

Template para entrega do projeto da disciplina
Ciência de Dados e Inteligência Artificial
Fase 2

Nome do estudante **Eduarda Pereira Blanco**

Desenvolva um processo de ciência de dados no Orange Data Mining, cobrindo os elementos abaixo. Para cada um dos itens solicitados é necessário inserir imagens que evidenciem o trabalho realizado.

Exploração dos dados

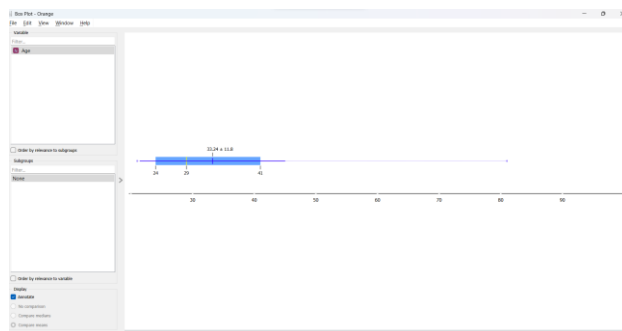
Que tipo de experimentos você fez na exploração dos dados (verificação de outliers, cálculos de médias, dados inválidos etc.).

O conjunto de dados escolhidos compreende um dataset de diabetes, que é originário do National Institute of Diabetes and Digestive and Kidney, disponível em: <https://www.kaggle.com/datasets/akshaydattatraykhare/diabetes-dataset> e considerado o padrão gold de dados.

O objetivo dos dados é prever diagnosticamente se o paciente tem diabetes, com base em medidas de diagnósticos incluídas no conjunto de dados. Todos os pacientes são mulheres com idade superior a 21 anos.

Alguns conceitos foram levantados durante a exploração de dados:

- Resultados: 500 com resultado de “não diabetes” e 268 com resultado de “diabetes”
- Idade:
 - Média: 33,2 anos
 - Moda: 22 anos
 - Encontrado o valor máximo de 81 anos, podendo ser considerado um único outlier como apresentado pelo box plot



Data Table - Orange

	Outcome	Age	BMI	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	abetesPedigreeFunci
460	0	81	25.9	9	134	74	33	60	0.46
454	0	72	19.6	2	119	0	0	0	0.832
667	1	70	32.5	4	145	82	18	0	0.235
124	0	69	26.8	5	132	80	0	0	0.186
685	0	69	0	5	136	82	0	0	0.64
675	0	68	35.6	8	91	82	0	0	0.587
490	0	67	26.1	8	194	80	0	0	0.551

- Quantidade de vezes grávida:
 - Média: 3,8
 - Moda: 1
 - Máximo: 17 porém não consideraria outlier visto que é um valor muito próximo de demais

Data Table - Orange

File Edit View Window Help

Info
768 instances (no missing data)
8 features
Target with 2 values
No meta attributes.

Variables
☒ Show variable labels (if present)
☒ Visualize numeric values
☒ Color by instance classes

Selection
☒ Select full rows

	Outcome	Age	BMI	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	diabetesPedigreeFunc
140	1	47	40.9	17	163	72	41	114	0.817
89	1	43	37.1	15	136	70	32	110	0.153
299	1	46	36.6	14	100	78	25	184	0.412
456	1	38	33.6	14	175	62	30	0	0.212
29	0	57	22.2	13	145	82	19	110	0.245
275	0	52	34.2	13	106	70	0	0	0.251
87	0	45	36.6	13	106	72	54	0	0.178
358	1	44	39.9	13	129	0	30	0	0.569
692	1	44	42.3	13	158	114	0	0	0.257
324	1	43	26.8	13	152	90	33	29	0.731
73	1	42	43.4	13	126	90	0	0	0.583
519	0	41	32.8	13	76	60	0	0	0.18
745	0	39	40.6	13	153	88	37	140	1.174
636	1	38	31.2	13	104	72	0	0	0.465
583	0	62	26.5	12	121	78	17	0	0.259
376	1	58	39.2	12	140	82	43	325	0.528
359	0	48	35.3	12	88	74	40	54	0.378
746	0	46	30	12	100	84	33	105	0.488

- BMI (Body Mass Index):
- O equivalente ao IMC aqui no Brasil, quanto maior o valor, maior a relação com sobrepeso e obesidade.
- Média: 31 – considerado obesidade
- Moda: 32 – considerado de obesidade
- Máximo: 67.1

Data Table - Orange

File Edit View Window Help

Info
768 instances (no missing data)
8 features
Target with 2 values
No meta attributes.

