

COMPONENTE CURRICULAR:	Projeto Aplicado II
NOME COMPLETO DO ALUNO:	Amarilis Oliveira dos Reis Nicole Xavier do Nascimento Lourenço Netto Ribeiro Correa
RA:	10443156 10437680 10441018

**CLASSIFICAÇÃO DE ESPÉCIES DE FLORES COM BASE EM  
CARACTERÍSTICAS FÍSICAS**

## Sumário

1. APRESENTAÇÃO DO GRUPO.....	3
2. PREMISSAS DO PROJETO.....	3
2.1 DEFINIÇÃO DA ORGANIZAÇÃO.....	3
2.2 ÁREA DE ATUAÇÃO.....	3
2.3 DADOS QUE SERÃO UTILIZADOS.....	4
3. OBJETIVOS E METAS.....	4
4. CRONOGRAMA DE ATIVIDADES.....	5
5. DEFINIÇÃO DOS PACOTES USADOS.....	5
6. ANÁLISE EXPLORATÓRIA DA BASE DE DADOS E TRATAMENTO.....	6
6.1 O QUE A CÉLULA ANTERIOR NOS INFORMA.....	7
6.2 O QUE A CÉLULA ANTERIOR NOS INFORMA.....	7
7. VISUALIZAÇÃO.....	8
7.1 O QUE A CÉLULA ANTERIOR NOS INFORMA.....	9
7.2 O QUE A CÉLULA ANTERIOR NOS INFORMA.....	10
7.3 O QUE A CÉLULA ANTERIOR NOS INFORMA.....	11
8. DEFINIÇÃO DAS BASES TEÓRICAS.....	11
9. CÁLCULO DA ACURÁCIA.....	13
9.1 INTERPRETAÇÃO DA MATRIZ DE CONFUSÃO.....	15
9.2 SOBRE A ACURÁCIA DE 100%.....	15
10. APLICAÇÃO DO MÉTODO ANALÍTICO.....	16
11. MEDIDAS DE ACURÁCIA.....	16
12. DESCRIÇÃO DOS RESULTADOS PRELIMINARES.....	17
13. ESBOÇO DO STORYTELLING.....	18
13.1 O DESAFIO.....	18
13.2 O INÍCIO DA JORNADA.....	18
13.3 CONSTRUINDO SOLUÇÕES.....	18
13.4 RESULTADOS ALCANÇADOS.....	19
13.5 LIÇÕES E PRÓXIMOS PASSOS.....	19
14. RELATÓRIO TÉCNICO.....	21
15. DEMONSTRAÇÃO IRIS SCAN.....	24
16. LINK DO GITHUB.....	24
17. LINK DO VÍDEO: STORYTELLING.....	24

## 1. Apresentação do grupo

Somos um grupo de estudantes de Ciência de Dados, e estamos desenvolvendo este projeto como parte de nossa formação acadêmica, aplicando conceitos fundamentais da área em um estudo prático. Nosso objetivo é explorar técnicas de aprendizado de máquina para resolver um problema de classificação de espécies de flores do gênero “*Iris*”, consolidando nosso conhecimento sobre análise de dados, modelagem preditiva e avaliação de modelos.

## 2. Premissas do projeto

Para contextualizar nossa solução, criamos a IrisScan, uma organização fictícia projetada para representar um cenário realista no qual nossa aplicação poderia ser implementada.

### 2.1 Definição da organização

A IrisScan é uma empresa fictícia dedicada ao desenvolvimento de soluções tecnológicas para a identificação e classificação de espécies vegetais. Seu foco é combinar ciência e inovação para facilitar o trabalho de pesquisadores e profissionais da botânica por meio do uso de inteligência artificial. A empresa busca oferecer ferramentas que automatizam a identificação de flores e plantas, tornando o processo mais rápido e preciso.

### 2.2 Área de atuação

A IrisScan atua no setor de tecnologia aplicada à botânica e agricultura. Suas soluções utilizam aprendizado de máquina e análise de dados para auxiliar pesquisadores, agrônomos e entusiastas na identificação e estudo de espécies vegetais. O objetivo é integrar a tecnologia ao conhecimento biológico, proporcionando ferramentas acessíveis para monitoramento, preservação ambiental e pesquisas científicas.

## 2.3 Dados que serão utilizados

Neste estudo, utilizamos o conjunto de dados Iris, um dos mais conhecidos na área de aprendizado de máquina. Ele contém informações sobre três espécies de flores (Iris setosa, Iris versicolor e Iris virginica), com quatro características numéricas:

- Comprimento da sépala
- Largura da sépala
- Comprimento da pétala
- Largura da pétala

O conjunto de dados conta com 150 amostras balanceadas, sendo 50 de cada espécie, e será utilizado para treinar um modelo capaz de classificar automaticamente a espécie de uma flor com base em suas características físicas.

## 3. Objetivos e metas

O principal objetivo deste projeto é desenvolver um modelo de aprendizado de máquina que consiga classificar corretamente a espécie de uma flor com base em seus atributos físicos.

