utilizada no knn, pois o tempo e o custo de processamento para gerar os resultados estava muito elevado. Assim, utilizei na terceira questão o valor de 5000 para cada uma das classes e na quarta questão o valor de 30000 mil para cada uma das classes, já que a implementação do knn nessa questão não seria necessariamente manual. Um ponto a ser destacado foi a dificuldade que tive para tratar os dados, já que as atividades exigem que o dados sejam numéricos e nem todos eram, assim não consegui realizar a última questão.

3. Resultados e discussão

Nesta seção deve-se descrever como foram as resoluções de cada questão. Crie sessões indicando a questão e discuta a implementação e resultados obtidos nesta. Explique o fluxograma do processo de cada questão, indicando quais processamentos são realizados nos dados. Sempre que possível, faça gráficos, mostre imagens, diagramas de blocos para que sua solução seja a mais completa possível. Discuta sempre sobre os números obtidos em busca de motivos de erros e acerto.

1ª Questão.

Nesta questão o requisito era fazer um pré-processamento do dataset. iniciei carregando o dataset, após isso utilizei o comando `df = df.dropna()` para excluir linhas que possuíam algum dado nulo ou branco. Em seguida precisei apenas realizar o salvamento do novo dataset criado na minha pasta de datasets

```
dataframe_ajustado = "../AV1/datasets/novo_creditcard.csv"
df.to_csv(dataframe_ajustado, index=False)
```

2ª Questão.

Nesta questão foi necessária a implementação da função KNN de forma manual:

```
def knn(X_train, y_train, X_test, k, dist_func, cov_inv=None):
    y_pred = []

    for test_point in X_test:
        # Calculando distâncias do ponto de teste para todos os pontos de treino
        distances = []
        for train_point, train_label in zip(X_train, y_train):
            # Passar cov_inv apenas se a função de distância for Mahalanobis
            if dist_func == dist_mahalanobis:
                dist = dist_func(test_point, train_point, cov_inv) # Passando cov_inv para Mahalanobis
            else:
                dist = dist_func(test_point, train_point)
            distances.append((dist, train_label[0]))

        distances.sort(key=lambda x: x[0])
        k_neighbors = [label for _, label in distances[:k]]
        majority_class = max(set(k_neighbors), key=k_neighbors.count)
        y_pred.append(majority_class)

    return np.array(y_pred)
```

Foi necessário fazer a separação do conjunto de treino e de teste utilizando a função `train_test_split`, onde o `X` armazenou os valores contidos nas demais colunas do dataset, excetuando-se a de `Class`, e `y` armazenou os valores de `Class`. Para `k=7` as acurácias obtidas para cada distância foram as seguintes:

```
Acurácia do modelo dist_euclidiana: 99.85%  
Acurácia do modelo dist_manhattan: 99.85%  
Acurácia do modelo dist_mahalanobis: 98.60%  
Acurácia do modelo dist_chebyshev: 99.75%
```

obs: Achei um pouco estranho os valores das distâncias euclidiana e Manhattan terem sido os mesmos, mas tentei outras abordagens e o valor delas continuava se repetindo.

3ª Questão.

Para essa questão, que pedia a melhor distância, eu escolhi a Euclidiana. Após isso fiz a normalização para verificar se iria interferir no valor da acurácia. O código utilizado foi o seguinte:

```
# Normalização Logarítmica  
X_train_pos = X_train - np.min(X_train) + 1  
X_test_pos = X_test - np.min(X_test) + 1  
X_train_log = np.log1p(X_train_pos)  
X_test_log = np.log1p(X_test_pos)  
  
# Normalização Z-score (média zero e variância unitária)  
scaler = StandardScaler()  
X_train_z = scaler.fit_transform(X_train)  
X_test_z = scaler.transform(X_test)
```

Em seguida utilizei a função KNN, ainda implementada de forma manual, para calcular o acertos e gerar o resultado a ser analisado:

