



Instituto Tecnológico
de Buenos Aires

Trabajo Final

82.05 - Análisis Predictivo

Ian Dalton - 62345

Caso de negocio



Introducción y caso de negocio

Predecir si un hongo es **comestible** o no.

Un hongo puede ser tanto:

- **Venenoso**
 - ◆ Muerte
 - ◆ Alucinogenos
 - ◆ Dolor de panza

- **Comestible**

Modelo de Predicción

Objetivo: Predecir si un hongo es **comestible** a partir de las características brindadas por el usuario

Unidad observacional: Características del hongo

Variable Target: “class”

Variables predictoras:

'cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor', 'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color', 'stalk-shape', 'stalk-root',
'stalk-surface-above-ring', 'stalk-surface-below-ring', 'stalk-color-above-ring',
'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number', 'ring-type', 'spore-print-color', 'population', 'habitat'

Modelo: Clasificación.

Dataset

Información del Dataset

Contiene información de las **características** de hongos, que pueden ser **dimensiones** o **descripciones**.

Cuenta con **8.124** registros y **22** variables.

Limpieza de la base

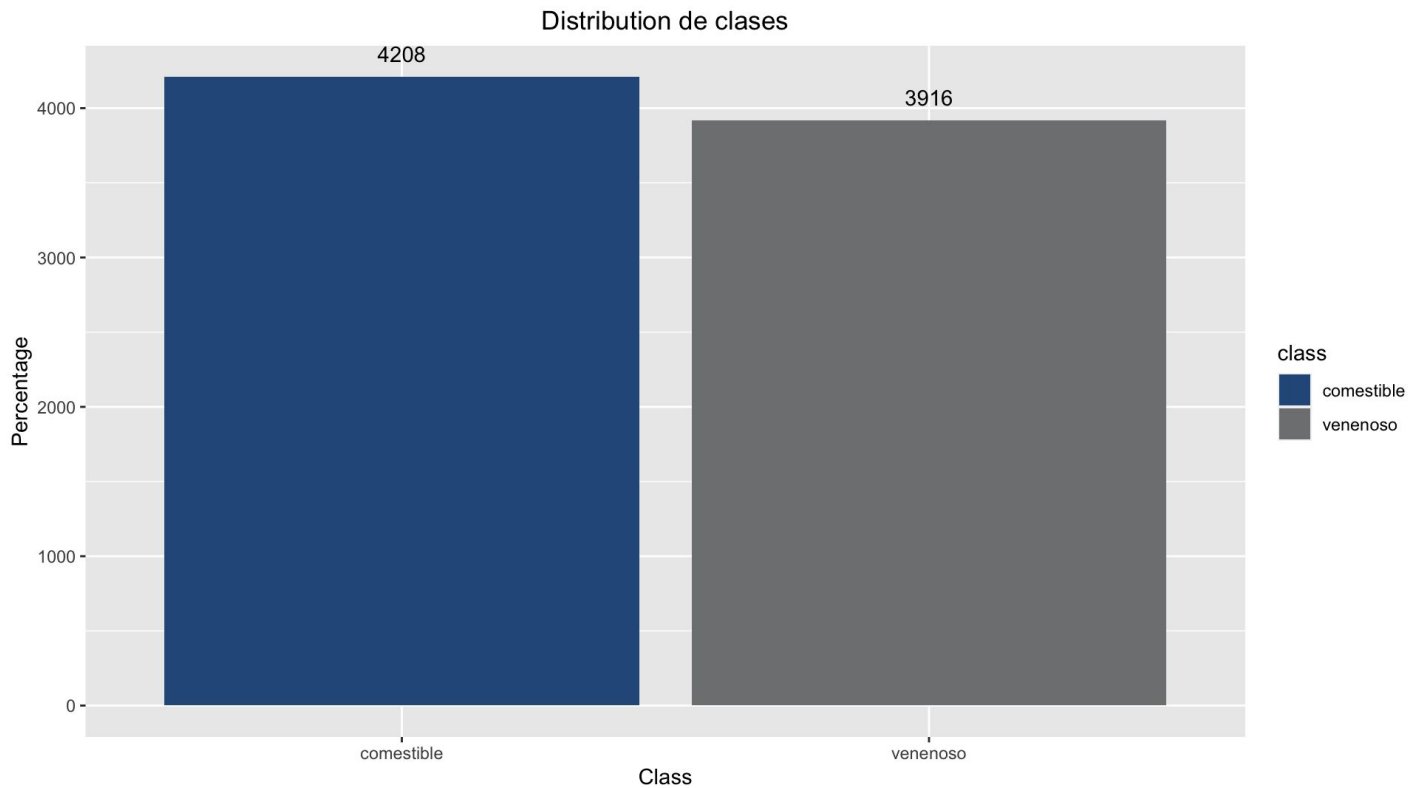
Todas las columnas son categoricas

Frecuencia de actualización: No hay

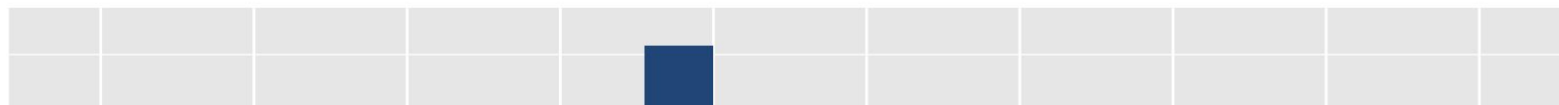
Análisis exploratorio



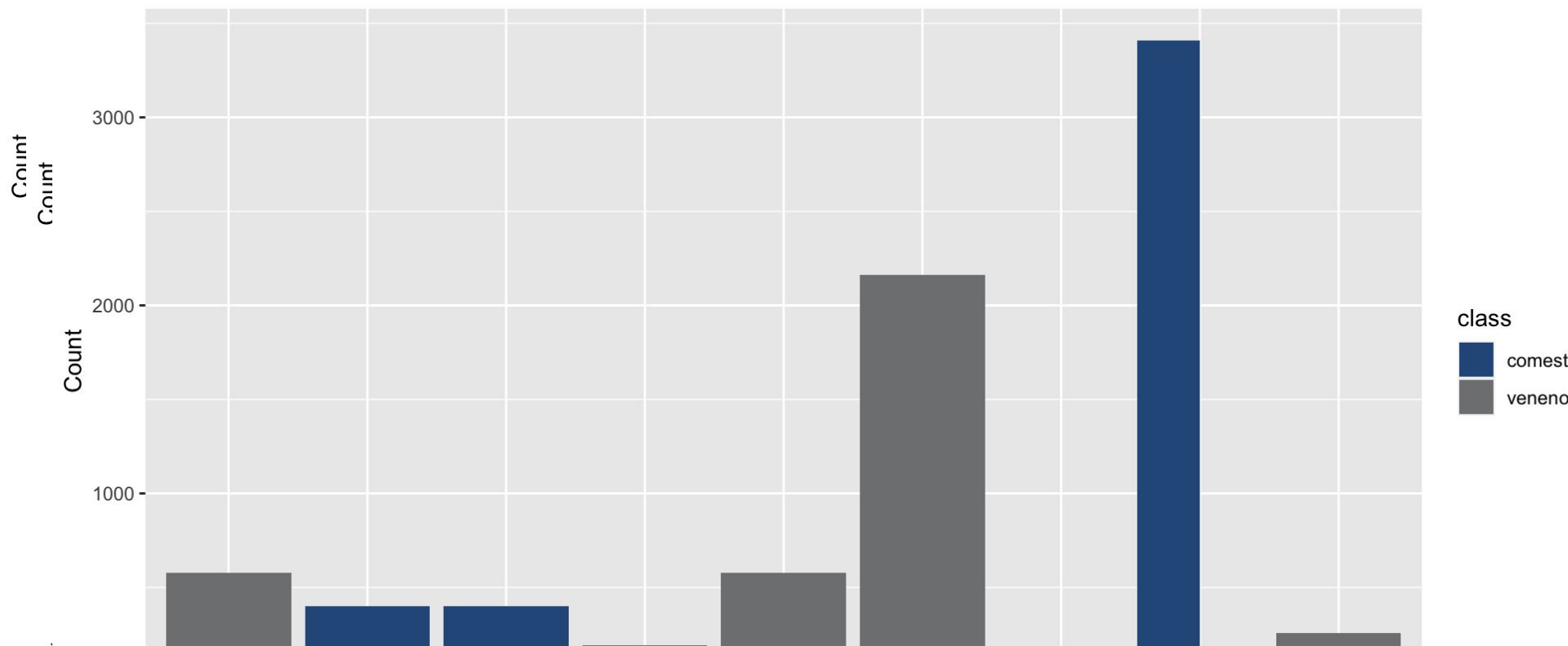
Tipos de hongos



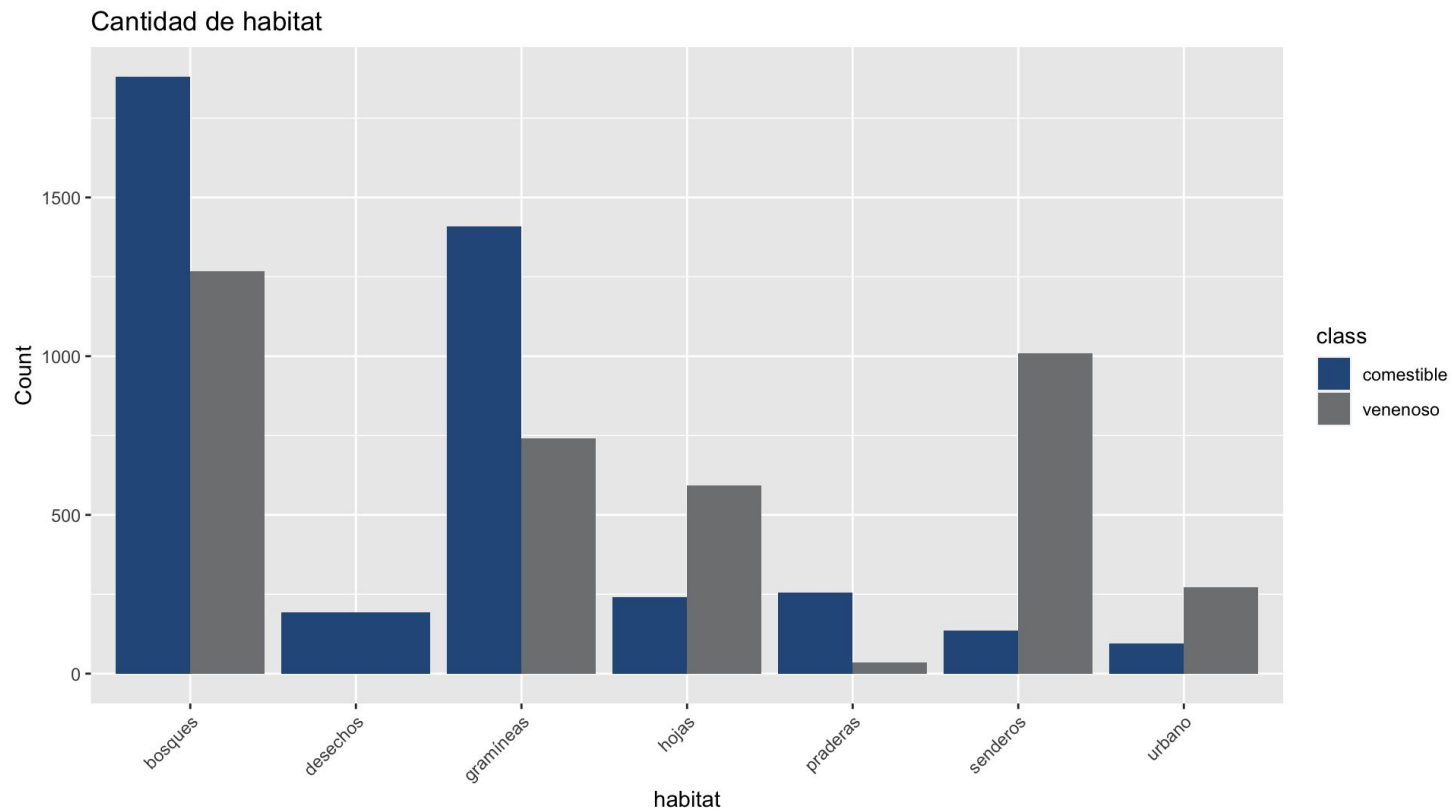
Cantidad de can.shana
Cantidad de can.surface
Cantidad de cap.color



