

Naive Bayes with Hierarchical MCA

Lakshmi Katrapati & Ipek Sayar
Dr.Yilmaz
Machine Learning Pd. 6

1. Motivation

In this project, we aim to enhance the performance of the Naive Bayes algorithm by addressing the challenge of feature dependencies within a dataset. Traditional Naive Bayes assumes that attributes are independent, which can lead to inaccuracies when attributes are correlated. To address this, we apply Multiple Correspondence Analysis (MCA) to reduce these correlations and make the features more independent.

Related Works

An article by Luo and Liu discusses how the Naive Bayes algorithm can be improved by combining it with an advanced version of Principal Component Analysis (PCA)—a dimensionality reduction method that reduces complex data into fewer parts while keeping the most important information [2]. The authors introduce an improved PCA method that calculates correlation coefficients for both quantitative and qualitative data. By transforming the data into principal components, the relationship between correlated variables is reduced, helping to meet Naive Bayes' independence assumption and improving classification performance. The result is a more accurate and robust algorithm, particularly when dealing with noisy data.

This idea of improving Naive Bayes with PCA is the basis of our approach as well. Instead of applying traditional PCA, we focus on reducing dependencies between attributes with Multiple Correspondence Analysis (MCA), which works better with categorical data. By handling feature dependence and ensuring the data is more independent, we aim to improve the accuracy and recall of Naive Bayes, similar to how PCA enhances its performance in this study.

2. Intended Experiments

Correlation Analysis: we will calculate correlation coefficients for numerical attributes using Pearson and then use techniques like Chi-square test for categorical features to

quantify relationships. This will help determine the extent of correlation that might violate Naive Bayes' independence assumption.

Classification Model Selection: based on the correlation analysis and the results of MCA, we will choose and run various classification models. We will experiment with multiple models to evaluate how well each performs. These could include:

1. Complement Naive Bayes
2. Naive Bayes with Full MCA
3. Naive Bayes with Hierarchical MCA

Evaluation Metrics: Accuracy and Recall

We will evaluate the models by measuring the overall accuracy to determine how well they predicted the class label.

Since the dataset may be imbalanced, we will also evaluate the recall scores of the class. These metrics will help assess how well the model performs, especially in the context of potential class imbalances. Recall is the best metric to use here because we want to minimize false negatives. Heart attacks are serious, and while misclassifying someone who isn't having one may result in an unnecessary hospital visit, misclassifying someone who is having a heart attack could cost them their life in certain cases.

3. Dataset

Heart Attack in Youth VS Adult in France [1]

Heart attacks are life-threatening emergencies that can cause permanent damage to the heart and even cause death. Analyzing factors that cause heart attacks can lead to better insight on whether or not an individual faces a high risk of getting a heart attack.

This dataset gives different attributes of individuals in France and whether they had a heart attack or not.

The class label of this dataset is "Heart Attack" denoted by the values "Yes" or "No". Out of all the instances, around 20% had a value of "Yes", while the other 80% had a value of "No". if the person had a heart attack, then the value would be "Yes", else "No".

It contains 26 features that describe characteristics of the individual. individuals are denoted as “patients” (one of the features is patient_D) and examples of features include “age”, “weight_kg”, “cholesterol_level”, and “Age_Group”.

