



***Instituto Politécnico Nacional***

***Escuela Superior de Cómputo***

***Profesor:***

***Andres Garcia Floriano***

***Alumno:***

***Hernández Jiménez Erick Yael***

***Patiño Flores Samuel***

***Robert Garayzar Arturo***

***5BV1***

***Practica:***

***Clasificador Bayesiano Ingenuo***

# Índice

<b>1. Introducción</b>	<b>2</b>
1.1. Clasificador Bayesiano . . . . .	2
1.2. Clasificador Bayesiano Ingenuo . . . . .	2
1.3. Métodos de división de datos . . . . .	2
1.3.1. Estratificación . . . . .	2
1.3.2. Validación Leave-One-Out (LOO) . . . . .	2
1.3.3. n-fold Cross-Validation . . . . .	3
<b>2. Desarrollo</b>	<b>3</b>
<b>3. Conclusiones</b>	<b>5</b>

# 1. Introducción

## 1.1. Clasificador Bayesiano

El **clasificador bayesiano** es un tipo de modelo de aprendizaje automático basado en el teorema de Bayes, un principio fundamental de la teoría de probabilidad que describe la probabilidad de un evento basado en el conocimiento previo de condiciones relacionadas con ese evento. En términos simples, el clasificador bayesiano determina la probabilidad de que una instancia pertenezca a una clase específica dado un conjunto de características observadas. Esta probabilidad se calcula combinando las probabilidades de las características con las probabilidades a priori de las clases.

Matemáticamente, el teorema de Bayes se expresa como:

$$P(C|X) = \frac{P(X|C) \cdot P(C)}{P(X)}$$

donde: -  $P(C|X)$  es la probabilidad posterior de la clase  $C$  dado los datos  $X$ . -  $P(X|C)$  es la probabilidad de observar los datos  $X$  dado que pertenecen a la clase  $C$ . -  $P(C)$  es la probabilidad a priori de la clase  $C$ . -  $P(X)$  es la probabilidad de los datos  $X$ .

## 1.2. Clasificador Bayesiano Ingenuo

El **clasificador bayesiano ingenuo** (Naive Bayes) es una variante simplificada y popular del clasificador bayesiano que asume que todas las características son independientes entre sí, dado el valor de la clase. Aunque esta suposición de independencia rara vez se cumple en la realidad, el clasificador bayesiano ingenuo ha demostrado ser altamente eficaz en una amplia gama de aplicaciones, incluyendo la clasificación de texto, el filtrado de spam, el análisis de sentimientos y el diagnóstico médico.

El adjetivo 'ingenuo' se refiere a la fuerte hipótesis de independencia que hace el modelo. A pesar de esta limitación, Naive Bayes funciona bien debido a su simplicidad, rapidez y eficacia en contextos en los que muchas otras técnicas se complicarían o resultarían menos efectivas.

El proceso de clasificación con Naive Bayes implica calcular la probabilidad de que una instancia dada pertenezca a una clase específica utilizando las probabilidades individuales de sus características. Se selecciona la clase con la mayor probabilidad. La eficiencia de Naive Bayes y su capacidad para manejar tanto características numéricas como categóricas lo convierten en una opción atractiva para problemas de clasificación supervisada.

En particular, los datasets usados ya presentan las etiquetas mínimas necesarias para realizar predicciones, por lo que nos apegamos a usar la misma cantidad de características 'k' que contienen. Los datasets son 'iris', 'wine' y 'mushroom' o 'agaricus-lepiota'.

## 1.3. Métodos de división de datos

En esta práctica se usaron 3 métodos de división de datos para entrenamiento y prueba:

### 1.3.1. Estratificación

La **estratificación** es una técnica utilizada para mantener la proporción de las clases en los conjuntos de entrenamiento y prueba durante la partición de los datos. Esto es especialmente útil cuando las clases están desbalanceadas, es decir, cuando algunas categorías tienen muchas más observaciones que otras. La estratificación garantiza que cada subconjunto de datos mantenga la misma distribución de clases que el conjunto original, proporcionando una evaluación más representativa del rendimiento del modelo.

### 1.3.2. Validación Leave-One-Out (LOO)

La **validación Leave-One-Out (LOO)** es un caso particular de validación cruzada en el que el número de subconjuntos es igual al número de observaciones en el conjunto de datos. En cada iteración, se utiliza una única observación como conjunto de prueba y el resto de las observaciones como conjunto de entrenamiento. Esto se repite tantas veces como observaciones haya, lo que da lugar a evaluaciones altamente exhaustivas. Aunque ofrece una estimación imparcial del error de generalización, LOO puede ser computacionalmente costoso para conjuntos de datos grandes debido a la gran cantidad de particiones.