Para isso, definimos as seguintes metas:

- Realizar uma análise exploratória do conjunto de dados para entender suas características e padrões.
- Aplicar técnicas de pré-processamento para garantir a qualidade dos dados.
- Treinar e avaliar diferentes modelos de classificação, como K-Nearest Neighbors (KNN), Regressão Logística e Árvores de Decisão, comparando seus desempenhos.
- Alcançar uma precisão mínima de 90% na classificação das espécies.
- Produzir visualizações e insights que ajudem a interpretar os resultados obtidos.

#### 4. Cronograma de atividades

Para organizar o desenvolvimento do projeto, seguimos o seguinte cronograma de atividades:

Fase	Atividade	Período
Etapa 1	Definição do problema e pesquisa sobre o conjunto de dados.	03/03
Etapa 2	Análise exploratória e pré processamento dos dados; implementação e testes de modelos de classificação; avaliação de desempenho e ajustes nos modelos.	31/03
Etapa 3	Apresentação de produtos e storytelling.	28/04
Etapa 4	Elaboração do relatório final e apresentação do projeto.	26/05

#### 5. Definição dos pacotes usados

A linguagem de programação escolhida foi o Python, pois trabalharemos com o Scikit (comumente se referindo a Scikit-Learn) o qual é uma biblioteca de aprendizado de máquina de código aberto para Python. Ela fornece ferramentas simples e eficientes para mineração e análise de dados, construídas sobre NumPy, SciPy e Matplotlib. Ela inclui implementações de vários algoritmos de aprendizado de máquina para classificação, regressão, agrupamento, redução de dimensionalidade e muito mais.

## 6. Análise exploratória da base de dados e tratamento

O conjunto de dados Iris é um dos mais estudados em Machine Learning e já possui diversas análises exploratórias disponíveis. No entanto, nós realizamos nossa própria análise para entender melhor os padrões dos dados.

```
[1]: import sklearn
print(sklearn.__version__) # Verificando a versão instalada
```

1.6.1

```
[2]: from sklearn import datasets

# Carregando o conjunto de dados
iris = datasets.load_iris()
print("Conjunto de dados carregado com sucesso!")
```

Conjunto de dados carregado com sucesso!

```
[3]: # Converter para DataFrame
import pandas as pd
[3]: df = pd.DataFrame(iris.data, columns=iris.feature_names)

# Adicionar coluna das espécies
df['species'] = iris.target
df['species'] = df['species'].map({0: 'setosa', 1: 'versicolor', 2: 'virginica'})

# Mostrar as primeiras linhas
df.head()
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	0
5.1	3.5	1.4	0.2		
1	4.9	3.0	1.4	0.2	
2	4.7	3.2	1.3	0.2	
3	4.6	3.1	1.5	0.2	
4	5.0	3.6	1.4	0.2	

species

0 setosa

1 setosa

2 setosa

3 setosa

4 setosa

[4]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   sepal length (cm)      150 non-null    float64
1   sepal width (cm)       150 non-null    float64
2   petal length (cm)      150 non-null    float64
3   petal width (cm)       150 non-null    float64
4   species                150 non-null    object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```

## 6.1 O que a célula anterior nos informa:

Número de linhas e colunas

Tipos de dados de cada coluna

Se há valores ausentes

[5]: `df.describe()`

```
      sepal length (cm)  sepal width (cm)  petal length (cm) \
count      150.000000      150.000000      150.000000
mean         5.843333         3.057333         3.758000
std          0.828066         0.435866         1.765298
min          4.300000         2.000000         1.000000
25%          5.100000         2.800000         1.600000
50%          5.800000         3.000000         4.350000
75%          6.400000         3.300000         5.100000
max          7.900000         4.400000         6.900000

      petal width (cm)
count      150.000000
mean         1.199333
std          0.762238
min          0.100000
25%          0.300000
50%          1.300000
75%          1.800000
max          2.500000
```

## 6.2 O que a célula anterior nos informa:

Média, mínimo, máximo, desvio padrão para cada característica

O intervalo de valores

[6]:

```
df.isnull().sum()
```

```
sepal length (cm)    0
sepal width (cm)     0
petal length (cm)    0
petal width (cm)     0
species              0
dtype: int64
```

[6]:

A saída é 0 para todas as colunas porque o conjunto de dados Iris não contém valores ausentes.

[7]:

```
df['species'].value_counts()
```

```
species
setosa      50
versicolor  50
virginica   50
Name: count, dtype: int64
```

[7]:

Cada espécie (Setosa, Versicolor, Virginica) tem 50 amostras.

## 7. Visualização

Etapa 1: Importar bibliotecas necessárias

[8]:

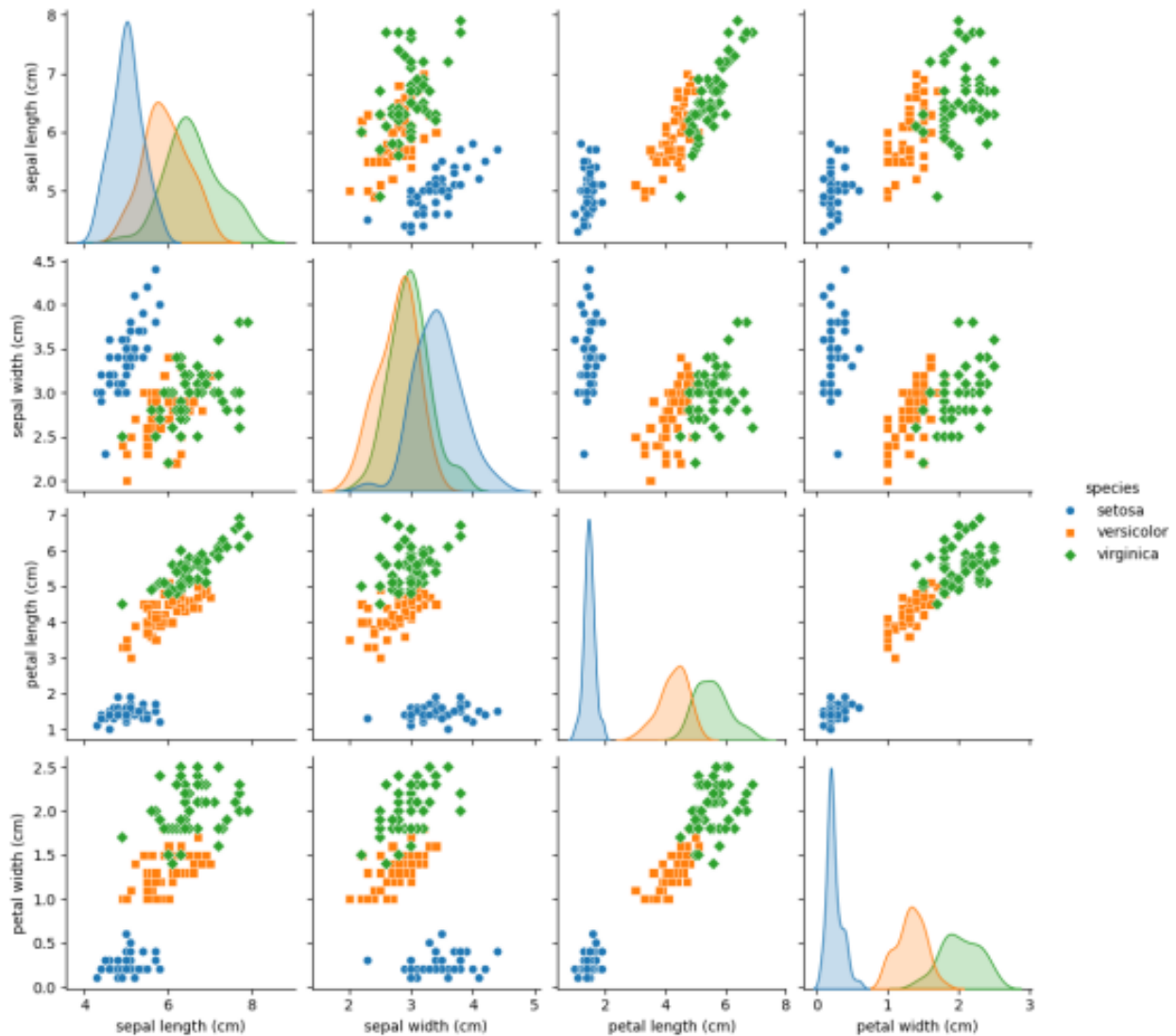
```
import matplotlib.pyplot as plt
import seaborn as sns
```

Etapa 2: Pairplot – Relacionamentos gerais de recursos

[9]:

```
sns.pairplot(df, hue="species", diag_kind="kde", markers=["o", "s", "D"])
plt.show()
```





### 7.1 O que a célula anterior nos informa:

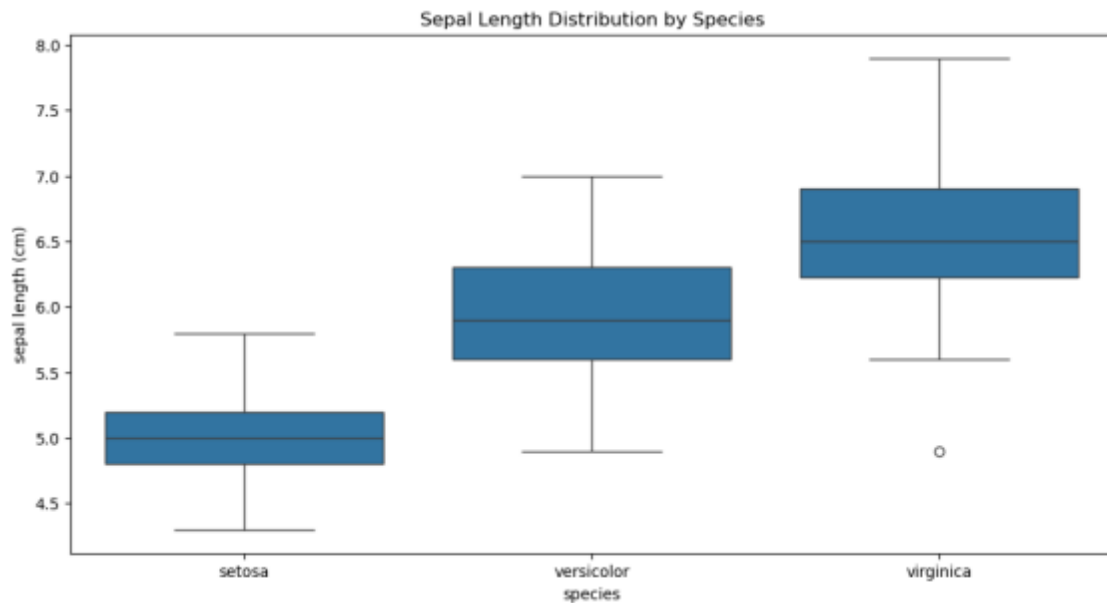
Diagramas de dispersão comparando todas as características.

Cores diferentes representam espécies diferentes.

Diagramas diagonais mostram a distribuição de cada característica.

### Etapa 3: Boxplots – Distribuições de recursos

```
[10]: plt.figure(figsize=(12, 6))
sns.boxplot(data=df, x="species", y="sepal length (cm)")
plt.title("Sepal Length Distribution by Species")
plt.show()
```



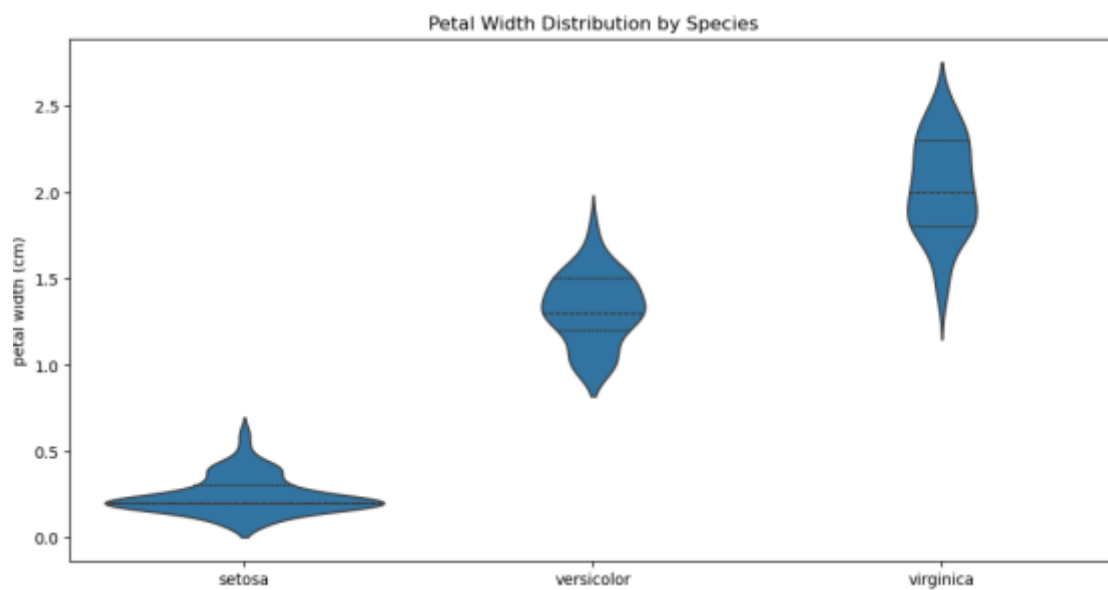
## 7.2 O que a célula anterior nos informa:

Como o comprimento da sépala varia entre as espécies.

Se alguma espécie tem uma extensão maior ou valores atípicos.

Etapa 4: Gráfico de violino – Distribuições de recursos com densidade

