

Support Vector Machines (SVM)

By Josh Houlding

Another machine learning algorithm is the support vector machine (SVM). It provides resourceful predictions for complex real-word problems, reduces redundant information, and is ideal for smaller datasets. Some real-world examples of when SVMs can be used are bioinformatics, face recognition, speech recognition, or cancer detection.

Part 1

To perform a classification using the support vector machine algorithm, complete the following:

1. Access the "UCI Machine Learning Repository," located in the topic Resources. Note: There are about 120 data sets that are suitable for use in a classification task. For this part of the exercise, you must choose one of these datasets, provided it includes at least 10 attributes and 10,000 instances.
2. Ensure that the datasets are suitable for classification using this method.
3. You may search for data in other repositories, such as Data.gov, Kaggle or Scikit Learn.
4. Examine the repository through which you accessed the dataset and discuss data management measures set in place, such as protecting the privacy of those accessing the site and protecting the intellectual property rights of the data owners/contributors.

Part 2

For your selected dataset, build a classification model as follows:

1. Explain the dataset and the type of information you wish to gain by applying a classification method.
2. Explain what makes SVM algorithm very special and very different from most other machine learning algorithms. Explain how you will be using it in your analysis (list the steps, the intuition behind the mathematical representation, and address its assumptions).
3. Explain the concepts: kernel, hyperplane, and decision boundary, and their role in SVM.
4. Explain the concepts: maximum margin, support vectors, and maximum margin hyperplane, and their role in SVM.
5. Import the necessary libraries, then read the dataset into a data frame and perform initial statistical exploration.

6. Clean the data and address unusual phenomena (e.g., normalization, feature scaling, outliers); use illustrative diagrams and plots and explain them.
7. Formulate two questions that can be answered by applying a classification method using the SVM
8. Split the data into 80% training and 20% testing sets using the train test split class.
9. Use a linear kernel to train the SVM classifier on the training set (e.g., fit the support vector regressor to the dataset). Explain the intuition behind each of the key mathematical steps.
10. Explain the choice of the optimal hyperplane.
11. Make classification predictions.
12. Interpret the results in the context of the questions you asked.
13. Validate your model using a confusion matrix, accuracy score, ROC-AUC curves, and k-fold cross validation. Then, explain the results.
14. Include all mathematical formulas used and graphs representing the final outcomes.

Part 1

Tasks 1.1-1.3

1.1: Access the "UCI Machine Learning Repository," located in the topic Resources. Note: There are about 120 data sets that are suitable for use in a classification task. For this part of the exercise, you must choose one of these datasets, provided it includes at least 10 attributes and 10,000 instances.

1.2: Ensure that the datasets are suitable for classification using this method.

1.3: You may search for data in other repositories, such as Data.gov, Kaggle or Scikit Learn.

Dataset Selected: "Secondary Mushroom" (UCI Machine Learning Repository)

Tagline: Dataset of simulated mushrooms for binary classification into edible and poisonous.

Link: <https://archive.ics.uci.edu/dataset/848/secondary+mushroom+dataset>

Task 1.4

Examine the repository through which you accessed the dataset and discuss data management measures set in place, such as protecting the privacy of those accessing the site and protecting the intellectual property rights of the data owners/contributors.

The UCI Machine Learning Repository site uses https, meaning it encrypts the data sent to the user's browser, meaning that malicious actors cannot eavesdrop on network packets as they are being delivered to the user. The site also enables users to browse and download datasets anonymously.

The "Secondary Mushroom" dataset is licensed under a Creative Commons Attribution 4.0 International (CC BY 4.0) license, which "allows for the sharing and adaptation of the datasets for any purpose, provided that the appropriate credit is given". This means the creators of the dataset allow liberal use of their data, as long as the user credits them properly.

Proper credit includes the following:

- A citation (APA-7 in this case because it is what GCU uses).
- An attribution statement (letting readers know the dataset has a CC BY 4.0 license).
- A link to the original dataset.

Part 2

Task 2.1

Explain the dataset and the type of information you wish to gain by applying a classification method.

This is the information presented in the "Dataset Information" section of the UCI ML Repo page for the dataset:

For what purpose was the dataset created?

Inspired by the Mushroom Data Set of J. Schlimmer: url:<https://archive.ics.uci.edu/ml/datasets/Mushroom>.

Additional Information

The given information is about the Secondary Mushroom Dataset, the Primary Mushroom Dataset used for the simulation and the respective metadata can be found in the zip.

