

Introduction to Machine Learning – Homework Assignment 1

Homework Title: Implementation of a Simple Naive Bayes Classifier in Python

Due Date: [In Google Classroom]

Objective

The purpose of this assignment is to deepen your understanding of the Naive Bayes algorithm by implementing it from scratch using Python. This exercise will help you:

- Grasp the foundational concepts behind probabilistic classifiers.
 - Develop skills in implementing machine learning algorithms without relying on pre-built libraries.
 - Apply your implementation to the "Play Tennis" toy dataset.
 - Analyze and interpret the results of your classifier.
-

Assignment Description

You are required to implement a **Naive Bayes classifier** from scratch. You **must not** use any machine learning libraries such as scikit-learn for the Naive Bayes implementation. However, you may use libraries for data manipulation and basic operations (e.g., NumPy, pandas).

Your implementation should:

- Handle both **categorical and numerical features** (though the "Play Tennis" dataset primarily contains categorical data).
- In the training stage, given the training set in json format, your code must calculate the prior probabilities and likelihoods (class conditional attribute value probabilities) based on the training data and persist them as the model in your HDD. You can also use json file to persist your model. Given a new training set, your code must update the model.
- In the classification stage, you will be given a single instance without class label in json format. Your code must read your model (e.g. json files) from hdd, load it to the memory and classify a single instance using your model and output the predicted class label along with all the calculations printed to the screen.

- If in the classification stage a number of instances are given along with their respective class labels (test set), your code must load the model from the HDD to the memory, classify each instance in the test set one by one and compare the predicted class label with the actual class label to calculate accuracy. All the classification operations, predicted class labels vs. actual class labels must be printed to a log file.
-

Dataset

You will use the "**Play Tennis**" dataset, a commonly used toy dataset in machine learning examples.

Dataset Description:

The dataset contains weather conditions and a target variable indicating whether tennis was played. It includes the following features:

- **Outlook:** Sunny, Overcast, Rain
- **Temperature:** Hot, Mild, Cool
- **Humidity:** High, Normal
- **Wind:** Weak, Strong
- **PlayTennis (Target Variable):** Yes, No

Dataset Instances:

	Day	Outlook	Temperature	Humidity	Wind	PlayTennis
--	-----	---------	-------------	----------	------	------------

1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Assignment Tasks

1. Data Preparation

- **Load the Dataset:**
 - Create the dataset in your program by manually inputting the data or reading from a CSV file you create.
- **Explore the Data:**
 - Display the dataset in a tabular format.
 - Summarize the dataset (e.g., count of instances per class).
- **Encode Features:**
 - Convert categorical features into a suitable numerical format if necessary for your implementation.
 - You can use techniques like label encoding or one-hot encoding, but remember that for Naive Bayes with categorical data, you can work with the categories directly.

2. Implement the Naive Bayes Classifier

- **Calculate Prior Probabilities:**
 - Compute the prior probabilities for each class (PlayTennis = Yes, PlayTennis = No).
- **Calculate Likelihoods:**
 - For each feature value given a class, calculate the likelihoods (conditional probabilities).
- **Implement the Algorithm:**
 - Create a function to train the Naive Bayes classifier using the training data.
 - Implement the prediction function that applies the Naive Bayes formula to classify new instances.
- **Handle Zero Probabilities:**
 - Implement **Laplace smoothing** to handle cases where a probability may be zero due to unseen feature-class combinations.

3. Testing the Classifier

- **Test Instances:**
 - Use the same dataset for testing due to its small size.
 - Alternatively, you can perform **Leave-One-Out Cross-Validation**:
 - For each instance in the dataset:
 - Use all other instances for training.
 - Test on the held-out instance.
- **Make Predictions:**
 - Predict the class label for each test instance.
 - Compare the predicted labels with the actual labels.

4. Evaluate the Model

- **Calculate Performance Metrics:**
 - **Accuracy:** The proportion of correctly classified instances.
 - **Confusion Matrix:** Show true positives, false positives, true negatives, and false negatives.
- **Interpret Results:**
 - Discuss the performance of your classifier.
 - Identify any patterns or observations in misclassified instances.

5. Write a Report

Prepare a report that includes:

- **Introduction:**
 - Briefly explain the Naive Bayes algorithm.
 - State the objectives of the assignment.
- **Methodology:**
 - Describe how you prepared the data.
 - Explain the implementation details of your Naive Bayes classifier.
 - Discuss how you handled any challenges (e.g., zero probabilities).
- **Results:**
 - Present the performance metrics.
 - Include any tables or figures that help illustrate your results.
- **Discussion:**
 - Analyze the results.
 - Suggest possible reasons for any misclassifications.
 - Discuss the limitations of your implementation.
- **Conclusion:**
 - Summarize your findings and what you learned.
- **References:**
 - Cite any resources you used.

Submission Guidelines

What to Submit

1. **Source Code:**
 - All Python source code files.
 - Json files for the datasets and for the model
 - Include a README file with instructions on how to run your code.
 - Ensure your code is well-commented to explain key sections.
2. **Report:**

- A PDF file named `Surname_Name_NaiveBayes_Report.pdf`.
- The report should be 3-5 pages, excluding references and appendices.

How to Submit

- **File Naming Convention:**
 - Compress all your files into a single ZIP file named `Surname_Name_NaiveBayes_HW.zip`.
- **Submission Platform:**
 - Submit the ZIP file via the university's learning management system (e.g., Moodle, Blackboard).

Code Requirements

- **Programming Language:**
 - Use Python 3.x for your implementation.
 - **Allowed Libraries:**
 - **Permitted:**
 - NumPy
 - pandas
 - matplotlib or seaborn (for plotting, if necessary)
 - **Not Permitted:**
 - scikit-learn (for model implementation)
 - Any other machine learning libraries that provide Naive Bayes functionality.
-

Academic Integrity

- **Original Work:**
 - The code and report must be your own work.
 - **Collaboration Policy:**
 - Discussions with classmates are allowed for understanding concepts.
 - Sharing of code or written content is prohibited.
 - **Plagiarism:**
 - Any form of plagiarism will result in a zero for the assignment and may lead to further disciplinary action.
-

Resources

- **Lecture Materials:**

- Review the slides and notes on Naive Bayes classifiers provided in class.
 - **Textbooks:**
 - *Introduction to Machine Learning* by Ethem Alpaydin
 - *Pattern Recognition and Machine Learning* by Christopher Bishop
 - **Online Tutorials:**
 - Articles and tutorials on Naive Bayes algorithms.
 - **Documentation:**
 - Python documentation for libraries like NumPy and pandas.
-
-

Questions and Clarifications

- Please ask all your questions from the google classroom only. Add a comment to this assignment
-

Good luck with your assignment! Implementing algorithms from scratch is a valuable exercise that enhances your understanding and prepares you for more advanced topics in machine learning.