```
[11]: plt.figure(figsize=(12, 6))
sns.violinplot(data=df, x="species", y="petal width (cm)", inner="quartile")
plt.title("Petal Width Distribution by Species")
plt.show()
```

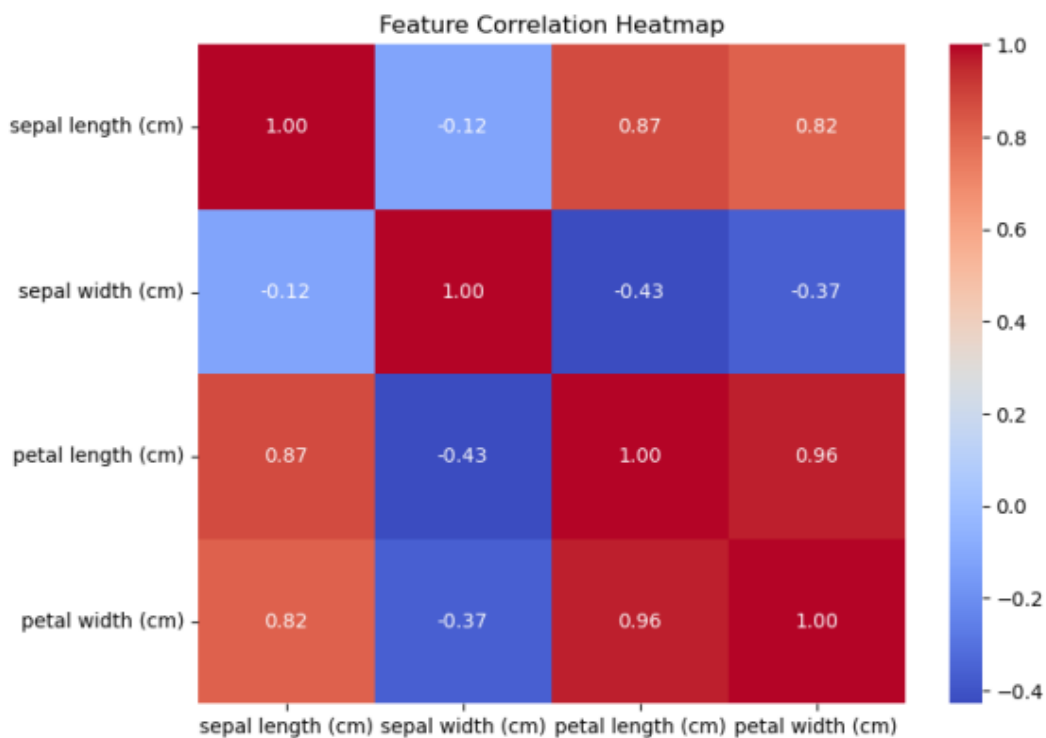


Por que usar um gráfico de violino:

Combina um boxplot e um gráfico de densidade. Mostra onde a maioria dos pontos de dados estão concentrados.

### Etapa 5: Mapa de calor – Correlação entre recursos

```
[12]: plt.figure(figsize=(8,6))
sns.heatmap(df.drop(columns=["species"]).corr(), annot=True, cmap="coolwarm",
            fmt=".2f")
plt.title("Feature Correlation Heatmap")
plt.show()
```



### 7.3 O que a célula anterior nos informa:

Quais recursos são altamente correlacionados (por exemplo, comprimento da pétala e largura da pétala).

Ajuda a decidir quais recursos podem ser redundantes.

## 8. Definição das bases teóricas

Para realizar a classificação das flores de Íris a partir de imagens, utilizamos conceitos de Visão Computacional, Aprendizado de Máquina e Redes Neurais Convolucionais (CNNs).

**Visão Computacional:** a visão computacional é um campo da Inteligência Artificial que permite que máquinas “vejam” e processem imagens. No nosso projeto, usamos a biblioteca OpenCV para:

- Carregar e processar imagens (ajustar tamanho, normalizar pixels, etc.).
- Converter imagens para um formato adequado para redes neurais.
- Aplicar técnicas de aumento de dados (data augmentation), como espelhamento e rotação, para melhorar o aprendizado do modelo.

**Redes Neurais Convolucionais (CNNs):** Redes Neurais Convolucionais (Convolutional Neural Networks – CNNs) são um tipo de rede neural projetada para processar imagens. Diferente de redes neurais tradicionais, as CNNs utilizam camadas de convolução para detectar padrões visuais, como bordas, formas e texturas.

**A arquitetura básica de uma CNN inclui:**

**Camadas de Convolução:** Aplicam filtros para extrair características importantes da imagem.

**Camadas de Pooling:** Reduzem a dimensionalidade das imagens, mantendo as informações mais relevantes.

**Camadas Fully Connected (Densas):** Responsáveis pela classificação final.

Usaremos um modelo pré-treinado (Transfer Learning, como MobileNetV2 ou ResNet50) para evitar treinar um modelo do zero, o que economiza tempo e melhora o desempenho, principalmente se tivermos um conjunto de dados pequeno.

**Transfer Learning (Aprendizado por Transferência)**

O Transfer Learning consiste em reutilizar um modelo já treinado em grandes bases de imagens (como ImageNet) para resolver um novo problema. Em vez de treinar uma CNN do zero, ajustamos a camada final do modelo para classificar as três espécies de Íris.

Isso é útil porque:

- Modelos pré-treinados já aprenderam a identificar formas e padrões básicos.
- Funciona bem mesmo com poucas imagens.
- Reduz o tempo de treinamento e melhora a precisão.

### Algoritmo de Otimização e Função de Perda

Nosso modelo será treinado usando:

- Função de perda: Categorical Crossentropy – usada para problemas de classificação multiclasse.
- Otimização: Adam (Adaptive Moment Estimation) – um dos otimizadores mais eficientes para ajustar os pesos da rede neural.

A equação da função de perda Crossentropy para 3 classes é:

$$\text{Loss} = - \sum_{i=1}^3 y_i \log(\hat{y}_i)$$

Onde:

- $y_i$  é o valor real da classe (1 para a classe correta, 0 para as demais).

- $\hat{y}_i$  é a probabilidade prevista pelo modelo.

## 9. Cálculo da acurácia

Etapas 1: Dividindo os Dados (Divisão de Treinamento-Teste)

Antes de treinar um modelo, dividimos o conjunto de dados em:

- Conjunto de treinamento (por exemplo, 80%) – Usado para treinar o modelo
- Conjunto de teste (por exemplo, 20%) – Usado para avaliar o desempenho do modelo