This dataset includes 61069 hypothetical mushrooms with caps based on 173 species (353 mushrooms per species). Each mushroom is identified as definitely edible, definitely poisonous, or of unknown edibility and not recommended (the latter class was combined with the poisonous class).

The related Python project contains a Python module `secondary_data_generation.py` used to generate this data based on `primary_data_edited.csv` also found in the repository. Both nominal and metrical variables are a result of randomization. The simulated and ordered by species version is found in `secondary_data_generated.csv`. The randomly shuffled version is found in `secondary_data_shuffled.csv`.

A classification method can be applied to this dataset to predict if a mushroom is edible or poisonous based on its physical characteristics.

Task 2.2

Explain what makes SVM algorithm very special and very different from most other machine learning algorithms. Explain how you will be using it in your analysis (list the steps, the intuition behind the mathematical representation, and address its assumptions).

Support Vector Machines (SVMs) separate data points into classes by finding the optimal hyperplane that maximizes the distance between classes in an n-dimensional space. Unlike k-nearest neighbors (kNN), for example, SVMs are highly effective for high-dimensional datasets because they handle sparse data and create margins between classes well.

These are the steps an SVM takes to perform classification:

1. Represent data points as vectors in an n-dimensional space, with each feature corresponding to a dimension.
2. Find the optimal hyperplane that maximizes the margin, or the distance between the hyperplane and the nearest data points (support vectors).
3. Find the support vectors, or the data points closest to the decision boundary. These determine the hyperplane's position.
4. If dealing with data that is not linearly separable, the "kernel trick" can be applied to map it into a higher-dimensional space, enabling the SVM to find nonlinear decision boundaries.
5. Formulate the optimization problem (maximize the margin subject to constraints).
6. Solve for Lagrange multipliers using optimization methods.
7. Compute the decision boundary using the support vectors and Lagrange multipliers.
8. Classify new data points based on the sign of the decision function.

- Evaluate the model using common metrics (accuracy, precision, recall, F1 score, etc.) and tune hyperparameters if necessary to improve performance.

SVMs assume that it is optimal to have the largest possible margin, and that support vectors are crucial.

Task 2.3

Explain the concepts: kernel, hyperplane, and decision boundary, and their role in SVM.

A **kernel** is a method that makes it possible to map nonlinear data into a higher-dimensional space, enabling the use of linear classifiers on nonlinear problems. SVMs benefit greatly from kernels because their goal is to find a linear decision boundary in a higher-dimensional space, and kernels provide an efficient way to do this that doesn't involve computing every transformation explicitly.

A **hyperplane** is a line or flat surface that divides an n-dimensional space into two parts. They are fundamental to the operation of SVMs because these algorithms use hyperplanes to make classifications for individual data points.

Finally, a **decision boundary** is the dividing line, curve or surface of a hyperplane that the SVM uses to determine which class a data point should be put in.

Task 2.4

Explain the concepts: maximum margin, support vectors, and maximum margin hyperplane, and their role in SVM.

The **maximum margin** for an SVM is the largest possible distance between the hyperplane and nearest data points (support vectors). In other words, it is part of the optimal solution to the optimization problem an SVM aims to solve.

As previously mentioned, **support vectors** are the data points that are closest to the decision boundary/hyperplane.

Finally, the **maximum margin hyperplane** is the hyperplane associated with the maximum margin, as the name suggests. This is the optimal solution to the SVM optimization problem.

Task 2.5

Import the necessary libraries, then read the dataset into a data frame and perform initial statistical exploration.

Loading the data

In [93]:

```
import pandas as pd

# Load and view data
df = pd.read_csv("mushroom.csv", delimiter=";")
df.head()
```

Out[93]:

	class	cap-diameter	cap-shape	cap-surface	cap-color	does-bruise-or-bleed	gill-attachment	gill-spacing	gill-color	stem-height	...	stem-root	stem-surface	stem-color	veil-type	veil-color	has-ring	ring-type	s
0	p	15.26	x	g	o	f	e	NaN	w	16.95	...	s	y	w	u	w	t	g	
1	p	16.60	x	g	o	f	e	NaN	w	17.99	...	s	y	w	u	w	t	g	
2	p	14.07	x	g	o	f	e	NaN	w	17.80	...	s	y	w	u	w	t	g	
3	p	14.17	f	h	e	f	e	NaN	w	15.77	...	s	y	w	u	w	t	p	
4	p	14.64	x	h	o	f	e	NaN	w	16.53	...	s	y	w	u	w	t	p	