Variables
☒ Show variable labels (if present)
☒ Visualize numeric values
☒ Color by instance classes

Selection
☒ Select full rows

	Outcome	Age	BMI	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	diabetesPedigreeFunc
178	1	26	67.1	0	129	110	46	130	0.319
446	1	23	59.4	0	130	78	63	14	2.42
674	0	22	57.3	3	123	100	35	240	0.88
126	1	26	55	1	89	30	42	99	0.496
221	1	25	53.2	0	162	76	56	100	0.759
304	1	28	52.9	5	115	98	0	0	0.209
194	1	40	52.3	11	135	0	0	0	0.578
248	0	23	52.3	0	165	90	33	680	0.427
156	1	36	50	7	152	88	44	0	0.337
100	1	31	49.7	1	122	90	51	220	0.325

- Tabela referência:



Escolha de, ao menos, três algoritmos de aprendizado para a modelagem

Apresente os algoritmos utilizados e justifique a escolha.

1. KNN

A utilização do modelo K-Nearest Neighbors (KNN) foi uma escolha para a classificação de diabetes em mulheres devido à sua simplicidade e à falta de suposições sobre a distribuição dos dados. Ele classifica novos casos com base na proximidade aos K vizinhos mais próximos no conjunto de treinamento, o que é intuitivo e eficaz para detectar padrões não lineares. Também, optei devido a sua pré preparação de dados realizada através do sistema Orange.

2. Logistic Regression

A regressão logística foi escolhida devido a eficácia para a classificação de diabetes devido à sua simplicidade e capacidade de fornecer probabilidades interpretáveis sobre a presença da doença. Ela assume uma relação linear entre as características e a probabilidade de uma classe, o que facilita a modelagem e a interpretação dos. A regressão logística também permite uma avaliação clara do desempenho por meio de métricas como acurácia – AUC.

3. Tree

As árvores de decisão são adequadas para a classificação de diabetes por sua capacidade de modelar relações complexas e não lineares de maneira visual e intuitiva, facilitando a interpretação dos critérios de decisão. Elas lidam bem com dados categóricos e contínuos, não exigem normalização dos dados e são robustas a outliers como os apresentados acima.

Preparação dos dados de acordo com as características dos algoritmos de aprendizado escolhidos

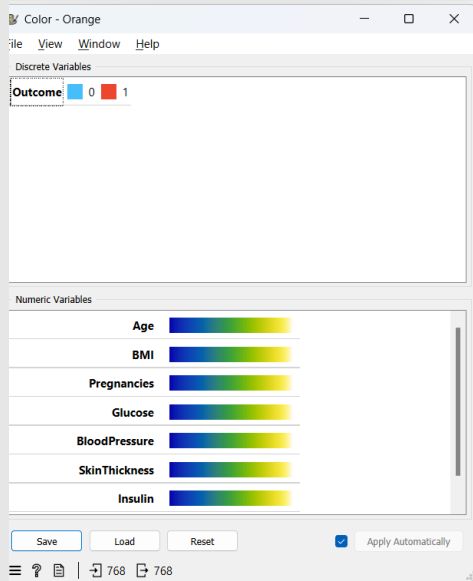
Descreva o processo realizado para essa etapa.

Retirei os dados nulos/vazios e realizei os primeiros experimentos com os dados não tratados, considerando que por padrão o modelo KNN, realiza um pré-processamento deste dados.