Cantidad de odor

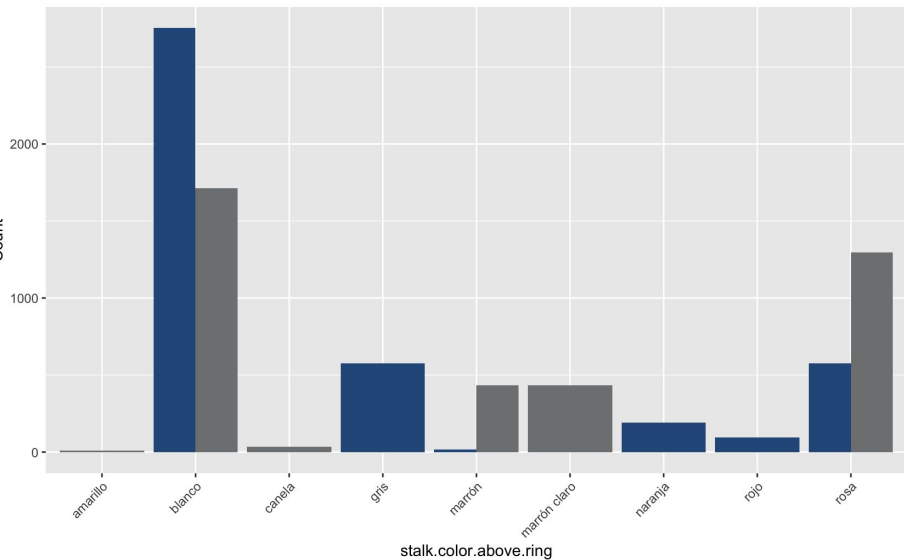


Factores geograficos

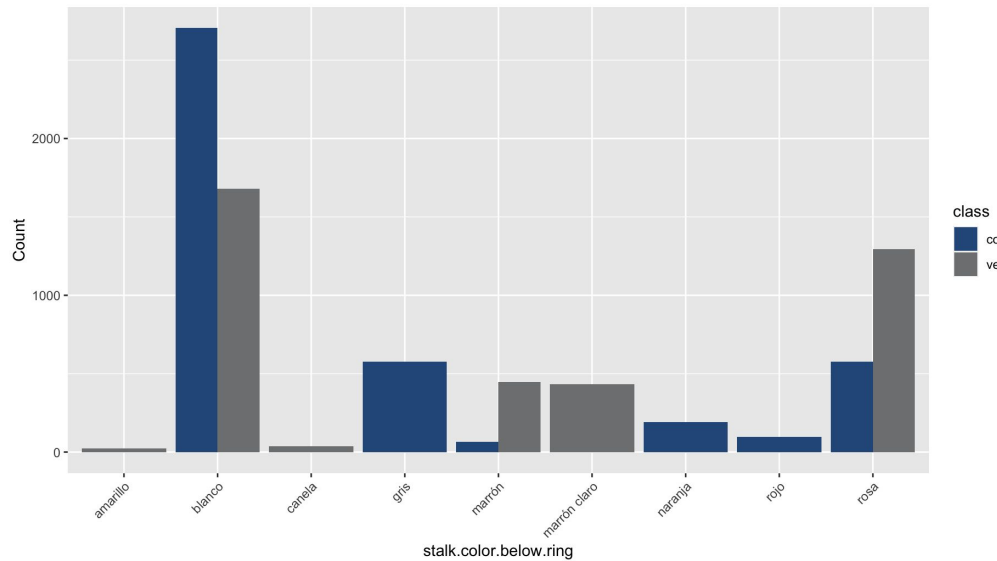