5 rows × 21 columns

Variable information (from UCI ML Repo page)

One binary class divided in edible=e and poisonous=p (with the latter one also containing mushrooms of unknown edibility). Twenty remaining variables (n: nominal, m: metrical)

1. cap-diameter (m): float number in cm
2. cap-shape (n): bell=b, conical=c, convex=x, flat=f, sunken=s, spherical=p, others=o
3. cap-surface (n): fibrous=i, grooves=g, scaly=y, smooth=s, shiny=h, leathery=l, silky=k, sticky=t, wrinkled=w, fleshy=e
4. cap-color (n): brown=n, buff=b, gray=g, green=r, pink=p, purple=u, red=e, white=w, yellow=y, blue=l, orange=o, black=k
5. does-bruise-bleed (n): bruises-or-bleeding=t,no=f
6. gill-attachment (n): adnate=a, adnexed=x, decurrent=d, free=e, sinuate=s, pores=p, none=f, unknown=?
7. gill-spacing (n): close=c, distant=d, none=f

8. `gill-color` (n): see `cap-color` + `none=f`
9. `stem-height` (m): float number in cm
10. `stem-width` (m): float number in mm
11. `stem-root` (n): `bulbous=b`, `swollen=s`, `club=c`, `cup=u`, `equal=e`, `rhizomorphs=z`, `rooted=r`
12. `stem-surface` (n): see `cap-surface` + `none=f`
13. `stem-color` (n): see `cap-color` + `none=f`
14. `veil-type` (n): `partial=p`, `universal=u`
15. `veil-color` (n): see `cap-color` + `none=f`
16. `has-ring` (n): `ring=t`, `none=f`
17. `ring-type` (n): `cobwebby=c`, `evanescent=e`, `flaring=r`, `grooved=g`, `large=l`, `pendant=p`, `sheathing=s`, `zone=z`, `scaly=y`, `movable=m`, `none=f`, `unknown=?`
18. `spore-print-color` (n): see `cap color`
19. `habitat` (n): `grasses=g`, `leaves=l`, `meadows=m`, `paths=p`, `heaths=h`, `urban=u`, `waste=w`, `woods=d`
20. `season` (n): `spring=s`, `summer=u`, `autumn=a`, `winter=w`

Adjusting column names

```
In [94]: # Replace hyphens with underscores
df.columns = df.columns.str.replace("-", "_")
```

```
In [95]: # View column names
df.columns
```

```
Out[95]: Index(['class', 'cap_diameter', 'cap_shape', 'cap_surface', 'cap_color',
       'does_bruise_or_bleed', 'gill_attachment', 'gill_spacing', 'gill_color',
       'stem_height', 'stem_width', 'stem_root', 'stem_surface', 'stem_color',
       'veil_type', 'veil_color', 'has_ring', 'ring_type', 'spore_print_color',
       'habitat', 'season'],
      dtype='object')
```

```
In [96]: # Rename columns with suboptimal names
df.rename(columns={"does_bruise_or_bleed": "bruises_or_bleeds", "class": "edible"}, inplace=True)
```

Moving label to the end of the dataframe

```
In [97]: # Move "edible" to the end
label = df.pop("edible")
```

```
df["edible"] = label
```

Showing dataset info

```
In [98]: # Display shape of data
print(f"Shape of dataset: {df.shape}\n")

# Display dataset info
print("Dataset info:")
print(df.info())
```

Shape of dataset: (61069, 21)

```
Dataset info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 61069 entries, 0 to 61068
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   cap_diameter    61069 non-null   float64
 1   cap_shape       61069 non-null   object 
 2   cap_surface     46949 non-null   object 
 3   cap_color       61069 non-null   object 
 4   bruises_or_bleeds 61069 non-null   object 
 5   gill_attachment 51185 non-null   object 
 6   gill_spacing    36006 non-null   object 
 7   gill_color      61069 non-null   object 
 8   stem_height     61069 non-null   float64
 9   stem_width      61069 non-null   float64
 10  stem_root       9531 non-null    object 
 11  stem_surface    22945 non-null   object 
 12  stem_color      61069 non-null   object 
 13  veil_type       3177 non-null    object 
 14  veil_color      7413 non-null    object 
 15  has_ring        61069 non-null   object 
 16  ring_type       58598 non-null   object 
 17  spore_print_color 6354 non-null   object 
 18  habitat          61069 non-null   object 
 19  season           61069 non-null   object 
 20  edible           61069 non-null   object 

dtypes: float64(3), object(18)
memory usage: 9.8+ MB
None
```

Our dataset contains mostly categorical/text columns, and so label encoding or one-hot encoding will be essential to prep the data for modeling.