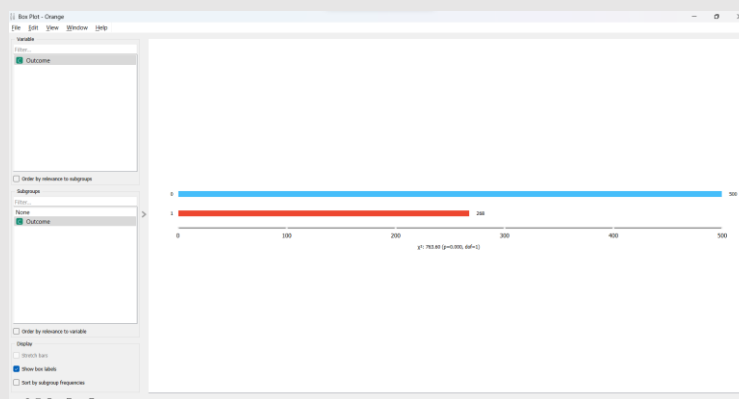
Já para a aplicação dos algoritmos Tree e Logistic Regression realizei um equilíbrio de casos positivos e negativos para acusação de diabetes, utilizando os métodos apresentados na aula 07, parte 05, evitando vieses. Assim normalizando os números entre pacientes com e sem diabetes, tornando os dados equilibrados para aplicação do algoritmo.

Apesar da normalização dos dados, realizei mais uma vez Tree e Logistic Regression nos dados não tratados para comparar seus resultados.

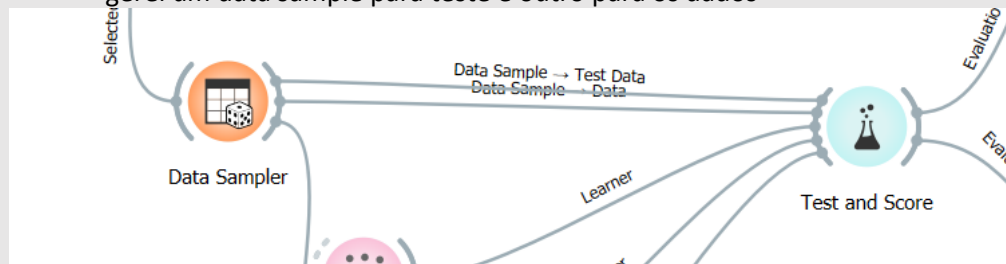
4. Para iniciar a exploração de dados coloquei a opção de coloração por resultado: azul para “não diabetes” e vermelho para “diabetes”.



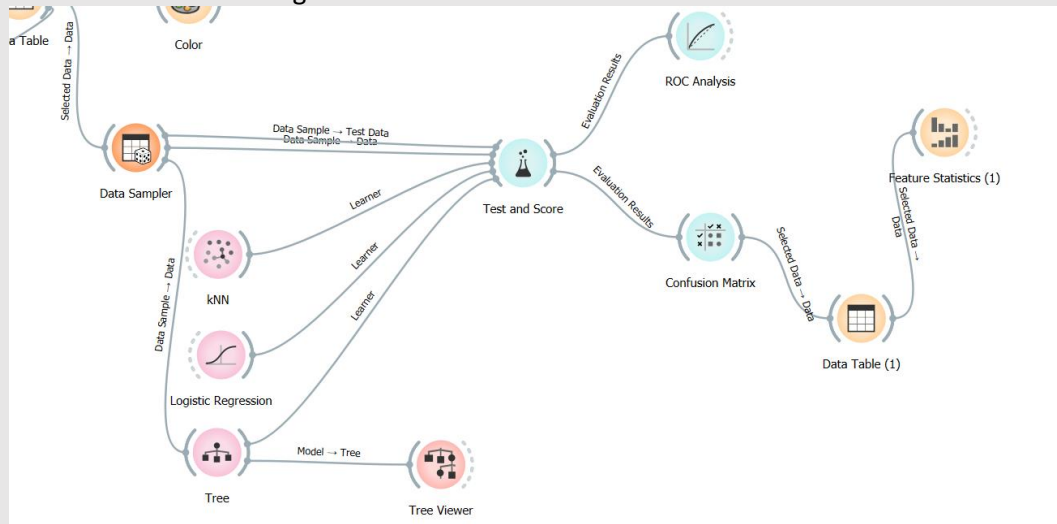
- 4.1. Seguindo no âmbito da exploração de dados foi incluído o visualização de feauture statistics e box plot, onde consegui visualizar os conceitos de média, moda, máxima, e os outliers, além de compreender o valor total por resultado.



