2. Model Selec	ation: able on Kaggle and provide a reference to your dataset in your report. Examine the dataset to understand its structure, including the types of variables and possible missing values. Clean the data if necessary, ensuring it's ready for analysis. Ition and Training: Ification algorithm. Justify your choice by considering the nature of the dataset and the problem at hand. Train your model using the training portion of the dataset, ensuring to split data appropriately between training and testing sets.
3. Model Evaluate the model's4. Report Writing	performance using appropriate metrics such as accuracy, precision, recall, and F1-score. Discuss the significance of each metric in assessing the model's performance. Use confusion matrices to better understand the model's classification accuracy.
Problem	Statement: Independent of the structured report. The report should include an introduction to the problem, methodology used, results with detailed explanations, and a conclusion. Visual aids such as graphs and tables are encouraged to support your analysis. Statement: Independent of the state of building a classification model to predict whether a mushroom is edible or poisonous based on its characteristics. We used the "Mushroom" dataset, which contains various features of mushrooms and their corresponding is.
<pre>import pandas as po import numpy as np from sklearn.model_</pre>	
<pre>from sklearn.metric import matplotlib.p import seaborn as s import pandas as po import numpy as np</pre>	
	C:/Users/Psycho Doc/OneDrive/Desktop/Machine Learning/machine_learning/Mush Classify/mushrooms.csv") Formation about the dataset
class 'pandas.core. RangeIndex: 8124 ent Data columns (total	ded and initial exploration complete.") frame.DataFrame'> ries, 0 to 8123 23 columns):
# Column 0 class 1 cap-shape 2 cap-surface 3 cap-color 4 bruises 5 odor	Non-Null Count Dtype
13 stalk-surface-b	8124 non-null object bove-ring 8124 non-null object elow-ring 8124 non-null object
<pre>14 stalk-color-abo 15 stalk-color-bel 16 veil-type 17 veil-color 18 ring-number 19 ring-type 20 spore-print-col 21 population</pre>	ow-ring 8124 non-null object
22 habitat dtypes: object(23) memory usage: 1.4+ M None class cap-shape ca p x L e x 2 e b 3 p x	8124 non-null object B p-surface cap-color bruises odor gill-attachment \
gill-spacing gill- c c c c c c c c c c c c c	s g f n f size gill-color stalk-surface-below-ring \ n k s b k s b n s c n n s b k s
	ring stalk-color-below-ring veil-type veil-color \ w
ring-number ring-t o o o o o o o o o o o t o o	ype spore-print-color population habitat p k s u p n n g p n n m p k s u p n n m p k s u e n a g
class cap-shape cap-surface cap-color oruises odor gill-attachment gill-spacing	
gill-size gill-color stalk-shape stalk-root stalk-surface-above- stalk-surface-below- stalk-color-above-ri	ring 0 ng 0
veil-type veil-color ring-number ring-type spore-print-color copulation nabitat dtype: int64	
The abo	ove output is obtained after the first step that is data Preparation which has returned the output as the dataset contains 8,124 entries with 23 s. The of object dtype, which means they contain categorical data.
# Reload and inspectimport pandas as po	
<pre># Load the dataset df = pd.read_csv('C import seaborn as s import matplotlib.p # Display the first df_head = df.head()</pre>	byplot as plt t few rows
<pre># Check the class of class_counts = df[' # Visualize the closes.set(style="white plt.figure(figsize="white")</pre>	distribution class'].value_counts() ass distribution tegrid") =(8, 5))
<pre>plt.title('Class Di plt.xlabel('Class (plt.ylabel('Count') plt.show() # Check for unique unique_values = df</pre>	values in features for potential visualization nunique()
<pre>selected_features = for feature in sele plt.figure(figs sns.countplot(c) plt.title('Dist plt.xlabel(feat plt.ylabel('Counter)</pre>	size=(10, 5)) data=df, x=feature, hue='class', palette='husl') tribution of ' + feature + ' by Class') ture) unt')
<pre>plt.legend(tit] plt.xticks(rota plt.show()</pre>	le='Class', labels=['Edible', 'Poisonous']) ation=45) ataset head and unique values
4000 3500	Class Distribution (Edible vs Poisonous)
3000 ——————————————————————————————————	
1500 —— 1000 —— 500 ——	
1200	Class (e=edible, p=poisonous) Distribution of cap-color by Class Class
1000 - 800 -	Edible Poisonous
400 -	
200	4 4 6 8 8 8 5 C (cap-color
3500	Distribution of odor by Class Class Edible Poisonous
2500 2000 1500	
500	
8	odor Distribution of population by Class Class
2500	Edible Poisonous
1500 1000	
0 %	population c
0 p x 1 e x 2 e b 3 p x 4 e x	p-surface cap-color bruises odor gill-attachment \ s
1 c 2 c 3 c 4 w stalk-color-above-	n