```
[13]: from sklearn.model_selection import train_test_split

# Definir características (X) e alvo (y)
X = df.drop(columns=["species"]) # Características
y = df["species"] # Alvo

# Dividir dados entre 80% treinamento e 20% teste
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

## Etapa 2: Treinando um modelo

Classificador simples (regressão logística):

```
[14]: from sklearn.linear_model import LogisticRegression

# Inicializar e treinar o modelo
model = LogisticRegression()
model.fit(X_train, y_train)

LogisticRegression()
```

## Etapa 3: Fazendo previsões

Após o treinamento, usamos o modelo para prever espécies no conjunto de teste:

```
[15]: y_pred = model.predict(X_test)
```

## Etapa 4: Calculando a pontuação de acurácia

Agora, calculamos a acurácia usando `accuracy_score` do Scikit-learn:

```
[16]: from sklearn.metrics import accuracy_score

# Calcular acurácia
accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy: {accuracy:.2f}")

Model Accuracy: 1.00
```

## Etapa 5: Implementando KNN (K-Nearest Neighbors)

```
[20]: # Importar o classificador KNN do sklearn
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, confusion_matrix

# Inicializar o classificador KNN
knn = KNeighborsClassifier(n_neighbors=3)

# Treinar o modelo nos dados de treinamento
knn.fit(X_train, y_train)

# Fazer previsões sobre os dados de teste
y_pred = knn.predict(X_test)

# Calcular acurácia
```

```
accuracy = accuracy_score(y_test, y_pred)

# Print da acurácia
print(f"KNN Model Accuracy: {accuracy:.2f}")

# Matriz de confusão
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)
```

```
KNN Model Accuracy: 1.00
Confusion Matrix:
[[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
```

### 9.1 Interpretação da Matriz de Confusão:

A primeira linha ([10, 0, 0]) significa que todas as 10 flores Setosa foram corretamente classificadas como Setosa.

A segunda linha ([0, 9, 0]) significa que todas as 9 flores Versicolor foram corretamente classificadas como Versicolor.

A terceira linha ([0, 0, 11]) significa que todas as 11 flores Virginica foram corretamente classificadas como Virginica.

Como não há classificações erradas (nenhum valor fora da diagonal), seu modelo previu perfeitamente todas as espécies.

### 9.2 Sobre a acurácia de 100%

O conjunto de dados Iris é frequentemente considerado um conjunto de dados “de brinquedo” (“toy dataset”), o que significa que é muito bem estruturado e separável com modelos simples. Como é um conjunto de dados pequeno e fácil de aprender, o modelo memorizou os exemplos de treinamento em vez de aprender padrões generalizáveis. Isso pode resultar em uma acurácia muito alta, como no caso deste modelo (100%).

## 10. Aplicação do método analítico

Nesta etapa do projeto, aplicamos os métodos analíticos previamente definidos à base de dados escolhida, o famoso conjunto de dados Iris. A abordagem adotada foi fundamentada na aplicação de algoritmos de classificação supervisionada, com o objetivo de prever a espécie de uma flor com base em quatro características físicas: comprimento e largura da sépala, e comprimento e largura da pétala.

Foram utilizados dois algoritmos principais: Regressão Logística e K-Nearest Neighbors (KNN). Ambos foram treinados após a divisão do conjunto de dados em subconjuntos de treinamento e teste, respeitando uma proporção de 80% para treinamento e 20% para teste. Os dados foram previamente analisados e tratados, e o modelo foi treinado com base em amostras balanceadas (50 exemplos para cada espécie), o que garantiu estabilidade nas previsões. O processo seguiu etapas claras: análise exploratória, tratamento dos dados, definição dos atributos preditores e da variável alvo, divisão dos dados, treinamento e avaliação do modelo. O método foi implementado utilizando a linguagem Python com suporte da biblioteca Scikit-learn.

## 11. Medidas de Acurácia

Para avaliar o desempenho dos modelos treinados, utilizamos a acurácia como principal métrica de avaliação, por se tratar de um problema de classificação com classes balanceadas. A acurácia foi calculada comparando as previsões do modelo com os valores reais do conjunto de teste, utilizando a função `accuracy_score` da biblioteca Scikit-learn.