5. Através da seleção de dados utilizei data sample que conectado ao test and score, gerei um data sample para teste e outro para os dados

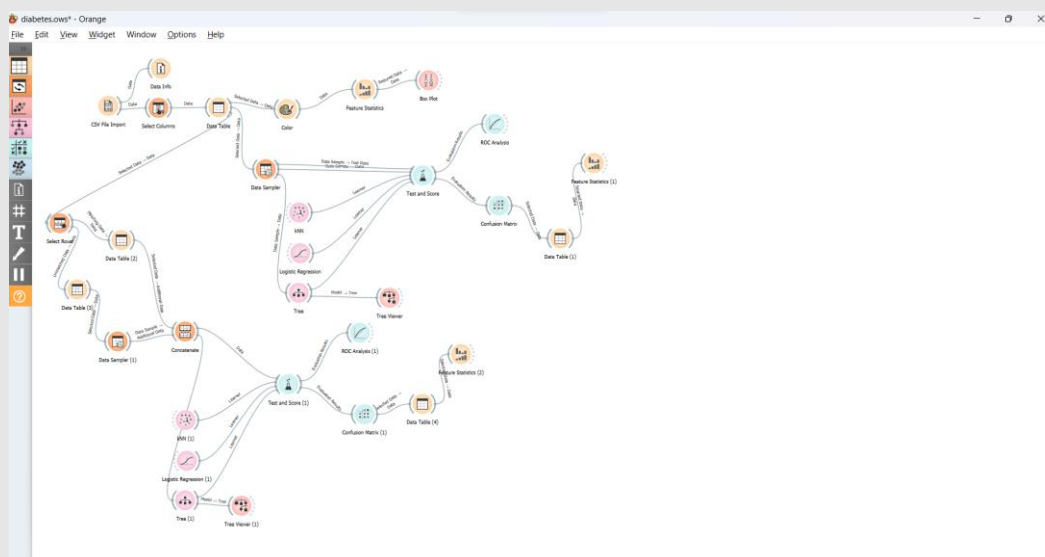


6. Conectei os 3 algoritmos utilizados ao test and score



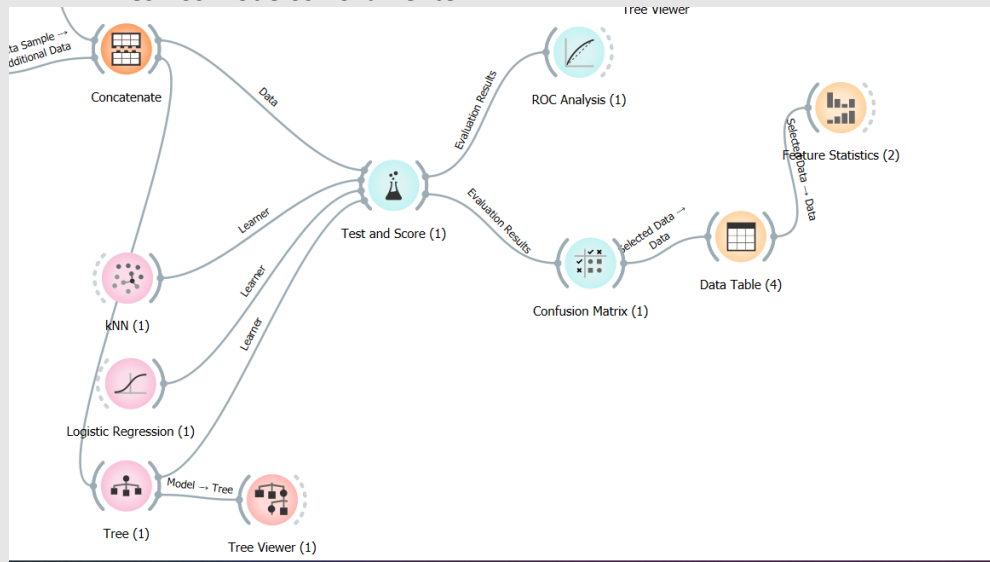
- 6.1. Conectei a visualização da árvore a própria árvore para análises
7. Conectei a visualização da confusion matrix e ROC analysis
8. Por fim, conectei uma visualização de tabelas na confusion matrix para identificar os dados que obtiveram erro durante o experimento.

