

### Assignment

Q1. A company conducted a survey of its employees and found that 70% of the employees use the company's health insurance plan, while 40% of the employees who use the plan are smokers. What is the

probability that an employee is a smoker given that he/she uses the health insurance plan?

Ans: To find the probability that an employee is a smoker given that he/she uses the health insurance plan, you can use conditional probability. The notation for this is often written as  $P(\text{Smoker} | \text{Uses health insurance})$ .

The formula for conditional probability is:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) =$$

$$P(B)$$

$$P(A \cap B)$$

In this case:

- Let A be the event "Employee is a smoker."
- Let B be the event "Employee uses the health insurance plan."

You are given:

$$P(B) = 0.70$$

$$P(B) = 0.70 \text{ (probability that an employee uses the health insurance plan)}$$

$$P(A \cap B) = 0.40$$

$$P(A \cap B) = 0.40 \text{ (probability that an employee is a smoker and uses the health insurance plan)}$$

Now you can plug these values into the formula:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = 0.40 / 0.70$$

$$P(A|B) =$$

$$P(B)$$

$$P(A \cap B)$$

=

$$0.70$$

$$0.40$$

Calculating this gives:

$$P(\text{Smoker} | \text{Uses health insurance}) = 0.40 / 0.70 \approx 0.5714$$

$$P(\text{Smoker} | \text{Uses health insurance}) =$$

$$0.70$$

$$0.40$$

$$\approx 0.5714$$

So, the probability that an employee is a smoker given that he/she uses the health insurance plan is approximately 0.5714, or 57.14%.

Q2. What is the difference between Bernoulli Naive Bayes and Multinomial Naive Bayes?

Ans: Both Bernoulli Naive Bayes and Multinomial Naive Bayes are algorithms used for text classification and other machine learning tasks. The primary difference between them lies in the type of data they are designed to handle:

Bernoulli Naive Bayes:

- Type of data: It is suitable for binary data, where features represent either presence (1) or absence (0) of a particular characteristic.
- Example application: Spam detection, where the presence or absence of certain keywords in an email is considered.

Multinomial Naive Bayes:

- Type of data: It is designed for discrete data, particularly for situations where features represent the frequency of occurrences within a fixed interval (e.g., word counts in a document).
- Example application: Text categorization, document classification, sentiment analysis, etc., where the frequency of words matters.

In summary, the key distinction lies in the nature of the features. If your features are binary (presence or absence), you might lean towards using Bernoulli Naive Bayes. If your features involve counts or frequencies, Multinomial Naive Bayes is a more appropriate choice. Both algorithms assume that the features are conditionally independent given the class label, which is the "Naive" assumption in Naive Bayes.

Q3. How does Bernoulli Naive Bayes handle missing values?

Ans: In the context of Bernoulli Naive Bayes, handling missing values can be important, especially when dealing with binary features where values are either present (1) or absent (0). The way missing values are typically handled in Bernoulli Naive Bayes depends on the specific implementation and the requirements of the task. Here are a few common approaches:

Ignoring Missing Values:

- Some implementations of Bernoulli Naive Bayes may simply ignore instances with missing values during training and classification. This means that any instance with a missing value for a particular feature is not used when estimating probabilities or making predictions involving that feature.

Imputation:

- Another approach is to impute missing values. In the context of Bernoulli features (binary), this could involve assigning the most common value (0 or 1) for the missing value in that feature.

Consider Missing as a Separate Category:

- In some cases, missing values are treated as a separate category. Instead of imputing a specific value, a missing value is considered as a distinct state, and the Naive Bayes model is trained accordingly.

The choice of how to handle missing values may depend on the nature of the data, the specific requirements of the problem, and the characteristics of the missing data. It's important to note

that the decision should be made carefully, considering the potential impact on the model's performance and the assumptions made during training and classification.

Q4. Can Gaussian Naive Bayes be used for multi-class classification?

Ans: Gaussian Naive Bayes can be used for multi-class classification. The Gaussian Naive Bayes algorithm is an extension of the Naive Bayes algorithm that assumes that the features follow a Gaussian (normal) distribution. This algorithm is commonly used when dealing with continuous data.

In the context of multi-class classification, Gaussian Naive Bayes can be adapted to handle more than two classes. The general idea remains the same as in binary classification: the algorithm calculates the conditional probability of each class given the input features and makes a prediction based on the class with the highest probability.

When dealing with multiple classes, the algorithm calculates the conditional probabilities for each class, and the class with the highest probability is assigned as the predicted class. The mathematical details involve estimating the mean and variance of each feature for each class based on the training data.

In summary, Gaussian Naive Bayes is a versatile algorithm that can be applied to both binary and multi-class classification problems, making it suitable for a variety of scenarios where the assumptions of the algorithm (independence of features) hold reasonably well.

Q5. Assignment:

Data preparation:

Download the "Spambase Data Set" from the UCI Machine Learning Repository (<https://archive.ics.uci.edu/ml/datasets/Spambase>). This dataset contains email messages, where the goal is to predict whether a message is spam or not based on several input features.

Implementation:

Implement Bernoulli Naive Bayes, Multinomial Naive Bayes, and Gaussian Naive Bayes classifiers using the scikit-learn library in Python. Use 10-fold cross-validation to evaluate the performance of each classifier on the dataset. You should use the default hyperparameters for each classifier.

Results:

Report the following performance metrics for each classifier:

Accuracy

Precision

Recall

F1 score

Discussion:

Discuss the results you obtained. Which variant of Naive Bayes performed the best? Why do you think that is

the case? Are there any limitations of Naive Bayes that you observed?

Conclusion:

Summarise your findings and provide some suggestions for future work.

Ans:the scikit-learn library in Python. Before you proceed, make sure you have the required libraries installed. You can install scikit-learn using:

bash

Copy code

Now, you can use the following code as a starting point for your implementation:

python

Copy code

```
import sys as sys
from sklearn import metrics
from sklearn import model_selection
from sklearn import datasets
```

```
"https://archive.ics.uci.edu/ml/machine-learning-databases/spambase/spambase.data"
    ,'
```

```
1
1
```

```
def evaluate_classifier
```

```
10
```

```
'accuracy'
```

```

        'precision'
        10
        10
        10
        'recall'
        'f1'

print f"\n{name} Naive Bayes:"
print f"Accuracy: {accuracy:.4f}"
print f"Precision: {precision:.4f}"
print f"Recall: {recall:.4f}"
print f"F1 Score: {f1_score:.4f}"

"Bernoulli"
"Multinomial"
"Gaussian"

```

This code uses the Spambase dataset from the UCI Machine Learning Repository, loads it, and then applies three different Naive Bayes classifiers (Bernoulli, Multinomial, and Gaussian) using 10-fold cross-validation. The performance metrics are then printed for each classifier.

Now, you can run this code and observe the results for each variant of Naive Bayes. Discussing the results involves analyzing which variant performed the best based on the provided performance metrics. You may observe that the choice of the best variant depends on the specific characteristics of the dataset.

Some points for discussion could include:

- How well each classifier performed in terms of accuracy, precision, recall, and F1 score.
- Whether the assumptions of each Naive Bayes variant (Bernoulli, Multinomial, Gaussian) align with the nature of the Spambase dataset.
- Limitations of Naive Bayes algorithms, such as the assumption of feature independence.

In the conclusion, summarize your findings and suggest potential areas for future work, such as exploring different feature engineering techniques or hyperparameter tuning to improve the performance of the classifiers.