Colores del tallo

Cantidad de stalk.color.above.ring

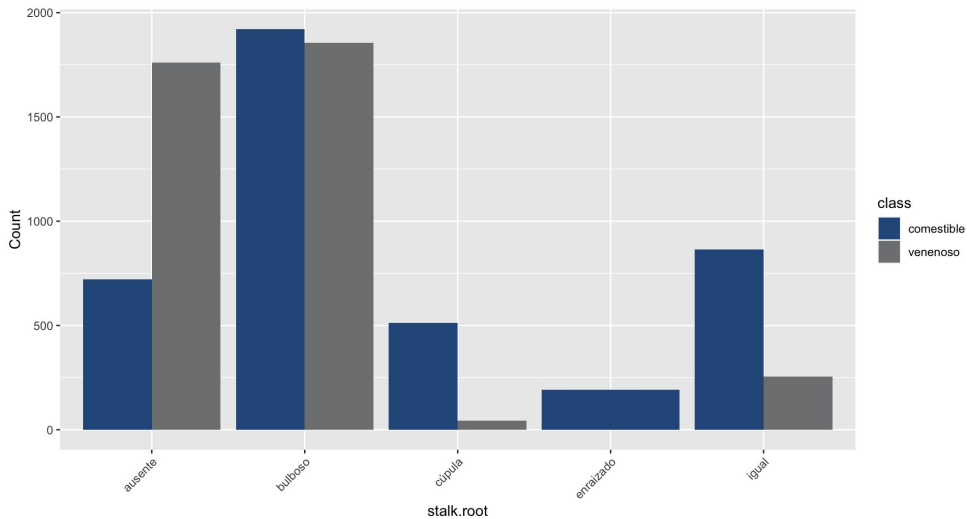