Parte 2 (dados normalizados):



9. Extraí uma seleção de linhas que possuíam o resultado de 1, no caso diabetes. Totalizando os 268 casos.

10. Conectei uma tabela de dados com esses 268 casos
11. Conectei uma tabela com o restante dos casos de resultado 0, não diabetes.
12. Conectei uma amostragem a tabela de dados com resultado 0, deixando apenas 268 casos.
13. Realizei a concatenação de ambas as tabelas gerando uma de 536 casos e rodei os mesmos modelos novamente.



Relato dos experimentos e lições aprendidas

Apresente uma reflexão acerca dos resultados obtidos com este projeto.

- KNN

Considerando os dados não normalizados o modelo apresentou a consistência entre os valores de acurácia (0.7089 e 0.6942) está mostrando um desempenho relativamente estável em diferentes conjuntos de validação ou diferentes execuções. A acurácia média (em torno de 0.70) indica que o modelo tem um desempenho razoável, mas há espaço para melhorias.

Considerados os dados normalizados a acurácia média dos valores fornecidos é bastante alta, especialmente o valor de 0.7876, que sugere que o modelo KNN(1) está alcançando um bom desempenho geral na maioria dos casos.

- Tree

Considerando os dados não normalizados o modelo de árvore de decisão apresenta uma acurácia média que varia entre 63% e 70% em diferentes execuções ou validações, com algum nível de consistência em torno de 70% em algumas validações. No entanto, o valor baixo de 0.3486 em uma métrica de erro sugere que o modelo pode ter problemas de desempenho em aspectos específicos. Em comparação com outros modelos, as árvores de decisão não está performando tão bem.

Considerando os dados normalizados a comparação entre os dois conjuntos de scores para o modelo de árvore de decisão mostra uma melhoria significativa na versão mais recente (Tree (1)). A acurácia aumentou de uma faixa de 63%-70% para 74.9%-78.8%, indicando um desempenho muito melhor e mais consistente. A métrica de erro também mostrou uma

variação, mas a versão mais recente do modelo demonstra um avanço em termos de desempenho geral.

- Logistic Regression

O modelo de Regressão Logística tem um desempenho robusto, com acurácia média alta e consistente na versão original, variando de 76.6% a 81.8%. A versão mais recente (Logistic Regression (1)) mostra uma acurácia máxima melhor (84.0%), mas com uma média um pouco menor e uma métrica de erro mais alta (0.4965). Apesar de a melhoria na acurácia máxima seja notável, a métrica de erro sugere que ainda há espaço para ajustes. No geral, ambos os modelos têm um desempenho competitivo, mas a versão (1) pode indicar que a acurácia foi aprimorada em alguns aspectos, enquanto o erro precisa ser melhorado.

Caso:

Optei por normalizar os dados visto que a logistic regression apresentou um número alto de falso negativo, levando em consideração o contexto dos dados, um falso negativo é muito mais prejudicial neste caso, imagine o seguinte cenário: paciente recebe o resultado de não ter diabetes, já considerando os dados vistos anteriormente que possui agravantes como sobrepeso e alto valor de insulina, quando olhamos para os casos que apresentaram falsos negativos, temos uma média de insulina de 77, considerando valores coletados em jejum o normal previsto é entre 5 a 29,0 $\mu\text{IU/mL}$. Sabendo as complicações da diabetes a médio e longo prazo um erro de 40,5% de pacientes apresentaria até mesmo a morte desses pacientes.