### 1.3.3. n-fold Cross-Validation

La **validación cruzada de n-fold** implica dividir el conjunto de datos en  $n$  subconjuntos (o 'folds') de aproximadamente el mismo tamaño. En cada iteración, uno de estos subconjuntos se utiliza como conjunto de prueba y el resto como conjunto de entrenamiento. Este proceso se repite  $n$  veces, con cada subconjunto utilizado como conjunto de prueba una vez. El rendimiento se evalúa promediando las métricas obtenidas en cada iteración. Este método proporciona un equilibrio entre la cantidad de datos de entrenamiento disponibles y el tiempo computacional, y su versión más común es la *10-fold cross-validation*.

## 2. Desarrollo

El clasificador bayesiano, como se observa en los resultados obtenidos, muestra un alto nivel de precisión en los conjuntos de datos Iris, Wine y Mushroom, tanto con la validación mediante particionado estratificado como con la técnica Leave-One-Out.

- **Precisión general:** En el caso del particionado estratificado, la precisión fue consistentemente alta: alrededor de 0,911 para Iris, 0,981 para Wine y 0,998 para Mushroom. Por otro lado, con Leave-One-Out, la precisión promedio alcanzó 0,953 para Iris, 0,966 para Wine y 0,997 para Mushroom, lo que confirma la robustez del modelo bayesiano en diversos contextos.
- **Matriz de confusión con particionado estratificado:** Esta matriz muestra los aciertos y errores de clasificación. Por ejemplo, para el conjunto de datos Iris, los valores indican la cantidad de elementos correctamente clasificados y los errores entre las clases. Las matrices de confusión para los conjuntos de datos Wine y Mushroom presentan un muy bajo número de errores de clasificación, evidenciando la alta capacidad predictiva.
- **Matriz de confusión con Leave-One-Out:** La matriz de confusión acumulada muestra cómo el clasificador mantiene un buen rendimiento al evaluar cada muestra individualmente. Los valores indican que las predicciones erróneas son mínimas en comparación con las predicciones correctas, confirmando su eficacia.

```
Iris - Estratificado
Matriz de Confusión:
[[15  0  0]
 [ 0 14  1]
 [ 0  3 12]]
Precisión: 0.9111111111111111

Precisión media (método Estratificado): 0.9111111111111111
-----

Wine - Estratificado
Matriz de Confusión:
[[18  0  0]
 [ 1 20  0]
 [ 0  0 15]]
Precisión: 0.9814814814814815

Precisión media (método Estratificado): 0.9814814814814815
-----

Mushroom - Estratificado
Matriz de Confusión:
[[1044  3]
 [  0 647]]
Precisión: 0.9982290436835891

Precisión media (método Estratificado): 0.9982290436835891
```

Figura 1: Resultados de la clasificación con el Bayesiano Ingenuo con estratificación 30/70.

```

***Leave One Out***
Iris - Leave-One-Out
Precisión media: 0.9533333333333334
Desviación estándar de la precisión: 0.21092389359408498
Matriz de Confusión Global Acumulada:
[[50.  0.  0.]
 [ 0. 47.  3.]
 [ 0.  4. 46.]]
Wine - Leave-One-Out
Precisión media: 0.9662921348314607
Desviación estándar de la precisión: 0.18047616184504506
Matriz de Confusión Global Acumulada:
[[57.  2.  0.]
 [ 1. 67.  3.]
 [ 0.  0. 48.]]
Mushroom - Leave-One-Out
Precisión media: 0.997165131112686
Desviación estándar de la precisión: 0.05316796409216457
Matriz de Confusión Global Acumulada:
[[3472.  16.]
 [  0. 2156.]]