Viewing class distribution for the label

As previously mentioned, the label in this case is `edible`, and it would be advantageous to check the class distribution to see if the SVM will be trained more on either edible or poisonous mushrooms, or if things are balanced. Training the model on a very unbalanced dataset could lead to a reduction in accuracy.

```
In [99]: # Check distribution of mushroom types  
df["edible"].value_counts()
```

```
Out[99]:  
edible  
p    33888  
e    27181  
Name: count, dtype: int64
```

There is a slight imbalance in favor of poisonous mushrooms, but it does not seem like anything severe enough to have a significant impact on the model's accuracy.

Viewing descriptive statistics for numeric columns

```
In [100...]  
# Show descriptive statistics  
df.describe().round(2)
```

Out[100]:

	cap_diameter	stem_height	stem_width
count	61069.00	61069.00	61069.00
mean	6.73	6.58	12.15
std	5.26	3.37	10.04
min	0.38	0.00	0.00
25%	3.48	4.64	5.21
50%	5.86	5.95	10.19
75%	8.54	7.74	16.57
max	62.34	33.92	103.91

Task 2.6

Clean the data and address unusual phenomena (e.g., normalization, feature scaling, outliers); use illustrative diagrams and plots and explain them.

Removing duplicates

In [101...]

```
# Find number of duplicate rows
print(f"Number of duplicate rows: {df.duplicated().sum()}")
df = df.drop_duplicates()
print(f"Number of duplicate rows after dropping: {df.duplicated().sum()}")
```

Number of duplicate rows: 146

Number of duplicate rows after dropping: 0

Handling missing values

In [102...]

```
# Find number of rows with missing values
print(f"Number of rows with missing values: {df.isna().any(axis=1).sum()}")
print(f"Total number of rows: {len(df)}\n")

# Find the number of missing values in each column
```

```
print("Missing variable count by column:")
print(df.isna().sum())
```

```
Number of rows with missing values: 60923
Total number of rows: 60923
```

```
Missing variable count by column:
```

```
cap_diameter          0
cap_shape              0
cap_surface            14120
cap_color                0
bruises_or_bleeds      0
gill_attachment        9855
gill_spacing            25062
gill_color                0
stem_height              0
stem_width                0
stem_root               51536
stem_surface             38122
stem_color                0
veil_type               57746
veil_color               53510
has_ring                  0
ring_type                 2471
spore_print_color       54597
habitat                   0
season                     0
edible                     0
dtype: int64
```

Every single row has at least one null value, so we will need to systematically go through the offending columns and replace these values. We will need to handle `cap_surface`, `gill_attachment`, `gill_spacing`, `stem_root`, `stem_surface`, `veil_type`, `veil_color`, `ring_type` and `spore_print_color`. My approach to these columns will be different depending on how many missing values there are.

- **Over 50% missing:** Remove the column. This is because imputing a large number of values can introduce significant noise to the data.
- **Up to 50% missing:** Impute a random selection from the top 3 most common values.

In [103...]

```
import random

# Drop columns with too many missing values
df.drop(columns={"stem_root", "stem_surface", "veil_type", "veil_color", "spore_print_color"}, inplace=True)
```

```
# Define function to impute a random value from the top 3 most common values in a column
def impute_missing(df, col, n=3, random_state=42):
    top_n_values = df[col].value_counts().index[:n]
    missing_indices = df[df[col].isnull()].index
    imputed_values = [random.choice(top_n_values) for _ in range(len(missing_indices))]
    df.loc[missing_indices, col] = imputed_values

# Impute missing values to remaining columns
for col in df.columns:
    if df[col].isnull().sum() > 0:
        impute_missing(df, col, n=3, random_state=42)

# Check number of rows with missing values after dropping and imputation
print(f"Number of rows with missing values: {df.isna().any(axis=1).sum()}")
```

Number of rows with missing values: 0

Thus, the problem with missing values has been resolved.