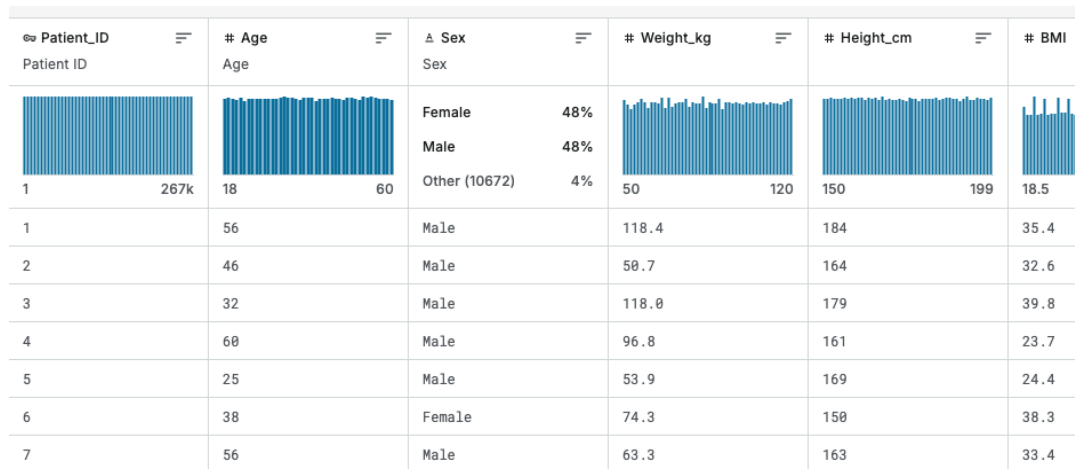


Figure 1. *A few of the features included in the dataset*

This dataset contains 266,785 instances where each instance represents a patient.

4. Preprocessing

Since our algorithm implements Multiple Correspondence Analysis (MCA), the values in our dataset need to be discretized. Some features in our dataset contain continuous, numerical values, and these values were sorted into different bins based on predefined classifications.

The following attributes were discretized: “Cholesterol_Level”, “Air_Pollution_Level”, “BMI”, “Alcohol_Consumption”, “Physical_Activity_Hours”, “Stress_Level”, and “Heart_Rate”. For each range, the lower number is inclusive while the upper number is exclusive. For example, if the range of values is ‘90-190’, this means that the category contains all numbers from 90 up until but not including 190. The sources where the values of categories are taken from are cited in the bibliography.

0-200	Healthy
200-240	At_Risk
≥ 240	'Dangerous'

Table 1: *Cholesterol level discretization values*

0-50	'Good'
51-100	'Moderate'
101-150	'Unhealthy_Sensitive'
151-200	'Unhealthy'
201-300	'Very_Unhealthy'
≥ 300	'Hazardous'

Table 2: *Air pollution level discretization values*

0-18.5	Underweight
18.5-25	Healthy_weight
25-30	Overweight
≥ 30	Obesity

Table 3: *BMI discretization values*

For the 'Alcohol_Consumption' attribute we assumed that small numbers meant low levels of consumption, as there are no predefined categories available for this attribute.

0	No Alcohol Consumption
1-4	Low Consumption

5-9	Moderate Consumption
10-15	High Consumption
>15	Very High Consumption

Table 4: *Alcohol consumption discretization values*

0	No Activity
0 - 3.7	Low Activity
3.7 - 7.5	Moderate Activity
7.5 - 11.3	Active
>11.3	Very Active

Table 5: *Physical activity hours discretization values*

1-3	Low Stress
4-5	Moderate Stress
6-8	High Stress
9-10	Very High Stress

Table 6: *Stress level discretization values*

50-59	Low Heart Rate (Bradycardia)
60-100	Normal Heart Rate
101-110	Elevated Heart Rate
111-119	High High Rate (Tachycardia)

Table 7: *Heart rate discretization values*

Several other numerical features also required discretization. However, this process involved combining multiple features to create a new attribute while discarding the original ones.

There is no need for two separate blood pressure attributes, therefore the attributes “Blood_Pressure_Systolic” and “Blood_Pressure_Diastolic” were combined to create a new attribute called “BP” (BP stands for Blood Pressure).

$$\begin{aligned}\text{Systolic} &= \text{Blood_Pressure_Systolic} \\ \text{Diastolic} &= \text{Blood_Pressure_Diastolic}\end{aligned}$$

Systolic < 120 & Diastolic < 80	Normal
120 ≤ Systolic ≤ 129 & Diastolic < 80	Elevated
130 ≤ Systolic ≤ 139 or 80 ≤ Diastolic ≤ 89	High_BP_Stage1
Systolic ≥ 140 or Diastolic ≥ 90	High_BP_Stage2
Systolic > 180 or Diastolic > 120	Hypertensive_Crisis

Table 8: *Blood Pressure (BP) discretization values*

After the ‘BP’ attribute was created, the ‘Blood_Pressure_Systolic’ and ‘Blood_Pressure_Diastolic’ attributes were deleted, because they no longer serve a purpose.