```

Figura 2: Resultados de la clasificación con el Bayesiano Ingenuo con método Leave-One-Out.

```

10-fold cross-validation
Iris - 10-fold Cross-Validation
Precisión media: 0.9533333333333334
Desviación estándar de la precisión: 0.04268749491621898
Matrices de Confusión por Fold:
Fold 1:
[[5 0 0]
 [0 4 1]
 [0 0 5]]

Fold 2:
[[5 0 0]
 [0 5 0]
 [0 1 4]]

Fold 3:
[[5 0 0]
 [0 5 0]
 [0 0 5]]

Fold 4:
[[5 0 0]
 [0 5 0]
 [0 1 4]]

Fold 5:
[[5 0 0]
 [0 4 1]
 [0 0 5]]

Fold 6:
[[5 0 0]
 [0 4 1]
 [0 0 5]]

Fold 7:
[[5 0 0]
 [0 5 0]
 [0 2 3]]

Fold 8:
[[5 0 0]
 [0 5 0]
 [0 0 5]]

Fold 9:
[[5 0 0]
 [0 5 0]
 [0 0 5]]

Fold 10:
[[5 0 0]
 [0 5 0]
 [0 0 5]]

```

Figura 3: Resultados de la clasificación con el Bayesiano Ingenuo con 10-fold cross-validation en Iris

```

Wine - 10-fold Cross-Validation
Precisión media: 0.9666666666666666
Desviación estándar de la precisión: 0.027216552697590882
Matrices de Confusión por Fold:
Fold 1:
[[6 0 0]
 [0 6 1]
 [0 0 5]]

Fold 2:
[[6 0 0]
 [0 7 0]
 [0 0 5]]

Fold 3:
[[6 0 0]
 [1 6 0]
 [0 0 5]]

Fold 4:
[[6 0 0]
 [0 6 1]
 [0 0 5]]

Fold 5:
[[5 1 0]
 [0 7 0]
 [0 0 5]]

Fold 6:
[[6 0 0]
 [0 7 0]
 [0 1 4]]

Fold 7:
[[6 0 0]
 [0 7 0]
 [0 0 5]]

Fold 8:
[[5 1 0]
 [0 7 0]
 [0 0 5]]

Fold 9:
[[6 0 0]
 [0 7 0]
 [0 0 4]]

Fold 10:
[[5 0 0]
 [0 8 0]
 [0 0 4]]

```

Figura 4: Resultados de la clasificación con el Bayesiano Ingenuo con 10-fold cross-validation en Wine

```

Mushroom - 10-fold Cross-Validation
Precisión media: 0.9443262411347518
Desviación estándar de la precisión: 0.09493749581839618
Matrices de Confusión por Fold:
Fold 1:
[[349 0]
 [ 0 216]]

Fold 2:
[[349 0]
 [ 0 216]]

Fold 3:
[[349 0]
 [ 0 216]]

Fold 4:
[[349 0]
 [ 0 216]]

Fold 5:
[[349 0]
 [ 0 215]]

Fold 6:
[[253 96]
 [ 0 215]]

Fold 7:
[[349 0]
 [ 0 215]]

Fold 8:
[[341 8]
 [ 0 215]]

Fold 9:
[[348 0]
 [ 45 171]]

Fold 10:
[[348 0]
 [165 51]]

```

Figura 5: Resultados de la clasificación con el Bayesiano Ingenuo con 10-fold cross-validation en Mushroom

### 3. Conclusiones

#### Hernández Jiménez Erick Yael:

Con los resultados anteriormente vistos, podemos notar la gran mejora en la precisión que nos otorga este modelo comparado con otros modelos implementados en prácticas anteriores y que, pese a su relativa simplicidad, la asunción que se hace de que las características se encuentran no relacionadas entre sí ayuda demasiado a no solo clasificar adecuadamente las clases, sino también encontrar más patrones entre sí.

#### Patiño Flores Samuel:

Naive Bayes se distingue por la suposición de independencia entre las características, lo cual permite simplificar los cálculos de probabilidad. Aunque en muchos casos esta suposición no se cumple completamente, el modelo sigue mostrando buenos resultados en varias aplicaciones prácticas, como la clasificación de texto, detección de spam, y

análisis de sentimientos. La simplicidad de Naive Bayes no solo lo hace computacionalmente eficiente, sino también fácil de interpretar, lo cual es una ventaja en situaciones donde la explicabilidad es importante.

### **Robert Garayzar Arturo:**

En esta práctica, implementamos y evaluamos el clasificador Naïve Bayes usando tres métodos de validación: Hold-Out estratificado, validación cruzada estratificada de 10 pliegues y Leave-One-Out. Esto permitió observar cómo varía el rendimiento del modelo en diferentes escenarios. El clasificador demostró ser efectivo, aunque su precisión depende de la estructura de cada conjunto de datos. Las métricas de Accuracy y la matriz de confusión brindaron una visión clara sobre el desempeño, confirmando que Naïve Bayes es una opción útil para tareas de clasificación básicas.