Converting data types

In [104...]

```
# Check datatypes and example values again
print(df.info())
print(df.head())
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 60923 entries, 0 to 61068
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   cap_diameter     60923 non-null   float64
 1   cap_shape        60923 non-null   object 
 2   cap_surface      60923 non-null   object 
 3   cap_color        60923 non-null   object 
 4   bruises_or_bleeds 60923 non-null   object 
 5   gill_attachment  60923 non-null   object 
 6   gill_spacing     60923 non-null   object 
 7   gill_color       60923 non-null   object 
 8   stem_height      60923 non-null   float64
 9   stem_width       60923 non-null   float64
 10  stem_color       60923 non-null   object 
 11  has_ring         60923 non-null   object 
 12  ring_type        60923 non-null   object 
 13  habitat          60923 non-null   object 
 14  season           60923 non-null   object 
 15  edible           60923 non-null   object 
dtypes: float64(3), object(13)
memory usage: 9.9+ MB
None
    cap_diameter  cap_shape  cap_surface  cap_color  bruises_or_bleeds  \
0            15.26        x            g          o                  f
1            16.60        x            g          o                  f
2            14.07        x            g          o                  f
3            14.17        f            h          e                  f
4            14.64        x            h          o                  f

    gill_attachment  gill_spacing  gill_color  stem_height  stem_width  stem_color  \
0                  e            f            w          16.95      17.09        w
1                  e            d            w          17.99      18.19        w
2                  e            c            w          17.80      17.74        w
3                  e            c            w          15.77      15.98        w
4                  e            f            w          16.53      17.20        w

    has_ring  ring_type  habitat  season  edible
0        t        g        d       w       p
1        t        g        d       u       p
2        t        g        d       w       p
3        t        p        d       w       p
4        t        p        d       w       p
```

First, we need to encode the label `edible` such that p \rightarrow 0 and e \rightarrow 1.

In [105...]

```
import numpy as np

# Encode Label
df["edible"] = np.where(df["edible"] == "p", 0, 1)
df["edible"].value_counts()
```

Out[105]:

```
edible
0    33742
1    27181
Name: count, dtype: int64
```

The numeric features `cap_diameter`, `stem_height` and `stem_width` all look fine, but every single categorical feature will need to be encoded to a numeric type. These categorical features all look nominal, so one-hot encoding seems to be the best approach.

In [106...]

```
# Select categorical columns
categorical_columns = df.select_dtypes(include=["object"]).columns

# Apply one-hot encoding
df = pd.get_dummies(df, columns=categorical_columns, drop_first=True)

# Move label "edible" to the end again
label = df.pop("edible")
df["edible"] = label

# Convert bool columns to int
for col in df.columns:
    if df[col].dtype == bool:
        df[col] = df[col].astype(int)

# View count of each datatype
df.dtypes.value_counts()
```

Out[106]:

```
int32      78
float64     3
Name: count, dtype: int64
```

We see that all columns are now appropriate datatypes for modeling.

Normalizing numeric features

`cap_diameter`, `stem_height` and `stem_width` are all of type `float64`, so we need to normalize these so they are on a [0, 1] scale alone with the other variables.

In [107]:

```
from sklearn.preprocessing import MinMaxScaler

# Initialize scaler
scaler = MinMaxScaler()

# Select columns to normalize
float_columns = ["cap_diameter", "stem_height", "stem_width"]

# Normalize columns
df[float_columns] = scaler.fit_transform(df[float_columns])

# Show final dataframe for modeling
df.head()
```

Out[107]:

	cap_diameter	stem_height	stem_width	cap_shape_c	cap_shape_f	cap_shape_o	cap_shape_p	cap_shape_s	cap_shape_x	cap_surface_e	...
0	0.240155	0.499705	0.164469	0	0	0	0	0	1	0	...
1	0.261782	0.530366	0.175055	0	0	0	0	0	1	0	...
2	0.220949	0.524764	0.170725	0	0	0	0	0	1	0	...
3	0.222563	0.464917	0.153787	0	1	0	0	0	0	0	...
4	0.230148	0.487323	0.165528	0	0	0	0	0	1	0	...

5 rows × 81 columns

The data is now ready for modeling.

Task 2.7

Formulate two questions that can be answered by applying a classification method using the SVM.

Question 1: Can we reliably predict whether a mushroom will be edible or poisonous based on its physical characteristics?

Question 2: What features of mushrooms have the strongest impact on a mushroom's edibility?

Task 2.8

Split the data into 80% training and 20% testing sets using the train test split class.