In order to discretize the attribute “Height” we considered the gender of the patient, because what is considered “short” for a female is not considered “short” for a male.

Female & height < 160 cm	Short
Female & 160cm ≤ height ≤ 170cm	Average
Female & 170cm < height ≤ 180 cm	Tall
Female & height > 180 cm	Very Tall
Male & height < 165 cm	Short

Male & 165 cm \leq height < 176 cm	Average
Male & 176cm \leq height < 190 cm	Tall
Male & height \geq 190 cm	Very Tall

Table 9: *Height attribute discretization values*

Attributes deleted: “Patient_ID”, “Age”, and “Weight_kg” .

- ‘Patient_ID’ was deleted because it gave no information relevant to the algorithm. It just contains the name-ID of the patient.
- ‘Age’ was deleted because there was another feature in the dataset called “Age_Group” that categories the patient as either a “youth” or “adult”. This feature essentially discretizes the “Age” attribute, so “Age” can be deleted.
- ‘Weight_kg’ was deleted because weight is primarily used to calculate the “BMI” and there is already an attribute for BMI.
 - Height is also used to calculate BMI, but we kept the ‘Height’ attribute because we wanted to see if height affected the chances of an individual receiving a heart attack, whereas weight was already accounted for in the calculation for BMI.

Lastly, the class label ‘Heart_Attack’ was moved to the last column from the 19th column (column numbering starts at 1) to indicate that ‘Heart_Attack’ is a class.

5. Method

We will improve the Naive Bayes algorithm by addressing its weakness, the assumption of attribute independence, using Multiple Component Analysis [9]. MCA is a dimensionality reduction technique specifically designed for categorical variables, similar to how PCA (Principal Component Analysis) works for continuous data. PCA isn't the right choice here to begin with because it combines different values into a single number. if we do PCA before discretizing, it merges different categories into a single, unclear number, leading us to lose the ability to interpret what each value means. Instead, we should discretize first to keep things clear. MCA produces new numerical variables called dimensions that summarize the patterns in the categorical data. These dimensions retain meaning because they are based on how different categories relate to each other. Additionally, Naive Bayes is inherently designed for categorical data since it calculates probabilities using frequency tables. By transforming continuous data into categorical

and applying MCA, we can create components that are not only tailored for this probabilistic approach but also easier to interpret compared to the continuous outputs of PCA. In conclusion, MCA's output is easier to interpret than PCA's continuous values in our context and also helps structure the data in a way that makes sense for Naive Bayes classification.

First, we will discretize our data so that we'll be able to use MCA. Next, we will use the Pearson Chi-Square Test [10] to determine if two categorical attributes are dependent or independent. if the p-value is small, it indicates a significant relationship. Then, Cramér's V [10] will be calculated to measure the strength of this relationship, with values closer to 1 indicating a stronger association. Together, they will help us tell whether two attributes are correlated to each other.

Chi-Square Test checks if the variables are related, Cramér's V tells you how strong that relationship is.

Chi-Square Test:

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} = N \sum_{i=1}^n \frac{(O_i/N - p_i)^2}{p_i}$$

χ^2 = Pearson's cumulative test statistic, which asymptotically approaches a χ^2 distribution.

O_i = the number of observations of type i .

N = total number of observations

$E_i = Np_i$ = the expected (theoretical) count of type i , asserted by the null hypothesis that the fraction of type i in the population is p_i

n = the number of cells in the table.

Cramér's V:

$$V = \sqrt{\frac{\varphi^2}{\min(k-1, r-1)}}$$

- φ is the phi coefficient.
- χ^2 is derived from Pearson's chi-squared test
- n is the grand total of observations and
- k being the number of columns.
- r being the number of rows.

