

## IS 428: Lab 6

### User

Task 1b: Build a logistic regression model

Using the dataset "student\_data.csv," write code to (1) create a visualization of the data, (2) fit a model using logistic regression, (3) output model coefficients and performance metrics such as accuracy and AUC and ROC;

### ChatGPT

Alright! Let's follow the steps:

1. Create a visualization of the data.
2. Fit a logistic regression model.
3. Output model coefficients and performance metrics such as accuracy, AUC, and ROC.

Let's start by loading the dataset and examining its structure.

### ChatGPT

It seems there was an issue processing the dataset. Let me try again to load the data and inspect its contents.

### ChatGPT

It seems that due to an unexpected issue, the environment has been reset. I apologize for the inconvenience. I'll need to reload the necessary libraries and then load the dataset once more. Let's proceed with that.

### ChatGPT

The dataset has three columns:

1. `Hours\_Studied`: The number of hours a student studied.
2. `Review\_Session`: A binary variable indicating whether a student attended the review session (1 for yes, 0 for no).

3. ``Results``: A binary outcome variable indicating whether a student passed (1 for pass, 0 for fail).

Let's start with the first task and create visualizations to understand the data distribution and relationships. We'll visualize:

1. The distribution of hours studied.
2. The relationship between hours studied and the results.
3. The impact of attending a review session on the results.

#### **ChatGPT**

The visualizations provide the following insights:

1. **\*\*Distribution of Hours Studied\*\***: The distribution shows that most students studied between 2 to 10 hours, with few students studying less than 2 hours.

2. **\*\*Relationship between Hours Studied and Results\*\***:

- Generally, as the number of hours studied increases, the likelihood of passing (a result of 1) also seems to increase.

- Students who attended the review session (indicated by the purple 'x' markers) appear to have a higher chance of passing, especially if they studied fewer hours.

Now, let's proceed with the second task: fitting a logistic regression model.

#### **ChatGPT**

The logistic regression model's performance metrics are as follows:

- **\*\*Accuracy\*\***: Approximately  $(86.67\%)$
- **\*\*AUC (Area Under the Curve)\*\***: 0.94

An AUC of 0.94 indicates a very good model performance.

Next, we'll display the model coefficients and visualize the ROC curve.

### ChatGPT

The model coefficients and intercept are as follows:

- Coefficient for `Hours\_Studied`:  $\approx 1.575$
- Coefficient for `Review\_Session`:  $\approx 1.294$
- Intercept:  $\approx -5.097$

The positive coefficients for both `Hours\_Studied` and `Review\_Session` indicate that:

1. An increase in the number of hours studied increases the log odds of passing the exam.
2. Attending the review session also increases the log odds of passing compared to not attending.

The ROC curve, which is plotted above, demonstrates a strong performance of the model, as it's significantly above the diagonal line (indicating random chance).

Would you like to explore or analyze anything further?

### User

NameError Traceback (most recent call last)

Cell In[4], line 8

```
5 import seaborn as sns

7 # Load the dataset

----> 8 data = pd.read_csv('student_data.csv')
```

```
10 # Set up the figure and axes

11 fig, ax = plt.subplots(nrows=2, ncols=1, figsize=(12, 10))
```

NameError: name 'pd' is not defined

#### ChatGPT

It appears you're trying to run a code snippet, but there's an error due to the missing import statement for the pandas library (`pd`).

If you'd like to run some code or analysis, please let me know your specific requirements or objectives, and I can assist you further! If you're pointing out an error from an earlier execution, I apologize for the oversight and will correct it accordingly.

#### User

```
-----
NameError                                Traceback (most recent call last)

Cell In[5], line 15

    12 fig, ax = plt.subplots(nrows=2, ncols=1, figsize=(12, 10))

    14 # Distribution of hours studied
--> 15 sns.histplot(student_data['Hours_Studied'], kde=True, ax=ax[0])

    16 ax[0].set_title('Distribution of Hours Studied')

    17 ax[0].set_xlabel('Hours Studied')
```

NameError: name 'student\_data' is not defined

#### ChatGPT

This error message indicates that the variable `student\_data` is not defined. It seems you're pointing out a previous error in the code execution.

I'm here to help! If you have any specific requirements or need clarifications on any part of the analysis, please let me know. Alternatively, if you'd like to explore another aspect of the dataset or any other topic, please guide me accordingly.