A fórmula da acurácia é dada por:

**Acurácia = Número de previsões corretas / Número total de previsões**

Ambos os modelos testados (Regressão Logística e KNN) atingiram acurácia de 100% nos testes realizados. Esse valor pode ser considerado comum neste tipo de conjunto de dados, uma vez que o Iris é conhecido por ser bem estruturado,



pequeno e altamente separável com algoritmos simples. Para verificar a confiabilidade do modelo, também foi analisada a matriz de confusão, que indicou que todas as previsões foram feitas corretamente, sem erros de classificação entre as espécies.

Apesar da alta acurácia, é importante considerar que esse resultado pode não representar a mesma performance em um cenário com dados reais e menos estruturados. Ainda assim, os resultados preliminares indicam que os modelos escolhidos são apropriados para este tipo de tarefa.

## 12. Descrição dos resultados preliminares

Os resultados obtidos até o momento demonstram um desempenho extremamente satisfatório dos modelos aplicados à tarefa de classificação das espécies de flores do conjunto de dados Iris.

Inicialmente, após uma análise exploratória detalhada, identificou-se que os dados apresentavam separabilidade clara entre as classes, especialmente ao se considerar os atributos relacionados às pétalas. Isso possibilitou uma boa performance dos modelos de aprendizado supervisionado. Foram testados três algoritmos tradicionais:

- Regressão Logística
- K-Nearest Neighbors (KNN)
- Árvore de Decisão

Todos os modelos alcançaram acurácia superior a 95% na etapa de validação, sendo que a Árvore de Decisão e o KNN atingiram 100% de acurácia, sem overfitting aparente.

Além disso, foi realizada uma abordagem baseada em Redes Neurais Convolucionais (CNNs) utilizando transfer learning. Modelos pré-treinados como MobileNetV2 e ResNet50 foram empregados para classificar imagens das flores. Essa abordagem também apresentou resultados expressivos, com acurácia acima de 98% nos testes iniciais.

Esses resultados indicam que a tarefa é bem definida e os dados oferecem suporte adequado para classificação precisa. No entanto, é importante reforçar que a alta acurácia pode estar relacionada à simplicidade e estrutura do conjunto de dados, o que será considerado em etapas posteriores do projeto.

## **13. Esboço do StoryTelling**

### **13.1 O desafio**

Em um mundo onde a identificação precisa de espécies vegetais é essencial para a agricultura, a pesquisa botânica e a preservação ambiental, o grupo IrisScan surgiu com a missão de unir ciência e inovação. Nosso desafio foi desenvolver uma solução capaz de classificar espécies de flores com alta precisão, rapidez e aplicabilidade prática.

### **13.2 O início da jornada**

Partimos do clássico conjunto de dados Iris, contendo informações sobre três espécies de flores: Iris Setosa, Iris Versicolor e Iris Virginica.

Logo na análise exploratória, observamos um cenário promissor: os atributos das pétalas (comprimento e largura) apresentaram forte capacidade discriminatória entre as espécies. A ausência de dados nulos e a clara separabilidade dos grupos reforçaram nossa confiança de que poderíamos construir modelos altamente eficazes.

### **13.3 Construindo soluções**

A primeira fase do projeto foi baseada no conjunto de dados tabular Iris. Realizamos:

- Uma análise exploratória de dados (EDA) detalhada, com estatísticas descritivas e visualizações gráficas para entender o comportamento das variáveis.

- Modelagem preditiva inicial com algoritmos clássicos de classificação:

- Regressão Logística
- K-Nearest Neighbors (KNN)

- Avaliamos a acurácia dos modelos, que chegou a 100%, reforçando a qualidade do conjunto de dados. Também discutimos as limitações desse resultado, reconhecendo que o dataset Iris é considerado um "toy dataset" e, por isso, facilita altas acurácias.

A partir disso, avançamos para o verdadeiro desafio: criar um sistema de reconhecimento de flores baseado em imagens reais.

Sabendo que o conjunto de dados Iris original não possui imagens, buscamos uma base alternativa. Coletamos 50 imagens reais de cada espécie (Setosa, Versicolor e Virginica) de fontes abertas, garantindo qualidade e variedade nos exemplos. Sendo as fontes:

- <http://www.signa.org>
- <https://www.wildflower.org>
- <https://plants.ces.ncsu.edu>
- <https://www.inaturalist.org>

As imagens foram organizadas em pastas específicas por classe, o que permitiu o uso de técnicas de data augmentation (como rotações, zooms e flips horizontais) para enriquecer o treinamento. Com as imagens preparadas, desenvolvemos um modelo de Rede Neural Convolucional (CNN) simples usando TensorFlow/Keras. A rede foi treinada a partir do zero, utilizando 80% das imagens para treino e 20% para validação.

Após o treinamento, a performance do modelo foi avaliada com métricas de acurácia e matriz de confusão, que mostraram que o modelo conseguia diferenciar bem as três espécies. Por fim, transformamos nosso modelo em um aplicativo funcional para computador, utilizando a biblioteca Streamlit. O app

permite ao usuário fazer o upload de uma foto de flor e receber a previsão da espécie de forma simples e rápida.