Our Naive Bayes model is Complement Naive Bayes. A model specifically designed to perform better on imbalance datasets, which is a common issue when dealing with real world data. In our case, the dataset has an 8-2 split between the two classes which is a significant imbalance that could greatly affect our classification. Traditional Naive Bayes can struggle in such situations because it tends to be biased towards the majority class. Complement Naive Bayes addresses this issue by modifying the way it calculates probabilities, making it more robust to class imbalance. Instead of assuming equal importance for all classes, it focuses on the complement of the minority class, effectively improving the model's sensitivity to underrepresented classes.

Approach 1: Full MCA Naive Bayes

We assume all the attributes have the utmost dependence with all the other attributes. We will use MCA on the whole entire dataset without doing any kind of processing to see which attributes are more or less correlated with each other than the others.

Approach 2: Hierarchical MCA

We assume that dependencies exist between attributes and then assess the strength of these dependencies. For each attribute, we normalize the Cramér's V values to a scale between 0 and 1. Using a threshold of around 0.5, we classify any value above it as indicating a moderate dependence, while values below 0.5 are considered sufficiently independent. Each attribute will have a list of other attributes it is dependent on based on these correlations.

Next, we apply MCA to generate combinations of these dependent attributes. After creating these individual combinations, we use MCA once again to create a global combination of all the data. This second step is essential because the subsets of data generated from the first MCA may still contain overlapping attributes, which could reintroduce the dependencies we aim to eliminate. By performing MCA a second time on the entire dataset, we can fully remove any remaining dependencies, ensuring that the final dataset is completely independent.

This hierarchical approach to MCA is beneficial because it allows us to break down complex dependencies step by step, ensuring that we don't overlook subtle interactions between attributes. By addressing dependencies at both the local (attribute-specific) and global (dataset-wide) levels, we maximize the clarity and interpretability of the data, ensuring that the relationships between attributes are captured accurately without introducing noise from hidden correlations.

6. Results

Complement Naive Bayes

Accuracy: 79.96% Recall: 0.9410

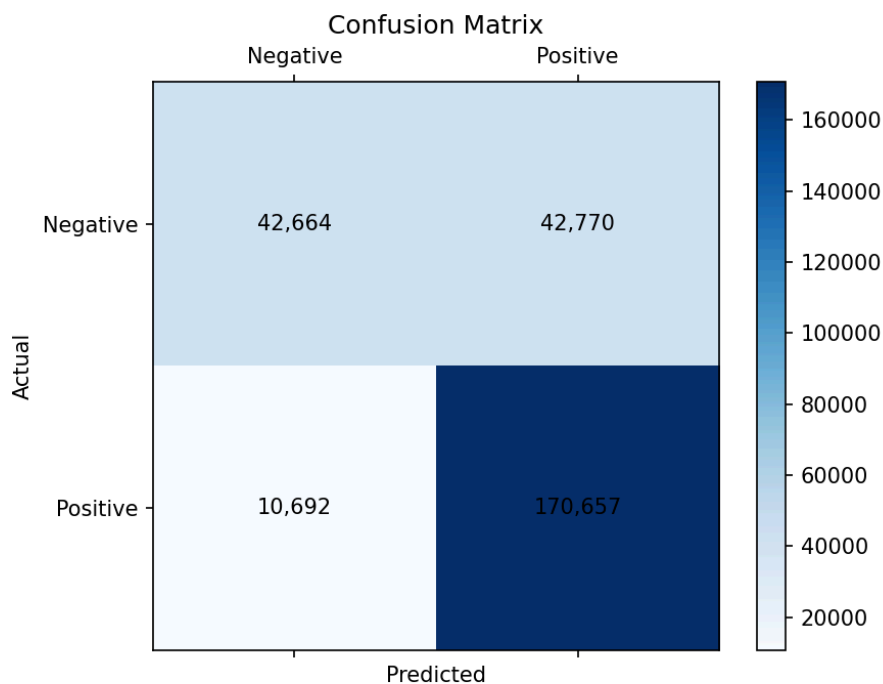


Figure 2. *Confusion matrix for Control: Conditional Naive Bayes. The matrix shows the classification results for the control group.*