NOTE: I found later on that attempting to train the SVM with 80% of the data leads to unacceptably slow runtimes, so I have opted to invert the percentages, IE 20% training and 80% testing, as this leads to the model being trained in a reasonable amount of time. I acknowledge that the accuracy of the model may be affected by this tweak.

In [108...]

```
from sklearn.model_selection import train_test_split

# Select features and label
x = df.drop(columns={"edible"})
y = df["edible"]

# Perform train-test split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.8, random_state=42)

# Show shapes of variables
vars = [x_train, x_test, y_train, y_test]
names = ["x_train", "x_test", "y_train", "y_test"]

for i in range(0, len(vars)):
    print(f"Shape of {names[i]}: {vars[i].shape}")
```

```
Shape of x_train: (12184, 80)
Shape of x_test: (48739, 80)
Shape of y_train: (12184,)
Shape of y_test: (48739,)
```

Task 2.9

Use a linear kernel to train the SVM classifier on the training set (e.g., fit the support vector regressor to the dataset). Explain the intuition behind each of the key mathematical steps.

The SVM's goal is to find the optimal hyperplane that maximizes the margin between the two classes present in the label, and it does this by solving the following optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (1)$$

Subject to:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad \forall i \quad (2)$$

Where w is the weight vector, b is the bias term, y_i are the labels, and x_i are the feature vectors.

The width of the optimal margin is

$$\text{Margin} = \frac{2}{\|\mathbf{w}\|} \quad (3)$$

Where w is the optimal weight vector.

In [109...]

```
from sklearn.svm import SVC

# Initialize SVM classifier with a Linear kernel
classifier = SVC(kernel="linear", probability=True)

# Train classifier
classifier.fit(x_train, y_train)
```

Out[109]:

```
SVC(kernel='linear', probability=True)
```

Using a linear kernel during classifier initialization ensures that the decision boundary between classes will be linear. When the classifier is being fit to the training data, it is trying to find the maximum margin between classes, and it does this by solving an optimization problem that involves finding the optimal hyperplane that maximizes the margin.

Task 2.10

Explain the choice of the optimal hyperplane.

As previously mentioned, the optimal hyperplane is the one which maximizes the margin between classes. A larger margin is better because it improves the model's generalizability to new data, as well as insulating it against the negative effects of noise and outliers.

Task 2.11

Make classification predictions.

In [110]:

```
# Make predictions with the model
y_pred = classifier.predict(x_test)
```

In [111]:

```
# Assemble dataframe of predictions
prediction_vs_actual = pd.DataFrame({"prediction": y_pred, "actual": y_test})
prediction_vs_actual.head().transpose()
```

Out[111]:

	2980	12960	6920	15694	19082
prediction	1	1	0	1	0
actual	1	0	1	0	0

In [112]:

```
# Find number of correct predictions
correct_prediction_count = len(prediction_vs_actual[prediction_vs_actual["prediction"] == prediction_vs_actual["actual"]])

# Find percentage of correct predictions
accuracy = (correct_prediction_count / len(y_test)) * 100

print(f"Number of test entries: {len(y_test)}")
print(f"Number of correct predictions: {correct_prediction_count}")
print(f"Accuracy: {round(accuracy, 2)}%")
```

Number of test entries: 48739
Number of correct predictions: 38572
Accuracy: 79.14%

Task 2.12

Interpret the results in the context of the questions you asked.

Question 1: Can we reliably predict whether a mushroom will be edible or poisonous based on its physical characteristics?

Our accuracy score indicates that the SVM classifier we have built correctly predicts whether a mushroom is edible or poisonous about 79.2% of the time, which confirms that it is indeed possible to reliably predict a mushroom's edibility based on its physical characteristics.

Question 2: What features of mushrooms have the strongest impact on a mushroom's edibility?

The second question will require a different type of classification model to answer. We will quickly apply a logistic regression model and use the feature coefficients to determine which features are the most important in determining a mushroom's edibility.