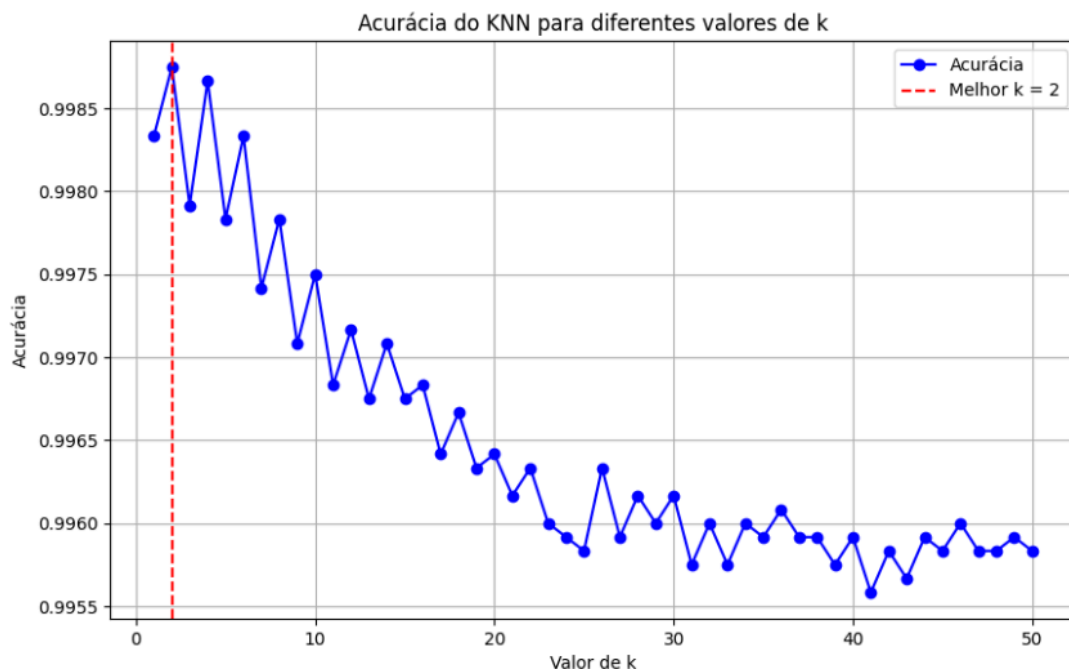
```
Acurácia do KNN com diferentes normalizações:  
Logarítmica: 92.60%  
Z-score: 99.35%
```

Após analisar o resultado, constatei que a normalização interfere no valor da acurácia.

4ª Questão.

Nesta questão é pedido que escolha a melhor parametrização, com base nos resultados a melhor foi a Z - score, então foi essa a escolhida. Após fazer novamente a separação e normalização dos conjuntos, utilizei a implementação mais convencional do KNN, da biblioteca sklearn.

Utilizando um intervalo de 1 a 51, para verificação do melhor K, obtive o resultado de que o valor melhor é o $k=2$, com uma acurácia de 99,88%.