Cantidad de stalk.color.below.ring

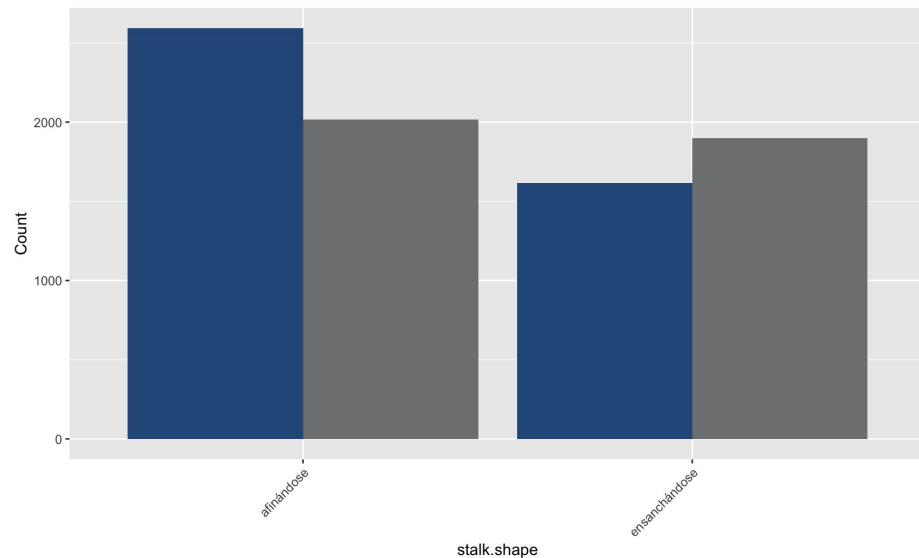


Detalles del tallo

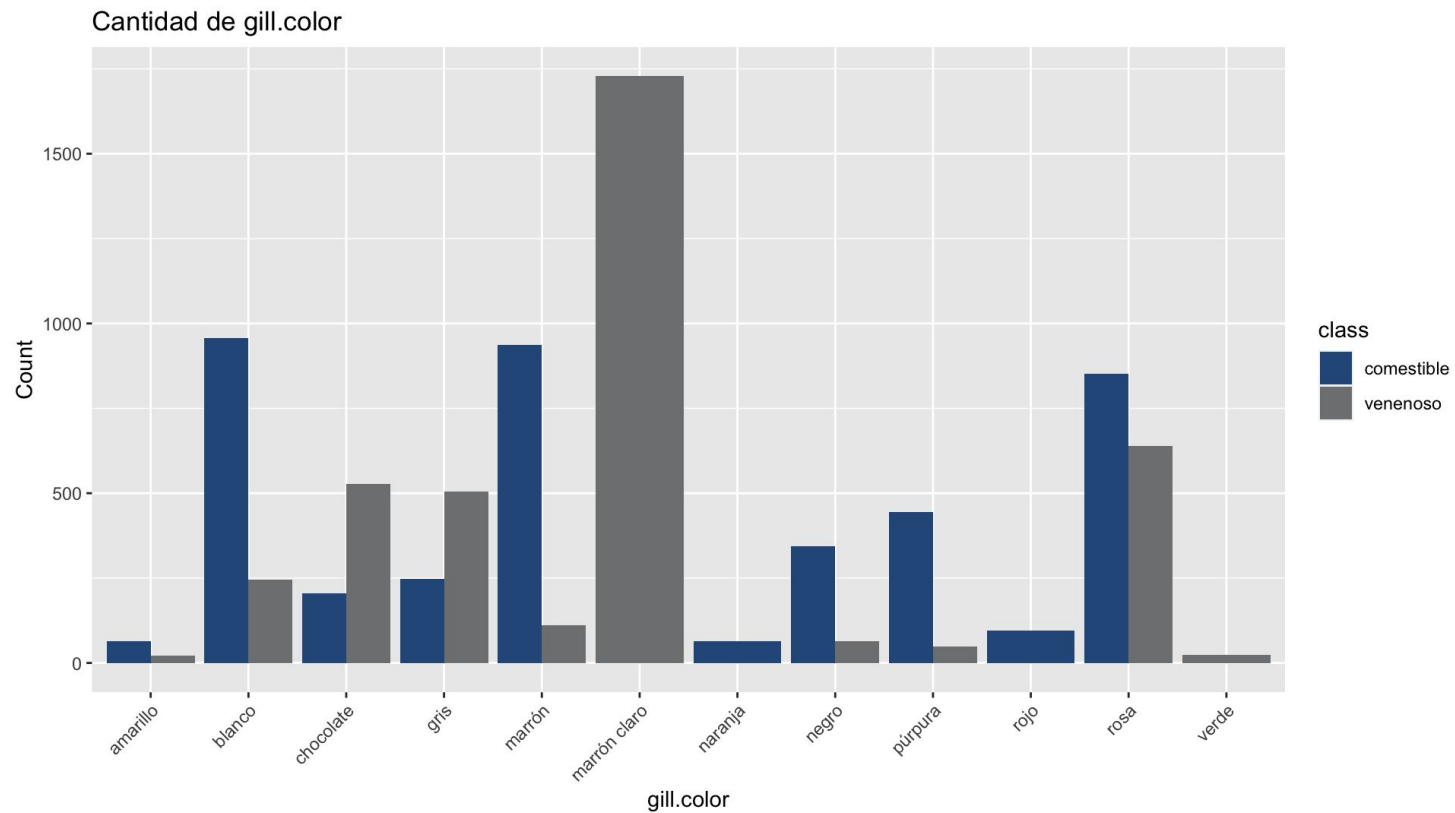