	Sex	Height_cm	BMI	BP	Cholesterol_Level	Smoking_Status	Alcohol_Consumption	Physical_Activity_Hours	Diabetes	Family_History	Diet_Type	Stress_Level	Heart_Rate	Exercise_Induced_Pain	Age_Group	Region	Air_Pollution_Level	Income_Level	Education_Level	Health_Insurance	Regular_Checkups	Medication_Adherence
Sex	1	0.5022173422	0.902556168	0.561796392	0.394010796	0.844263038	0.569626293	0.712110103	0.599191568	0.254543009	0.906322986	0.196269189	0.515591154	0.765568822	0.482444483	0.7728769757	0.4496828505	0.7580070718	0.6523698422	0.7262583683	0.4565024323	0.5376428551
Height_cm	0.02201502987	1	0.02369784547	0.0228755162	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
BMI	0.02201502987	0.02369784547	1	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547	0.02369784547
BP	0.0228755162	0.0200042452	0.02369784547	1	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Cholesterol_Level	0.0228755162	0.0228755162	0.02369784547	0.0200042452	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Smoking_Status	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Alcohol_Consumption	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Physical_Activity_Hours	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Diabetes	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Family_History	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Diet_Type	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Stress_Level	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Heart_Rate	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Exercise_Induced_Pain	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Age_Group	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Region	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Air_Pollution_Level	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Income_Level	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Education_Level	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1	0.0228755162	0.0228755162
Health_Insurance	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	1
Regular_Checkups	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162
Medication_Adherence	0.0228755162	0.0228755162	0.02369784547	0.0200042452	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162	0.0228755162

Table 10. [Normalized Cramer's V outputs](#)

Plain Naive Bayes performed pretty well on this dataset. This tells us that most of the attributes in the dataset are already mostly independent of one another. We can also see this by the calculated Cramer's V outputs, most of them are under 0.015, meaning that they can be considered to be independent. However, in this project, we will try to further improve this accuracy and recall.

Complement Naive Bayes with Full MCA:

Accuracy: 85.61% Recall: 0.9597

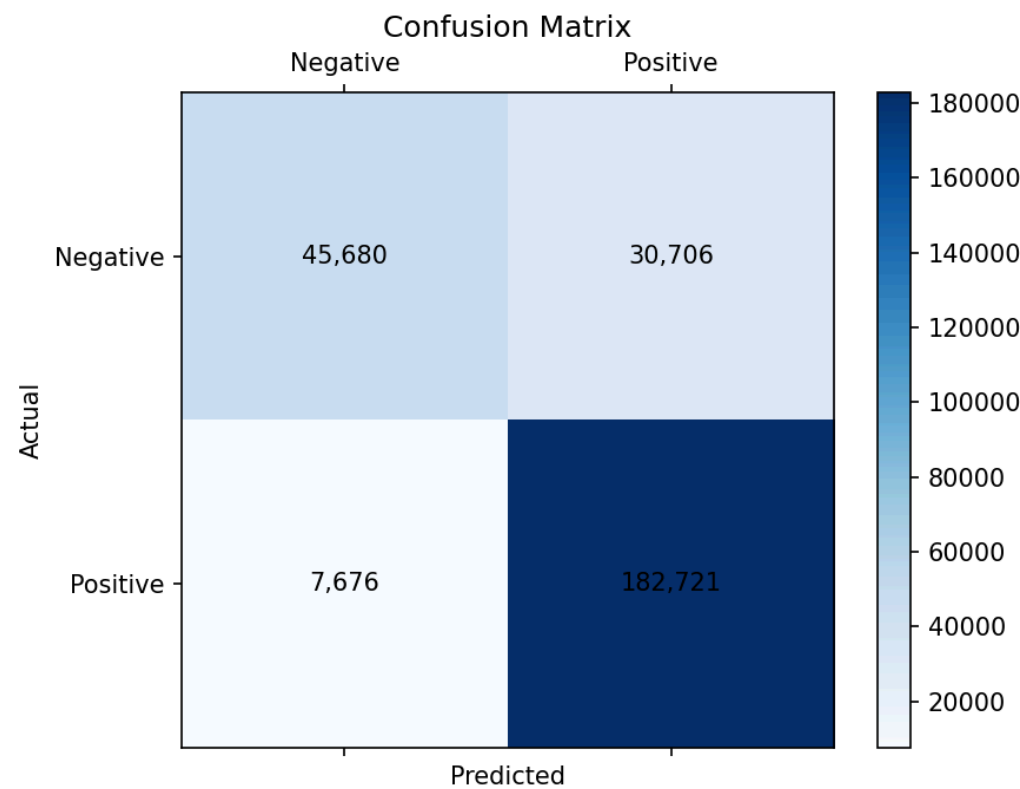


Figure 3: *Confusion Matrix for Approach 2: Full MCA*

Complement Naive Bayes with Cramer's V and Hierarchical MCA:

Accuracy: 88.73% Recall: 0.9692

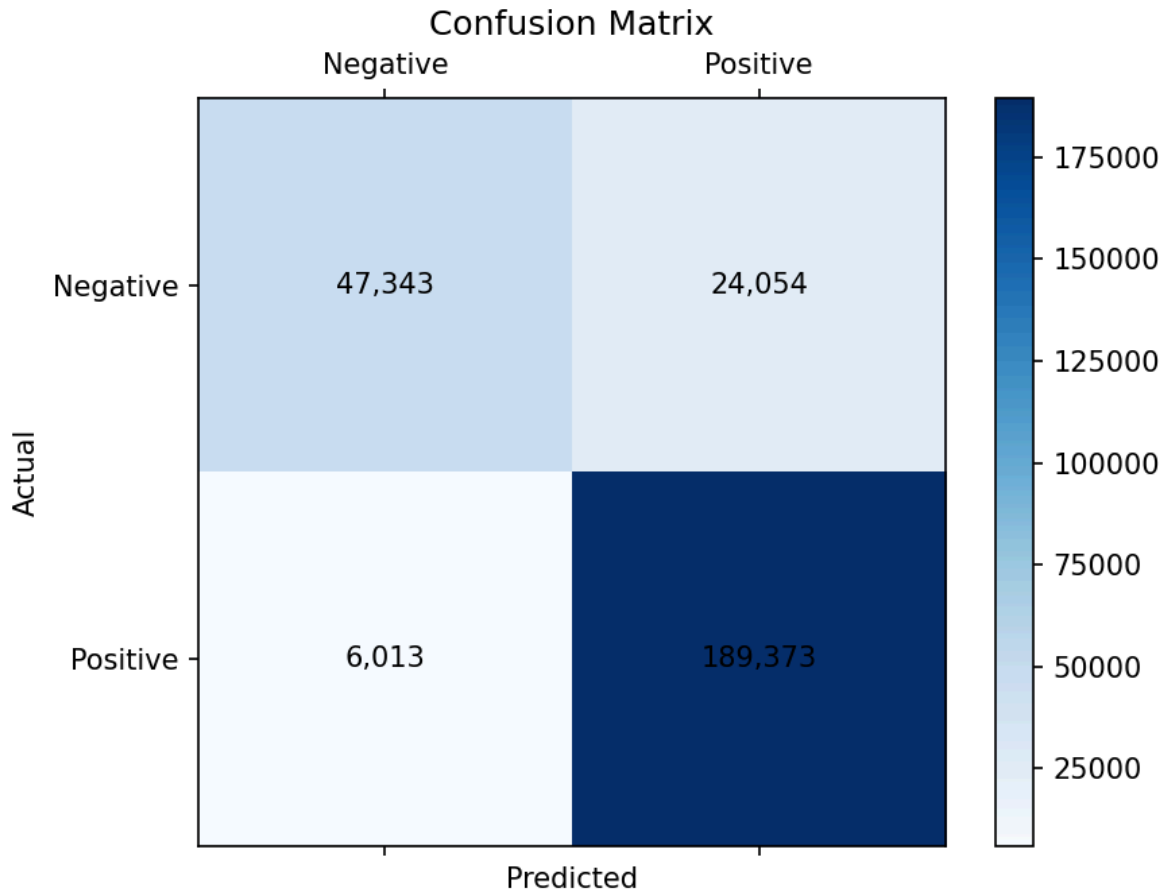


Figure 4: *Confusion Matrix for Approach 3: Hierarchical MCA*

7. Analysis

Overall, each modification of the Naive Bayes model—whether using full MCA or hierarchical MCA—resulted in improvements in performance compared to the original model, especially in terms of recall.

Initially, I had concerns that hierarchical MCA might reduce accuracy and recall, as the second round of MCA could potentially introduce noise by losing subtle relationships during the process of combining dependencies. Specifically, if the global combination fails to fully capture these relationships, the model could miss important patterns, leading to lower performance. However, contrary to this concern, the hierarchical MCA version performed slightly better than the full MCA version, achieving an accuracy of 89 and a recall of 96, compared to the full MCA model's

accuracy of 86 and recall of 95. This suggests that the hierarchical approach helped to further refine the dataset, allowing the model to better handle complex interactions between attributes.

Looking at the numbers, it's clear that each modification didn't dramatically improve the accuracy of the original Naive Bayes model (which had an accuracy of around 80 and recall of 94). However, considering the original model was already performing quite well, the improvements in accuracy and recall with each modification are notable. The full MCA and hierarchical MCA both contributed to slight but meaningful improvements in recall, which is critical when aiming to minimize false negatives, particularly in high-stakes scenarios like detecting heart attacks.

The model's performance suggests that while the dataset we used already contained many independent attributes, further testing on a dataset with more dependent attributes could reveal whether these improvements are consistent across different types of data. It's possible that hierarchical MCA could offer even more substantial gains when dealing with more complex, dependent relationships between attributes.

8. References

- [1] Panday, A. (2021). *Heart attack in youth vs adult in France*. Kaggle.
<https://www.kaggle.com/datasets/ankushpanday1/heart-attack-in-youth-vs-adult-in-france>
- [2] Luo, L.; Liu, T. (2024). Integrating advanced principal component analysis into naive bayes for enhanced classification performance. *Advances in Operation Research and Production Management*, 3, 27-31.
- [3] AirNow.gov, U.S. EPA. (n.d.). *Aqi Basics*. AQI Basics | AirNow.gov.
<https://www.airnow.gov/aqi/aqi-basics/>
- [4] Centers for Disease Control and Prevention. (n.d.). *Adult BMI categories*. Centers for Disease Control and Prevention.
<https://www.cdc.gov/bmi/adult-calculator/bmi-categories.html>
- [5] MediLexicon International. (n.d.). *What is the average height for men?*. Medical News Today.
<https://www.medicalnewstoday.com/articles/318155#what-is-the-link-between-height-and-weight>
- [6] professional, C. C. medical. (2025, February 6). *What should my cholesterol levels be?*. Cleveland Clinic.
<https://my.clevelandclinic.org/health/articles/11920-cholesterol-numbers-what-do-they-mean>

- [7] *Target heart rates chart*. www.heart.org. (2024a, August 12).
<https://www.heart.org/en/healthy-living/fitness/fitness-basics/target-heart-rates>
- [8] *Understanding blood pressure readings*. www.heart.org. (2024b, December 18).
<https://www.heart.org/en/health-topics/high-blood-pressure/understanding-blood-pressure-readings>
- [9] Abdi, Hervé & Valentin, Dominique. (2007). *Multiple Correspondence Analysis*. Encyclopedia of Measurement and Statistics.
- [10] McHugh M. L. (2013). *The chi-square test of independence*. Biochemia medica, 23(2), 143–149. <https://doi.org/10.11613/bm.2013.018>