Confusion Matrix - Orange

File View Window Help

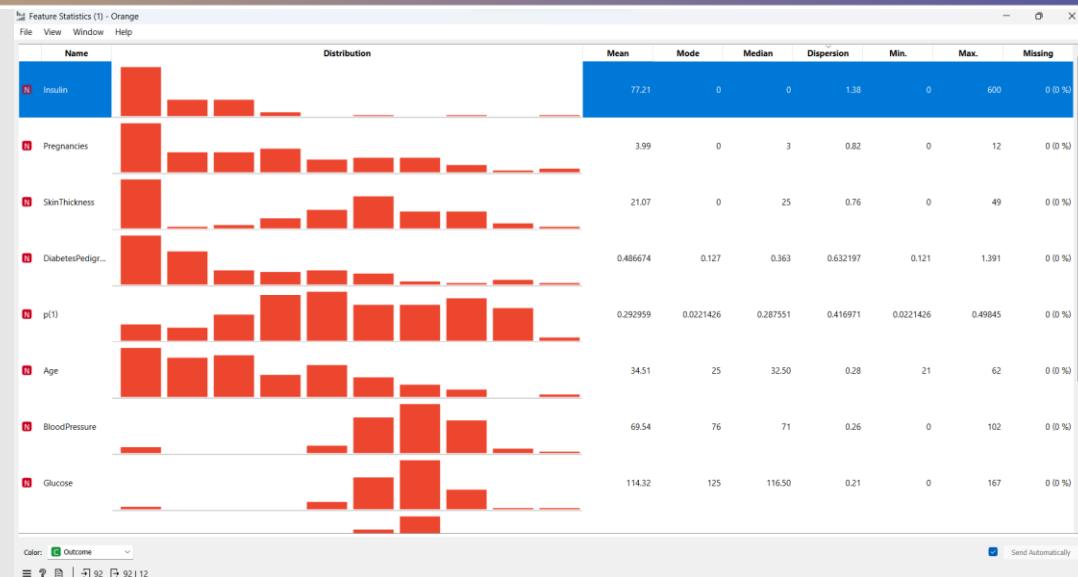
Learners

kNN

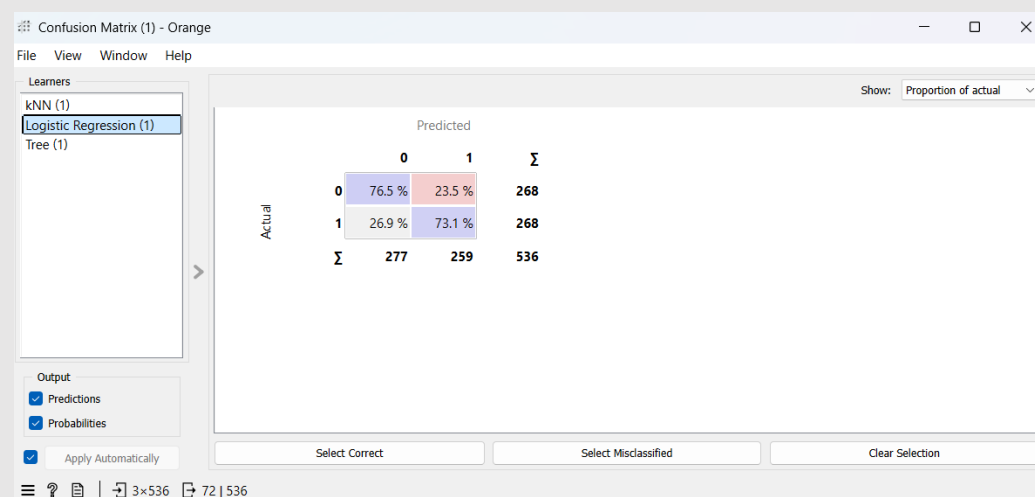
Logistic Regression

Tree

		Predicted		Σ
		0	1	
Actual	0	86.6 %	13.4 %	388
	1	40.5 %	59.5 %	227
Σ		428	187	615



Já ao realizar a normalização dos dados o valor de falso negativo cai mais de 10 p.p proporcional a amostragem.



Considero importante para os modelos a normalização e amostragem de dados com valores muito próximos de classificação, ainda mais quando pensando no contexto proposto de saúde/doença, mesmo que os modelos sejam treináveis, é compreensível um olhar humano no impacto dessa margem de erro, aqui, poderia significar 10% a menos de possíveis mortes ocasionadas pela diabetes com diagnóstico errado.

LINK PARA O ARQUIVO DO PROJETO DO ORANGE E DOS DADOS UTILIZADOS

Insira os links para os arquivos.

Dados: <https://www.kaggle.com/datasets/akshaydattatraykhare/diabetes-dataset>

Orange: <https://drive.google.com/file/d/1sMGi-thQMKbHCz9TrgMnxhwVeF-K-3Eh/view?usp=sharing>