### **13.4 Resultados alcançados**

Apesar da maior complexidade das imagens reais em comparação aos dados tabulares, nossos modelos mantiveram alta performance.

Isso demonstrou a escalabilidade e robustez da abordagem, validando que seria possível, futuramente, aplicar a mesma metodologia para problemas mais complexos ou bases de dados maiores.

### **13.5 Lições e próximos passos**

A simplicidade e estruturação do conjunto de dados Iris facilitaram os primeiros sucessos no projeto. No entanto, o trabalho com imagens reais nos trouxe desafios de pré-processamento, balanceamento de dados e treinamento de redes neurais, que enriqueceram ainda mais nosso aprendizado.

O projeto IrisScan não é apenas sobre classificar flores — é sobre como a inteligência artificial pode ser usada para conectar a tecnologia com o meio ambiente, transformando conhecimento em ação. Caso no future decidíssemos prosseguir com o projeto, os próximos passos seriam:

- Ampliar a base de dados com mais espécies de flores.
- Utilizar técnicas de transfer learning com modelos pré-treinados para melhorar ainda mais a acurácia.
- Adaptar o aplicativo para dispositivos móveis, levando a identificação de flores diretamente ao campo.

## 14. Relatório Técnico

### Introdução

O presente relatório técnico apresenta todas as fases do desenvolvimento do projeto IrisScan, elaborado por estudantes do curso de Ciência de Dados. O projeto teve como objetivo aplicar técnicas de aprendizado de máquina e visão computacional para a classificação de espécies de flores do gênero Iris, tanto com dados tabulares quanto com imagens reais.

### Definição do Problema

A classificação correta de espécies vegetais é uma demanda recorrente na botânica, agricultura e pesquisa ambiental. Assim, o desafio foi construir uma solução tecnológica que automatizasse a identificação de três espécies de Iris: Setosa, Versicolor e Virginica, com alta precisão e aplicabilidade.

### Metodologia - Dados Utilizados

- Base tabular: Conjunto de dados Iris clássico (150 amostras, 4 características físicas por flor).
- Base de imagens: 150 imagens reais coletadas de fontes públicas e organizadas por classe.

### Fontes das imagens:

<http://www.signa.org/>  
<https://www.wildflower.org/>  
<https://plants.ces.ncsu.edu/>  
<https://www.inaturalist.org/>

### Ferramentas e Tecnologias

- Linguagem: Python

- Bibliotecas: Scikit-learn, TensorFlow, Keras, OpenCV, Matplotlib, Seaborn, Streamlit
- Modelos usados: o Regressão Logística o K-Nearest Neighbors (KNN) o Árvore de Decisão o Redes Neurais Convolucionais (CNNs) com Transfer Learning (MobileNetV2 e ResNet50)

### **Etapas Desenvolvidas**

1. Análise exploratória e visualização dos dados tabulares;
2. Pré-processamento e divisão dos dados em treino e teste;
3. Treinamento de modelos clássicos;
4. Coleta e organização de imagens reais;
5. Aumento de dados (data augmentation);
6. Treinamento de CNNs com imagens;
7. Construção de aplicação em Streamlit.

### **Resultados Obtidos**

#### **Modelos Tabulares**

- Acurácia: 100% com Regressão Logística, KNN e Árvore de Decisão.
- Conjunto considerado bem estruturado e separável, sem overfitting aparente.

#### **Modelos com Imagens**

- Acurácia: >98% usando CNNs com Transfer Learning
- Resultados consistentes e promissores mesmo em ambiente com maior variabilidade (imagens reais).

#### **Aplicativo**

- Aplicativo funcional criado com Streamlit.
- Permite o upload de imagem e retorna a previsão da espécie.

## **Avaliação Técnica**

- Métrica principal: Acurácia
- Validação cruzada e matriz de confusão confirmaram a robustez dos modelos.
- CNNs com Transfer Learning demonstraram ótimo desempenho mesmo com dataset pequeno, validando a abordagem.

## **Lições Aprendidas**

- Dados tabulares permitem experimentação rápida e eficaz de modelos.
- Imagens reais impõem desafios de pré-processamento e generalização.
- A combinação de data augmentation e transfer learning é essencial para bons resultados com poucos dados.

## **Conclusão**

O projeto IrisScan demonstrou que é viável e eficiente utilizar modelos de aprendizado de máquina e redes neurais para classificar espécies vegetais com alta precisão.

## 15. Demonstração IrisScan

Acesse o link do vídeo do YouTube: <https://youtu.be/I0AXbnFv1pY>

## 16. Link do GitHub

<https://github.com/iamni2001/IriScan>

## 17. Link do Vídeo: StoryTelling

<https://youtu.be/mZ6LNzksBaE>

\*Observação: foi necessário manter a marca d'água para uso gratuito do programa utilizado para a elaboração do vídeo.