2 3 4 ring-number ring-t 0	w w p w w w p w w w p w w w p w wype spore-print-color population habitat p k s u p n n g p n n m
3 o 4 o [5 rows x 23 columns class cap-shape cap-color	p k s u e n a g 2 6 4 10
bruises odor gill-attachment gill-spacing gill-size gill-color stalk-shape stalk-root	2 9 2 2 2 12 2 5
stalk-surface-above- stalk-surface-below- stalk-color-above-ri stalk-color-below-ri veil-type veil-color ring-number ring-type	ring 4 .ng 9
spore-print-color population habitat dtype: int64	graph shows the Edible vs Poisonous Mushrooms in dataset, along with other features.
	olumns: e.fit_transform(df[column]) anto features (X) and target (y)
<pre># Split the data in X_train, X_test, y_ # Train the Decision</pre>	nTreeClassifier(random_state=42)
<pre>precision = precisi recall = recall_sco </pre>	tion metrics y_score(y_test, y_pred) tion_score(y_test, y_pred) tore(y_test, y_pred)
<pre>f1 = f1_score(y_tes print("Model Evalua print(f"Accuracy: { print(f"Precision: print(f"Recall: {re print(f"F1-score: { # Create and plot of</pre>	ation Metrics:") {accuracy:.4f}") {precision:.4f}") ecall:.4f}") {f1:.4f}")
<pre>cm = confusion_matr plt.figure(figsize=</pre>	rix(y_test, y_pred) =(8, 6)) not=True, fmt='d', cmap='Blues') on Matrix') ted') ')
<pre># Feature importance feature_importance feature_importance plt.figure(figsize= sns.barplot(x='importance) plt.title('Top 10 Figure)</pre>	<pre>= pd.DataFrame({'feature': X.columns, 'importance': dt_model.feature_importances_}) = feature_importance.sort_values('importance', ascending=False).head(10) =(10, 6)) ortance', y='feature', data=feature_importance) Feature Importances')</pre>
<pre>plt.xlabel('Importa plt.ylabel('Feature plt.tight_layout() plt.show() plt.savefig('featur plt.close() plt.figure(figsize=</pre>	ance') e') re_importance.png') =(8, 6))
<pre>sns.heatmap(cm, and plt.title('Confusion plt.xlabel('Predict plt.ylabel('Actual' plt.savefig('confus # Create and plot of cm = confusion_matr plt.figure(figsize=</pre>	not=True, fmt='d', cmap='Blues') on Matrix') ted') ') sion_matrix.png') confusion matrix rix(y_test, y_pred) =(8, 6))
	not=True, fmt='d', cmap='Blues') on Matrix') ted') ') sion_matrix.png')
<pre>feature_importance feature_importance plt.figure(figsize= sns.barplot(x='importance)</pre>	<pre>= pd.DataFrame({'feature': X.columns, 'importance': dt_model.feature_importances_}) = feature_importance.sort_values('importance', ascending=False).head(10) =(10, 6)) ortance', y='feature', data=feature_importance) Feature Importances') ance')</pre>
<pre>plt.tight_layout() plt.savefig('featur plt.close()</pre>	re_importance.png') ing, evaluation, and visualization complete.")
Accuracy: 1.0000 Precision: 1.0000 Recall: 1.0000 F1-score: 1.0000	Top 10 Feature Importances
spore-print-o	color -
gill	odor - nises -
Featur Prod	
Featur Prod	O.00 0.05 0.10 0.15 0.20 0.25 0.30 0.35 Importance
hal stalk-sh stalk-color-below- stalk-	uation, and visualization complete. Confusion Matrix
hal stalk-sh stalk-color-below- stalk-	- 800 - 700
hall stalk-sh stalk-color-below- stalk- Model training, eval	843 0 - 800 - 700 - 600 - 500
hal stalk-sh stalk-color-below-stalk-	- 800 - 700 - 600
hall stalk-sh stalk-color-below-stalk-	- 800 - 700 - 600 - 500 - 400 - 300
hall stalk-sh stalk-color-below- stalk- Model training, eval	843 0 -800 -700 -600 -500 -400 -300 -300 -100 -100 -0 Predicted was properly prepared, encoded, and split into training (80%) and testing (20%) sets.
The data Key Hig Model Choice: Performance N Accuracy: 1.00	assume that the state of the st
The data Key Hig Model Choice: Performance Notes and the second of th	was properly prepared, encoded, and split into training (80%) and testing (20%) sets. A Decision Tree was selected for its effectiveness with categorical data and interpretability.
The data Key Hig Model Choice: Performance N Accuracy: 1.00 Precision: 1.000 Recall: 1.0000 F1-score: 1.00 The confusion Conside Overfittin	843 0 782 400 782 400 - 300 - 100 -
The data Key Hig Model training, eval Performance N Accuracy: 1.00 Precision: 1.00 Recall: 1.0000 F1-score: 1.00 The confusion Conside Overfittin Need to d Insights g	843 0 -700 -600 -700 -600 -700 -600 -700 -600 -700

To validate the model's effectiveness further, future work should focus on implementing cross-validation techniques and testing the model against new, unseen datasets. Additionally, exploring alternative algorithms and leveraging expert knowledge for feature selection may

enhance the model's robustness and interpretability. As the domain of mushroom classification continues to evolve, collecting more diverse data and rigorously assessing the model will be imperative to ensure that it can reliably support critical applications in the real world.

Lab 2: Classification Model Development