### **User**

When I run this code you gave me with:

```
import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

# Load the dataset

data = pd.read_csv('student_data.csv')

# Set up the figure and axes

fig, ax = plt.subplots(nrows=2, ncols=1, figsize=(12, 10))

# Distribution of hours studied

sns.histplot(student_data['Hours_Studied'], kde=True, ax=ax[0])

ax[0].set_title('Distribution of Hours Studied')

ax[0].set_xlabel('Hours Studied')

ax[0].set_ylabel('Number of Students')

# Relationship between hours studied and results, with distinction for review
session

sns.scatterplot(data=student_data, x='Hours_Studied', y='Results',
hue='Review_Session', ax=ax[1], alpha=0.6)

ax[1].set_title('Relationship between Hours Studied and Results')
```

```
ax[1].set_xlabel('Hours Studied')

ax[1].set_ylabel('Results (Pass/Fail)')


plt.tight_layout()

plt.show()
```

I get the following errors:

```
-----

NameError                                Traceback (most recent call last)

Cell In[5], line 15

    12 fig, ax = plt.subplots(nrows=2, ncols=1, figsize=(12, 10))

    14 # Distribution of hours studied
--> 15 sns.histplot(student_data['Hours_Studied'], kde=True, ax=ax[0])

    16 ax[0].set_title('Distribution of Hours Studied')

    17 ax[0].set_xlabel('Hours Studied')
```

NameError: name 'student\_data' is not defined

### ChatGPT

I see the issue. In the code you provided, you're loading the dataset into a variable named `data`:

```
```python

data = pd.read_csv('student_data.csv')

```
```

However, in the subsequent lines, you're trying to access the dataset using the variable `student\_data`:

```
```python

sns.histplot(student_data['Hours_Studied'], kde=True, ax=ax[0])

```
```

To fix this, you should either rename the variable `data` to `student\_data` when you load the dataset or modify the subsequent lines to use the variable `data`.

Here's the corrected code:

```
```python

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

# Load the dataset

student_data = pd.read_csv('student_data.csv') # <-- Changed 'data' to
'student_data' here

# Set up the figure and axes

fig, ax = plt.subplots(nrows=2, ncols=1, figsize=(12, 10))

# Distribution of hours studied
```

```

sns.histplot(student_data['Hours_Studied'], kde=True, ax=ax[0])

ax[0].set_title('Distribution of Hours Studied')

ax[0].set_xlabel('Hours Studied')

ax[0].set_ylabel('Number of Students')


# Relationship between hours studied and results, with distinction for review
session

sns.scatterplot(data=student_data, x='Hours_Studied', y='Results',
hue='Review_Session', ax=ax[1], alpha=0.6)

ax[1].set_title('Relationship between Hours Studied and Results')

ax[1].set_xlabel('Hours Studied')

ax[1].set_ylabel('Results (Pass/Fail)')


plt.tight_layout()

plt.show()

...

```

This should resolve the `NameError` you encountered.

### **User**

Task 2: Understanding and Prevent Overfitting in the context of SVM

Write code to fit a Support Vector Machine model using (1) linear kernel and (2) RBF kernel. For the RBF kernel, use grid search to find the best gamma parameter using k-fold cross-validation.

### **ChatGPT**

Sure, let's proceed step by step:



1. Split the data into training and test sets.
2. Fit a Support Vector Machine (SVM) model using a linear kernel.
3. Fit an SVM model using an RBF kernel and utilize grid search to find the best gamma parameter with k-fold cross-validation.

Let's start by splitting the data into training and test sets.

### **ChatGPT**

Now that we have our training and test sets, let's proceed with fitting the SVM models.

1. We'll start by fitting an SVM model using a linear kernel.
2. Then, we'll fit an SVM model using the RBF kernel and utilize grid search with k-fold cross-validation to find the best gamma parameter.

### **ChatGPT**

The results for the Support Vector Machine (SVM) models are as follows:

1. **SVM with Linear Kernel**:

- Accuracy: \ (86.67\%\)

2. **SVM with RBF Kernel**:

- Best Gamma: `scale` (This means the gamma is computed as the inverse of the variance of the training data.)

- Accuracy with Best Gamma: \ (86.67\%\)

Both the linear kernel and the RBF kernel with the best gamma parameter provide the same accuracy on the test set.

Would you like to explore or analyze anything further related to SVM or any other topic?