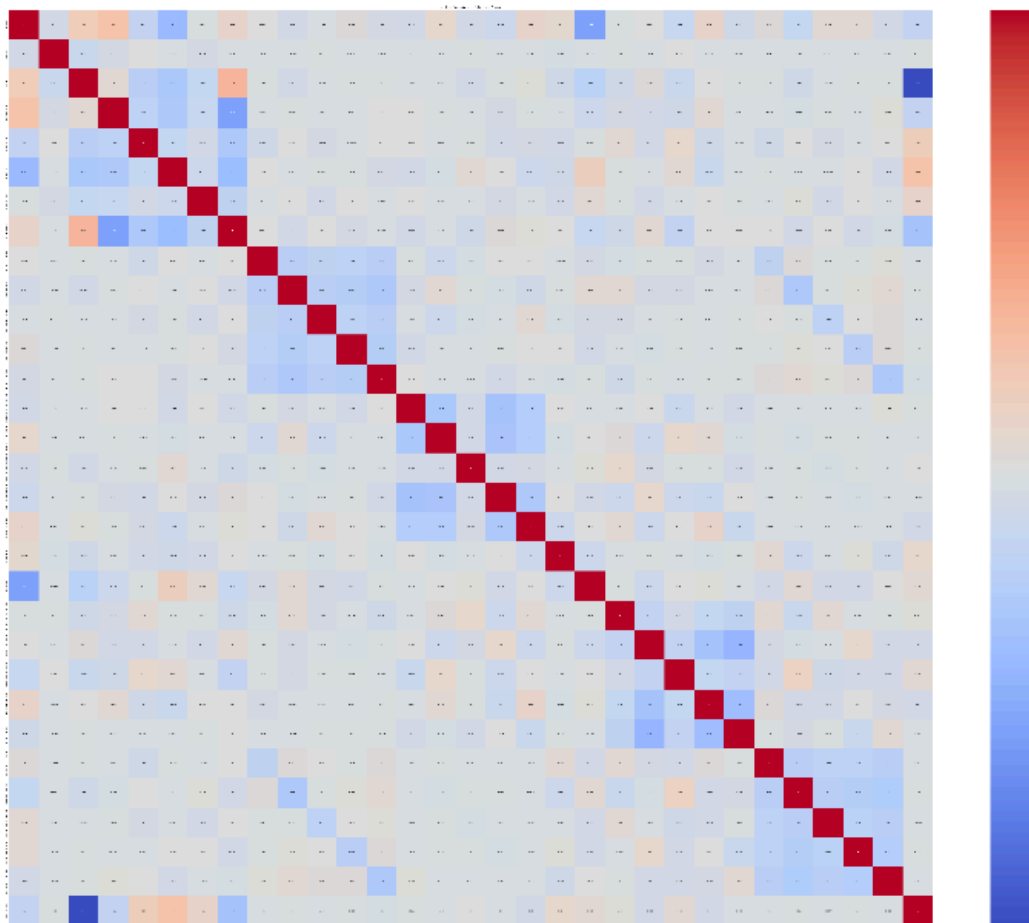
5ª Questão.

Nesta questão foi utilizada uma outra metodologia, a de regressão. Inicialmente comecei com o pré-processamento do dataset. Excluí linhas com valores NaN e também algumas colunas que não tinham relevância para a análise. As colunas foram a primeira que não tinha nome e a terceira que tinha o nome flight.

Tive que realizar uma conversão das colunas categóricas para gerar a matriz de correlações e verificar qual tinha mais relação com a coluna price. Para isso utilizei uma abordagem chamada one_hot, onde converte dados categóricos do dataset para números:

```
df_one_hot = pd.get_dummies(df, columns=colunas_categoricas, drop_first=True)
matriz_correlacao = df_one_hot.corr()
```

Em seguida gerei a matriz de correlações, para analisar.



6ª Questão.

Nesta questão, utilizando-se do atributo mais relevante já previsto na questão anterior, foi implementada a regressão linear. Para tal foi utilizado novamente o one-hot encoding para converter os valores da coluna airline e também implementadas as métricas:

```
RSS = np.sum((y - y_pred) ** 2)
MSE = mean_squared_error(y, y_pred)
RMSE = np.sqrt(MSE)
R_squared = r2_score(y, y_pred)
```

A geração dos resultados das métricas foi a seguinte:

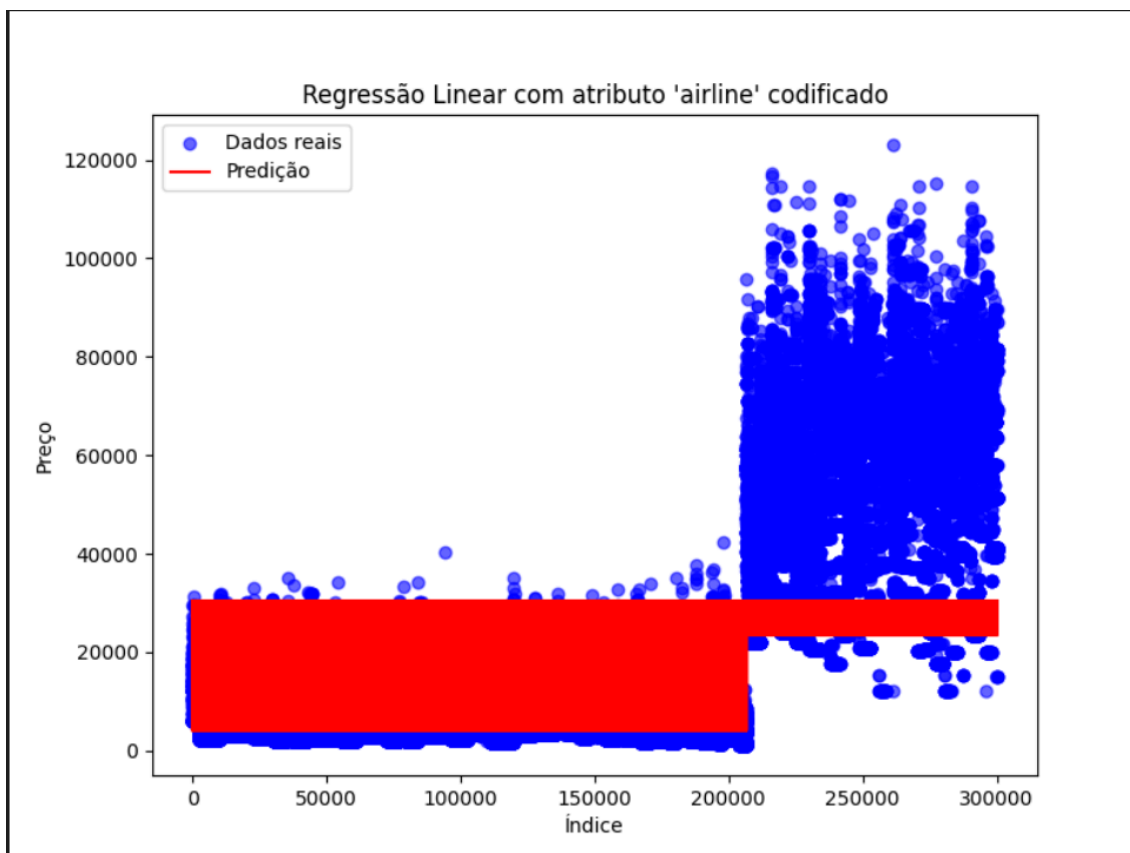
RSS: 120204394446504.91

MSE: 400477071.52

RMSE: 20011.92

R^2 : 0.22

E o gráfico da reta de regressão e a nuvem de atributo foi o seguinte:



Com base nos dados da Matriz, pude ver que o elemento que melhor se relaciona com o price é uma linha aérea chamada Vivara. Assim, a coluna airline se torna a mais relevante com base na matriz.

Obs: Devido à grande quantidade de dados, é necessário dar o zoom para ver os dados da matriz. Porém só é possível no console, quando o código é gerado, assim a imagem é apenas para comprovar graficamente.

4. Conclusões

Nem todos os resultados esperados foram satisfeitos, pois não consegui concluir a última questão. Mas acredito que os resultados obtidos nas demais questões foram satisfatórios, pelo menos para mim. Considero um avanço no entendimento e execução dos conteúdos já vistos. Vale também ressaltar a importância da atividade proposta e



do relatório para o entendimento e aprendizado dos alunos com relação ao conteúdo abordado.

5. Próximos passos

Tenho como sugestão pegar um dataset de alguém aleatoriamente e refazer juntamente com a turma expondo os pontos importantes e corrigindo juntamente com a turma o código.