In [113...]

```
from sklearn.linear_model import LogisticRegression

# Fit Logistic regression model
model = LogisticRegression()
model.fit(x_train, y_train)

# Get model coefficients
feature_coefficients = model.coef_[0]

# Get feature names
feature_names = df.drop(columns={"edible"}).columns

# Assemble dataframe of features and their coefficients
feature_importance_df = pd.DataFrame({
    "feature_name": feature_names,
    "coefficient": feature_coefficients,
    "abs_coefficient": np.abs(feature_coefficients)
})

# Show top 5 most impactful features and their coefficients
feature_importance_df.sort_values(by="abs_coefficient", ascending=False).head()
```

Out[113]:

	feature_name	coefficient	abs_coefficient
66	ring_type_m	4.688757	4.688757
51	stem_color_f	-3.889348	3.889348
76	habitat_w	3.071529	3.071529
0	cap_diameter	2.922636	2.922636
26	cap_color_r	-2.882229	2.882229

A positive coefficient means that the presence of that feature makes it more likely for the mushroom to be edible (`edible = 1`), and a negative coefficient makes it less likely to be edible.

We can see from the logistic regression results that these 5 features are the most important in determining a mushroom's edibility:

1. Whether or not the mushroom has a movable ring type. Mushrooms with movable ring types are more likely to be edible.
2. Whether or not the mushroom has a stem color. Mushrooms with no stem color are less likely to be edible.
3. Whether or not the mushroom lives in waste. Mushrooms that live in waste are more likely to be edible (surprisingly).
4. The mushroom's cap diameter. Mushrooms with larger caps are more likely to be edible.
5. Whether or not the mushroom has a red cap. Mushrooms with red caps are less likely to be edible.

Thus, we have successfully answered question 2 as well.

Task 2.13

Validate your model using a confusion matrix, accuracy score, ROC-AUC curves, and k-fold cross validation. Then, explain the results.

Confusion matrix

In [114...]

```
from sklearn.metrics import confusion_matrix

# Create confusion matrix
cm = confusion_matrix(y_test, y_pred)

# Find order of Label values
print("Confusion matrix:")
```

```
# Print confusion matrix
print(cm)
```

```
Confusion matrix:
[[21871  5066]
 [ 5101 16701]]
```

The labeled matrix is as follows:

	predicted negative	predicted positive
actual negative	21,869	5,068
actual positive	5,049	16,753

It is clear that most of our model's predictions are either true negatives or true positives, so the confusion matrix validates what we already determined from our accuracy score.

Accuracy score

Our SVM classifier's accuracy score is about 79.2% as previously mentioned, meaning it will correctly predict whether a mushroom is edible or poisonous about 4 out of 5 times.

ROC-AUC curve

In [115...]

```
from sklearn.metrics import roc_curve, roc_auc_score
import matplotlib.pyplot as plt

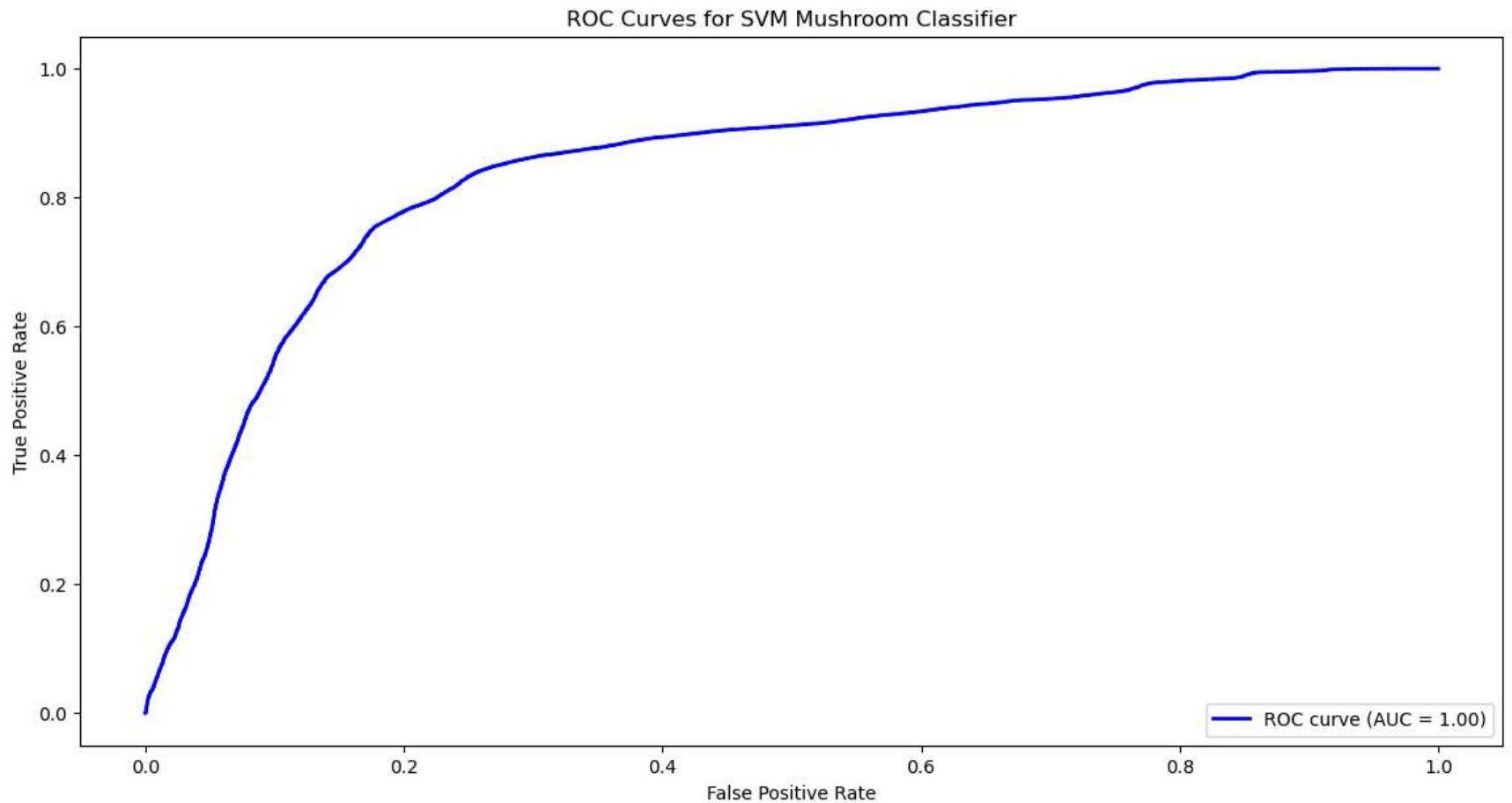
# Calculate ROC-AUC metrics
y_prob = classifier.predict_proba(x_test)[:, 1]
fpr, tpr, thresholds = roc_curve(y_test, y_prob)
auc = roc_auc_score(y_test, y_test)

# Set up plot figure
plt.figure(figsize=(14, 7))

# Plot ROC-AUC curves
plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {auc:.2f})')

# Configure and show plot
```

```
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curves for SVM Mushroom Classifier")
plt.legend(loc="lower right")
plt.show()
```



The ROC-AUC curve rises to high vertical axis values quickly at low values of the horizontal axis, meaning that the true positive rate is very high compared to the false positive rate. This demonstrates that the model is effective at correctly classifying mushrooms.

K-fold cross-validation

Once again, I found this step to be taking too long to execute, so I chose to perform k-fold cross-validation on smaller subsets of the data.

In [118...]

```
from sklearn.model_selection import cross_val_score

# Take samples to improve performance
x_sample = x.sample(10000, random_state=42)
y_sample = y.sample(10000, random_state=42)

# Calculate k-fold cross-validation scores
scores = cross_val_score(classifier, x_sample, y_sample, cv=5, scoring="accuracy")
```

In [122...]

```
# Show cross-validation scores
print("Cross-validation scores:", scores)
print(f"Mean score: {round(scores.mean(), 3)}")
print(f"Accuracy score: {round(accuracy, 2)}%")
```

Cross-validation scores: [0.7955 0.7915 0.783 0.7875 0.7925]

Mean score: 0.79

Accuracy score: 79.14%

The mean cross-validation score hovers very close to the accuracy score from earlier at about 79%, showing that the SVM classifier is highly generalizable and effective even with unseen data.

Task 2.14

Include all mathematical formulas used and graphs representing the final outcomes.

All relevant formulas and graphs have been included in their respective tasks.

References

Wagner,Dennis, Heider,D., and Hattab,Georges. (2023). Secondary Mushroom. UCI Machine Learning Repository.
<https://doi.org/10.24432/C5FP5Q>.

Nik. (2023, April 22). Support Vector Machines (SVM) in Python with Sklearn. Datagy. <https://datagy.io/python-support-vector-machines/>

GeeksforGeeks. (2023, February 2). Introduction to Support Vector Machines (SVM). GeeksforGeeks.
<https://www.geeksforgeeks.org/introduction-to-support-vector-machines-svm/>

ChatGPT. (n.d.). <https://chat.openai.com/>