Cantidad de stalk.root



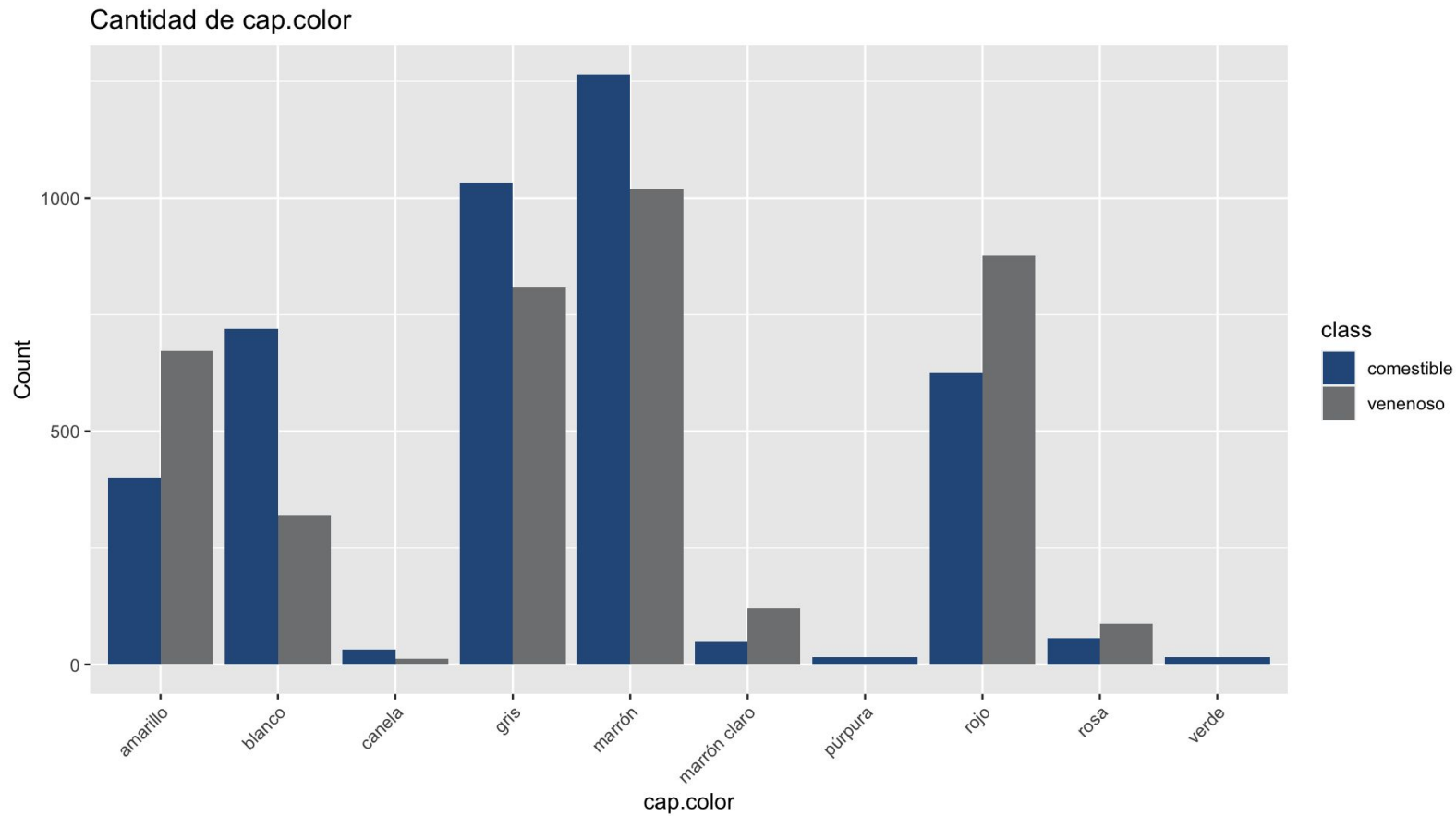
Cantidad de stalk.shape



Colores



Colores

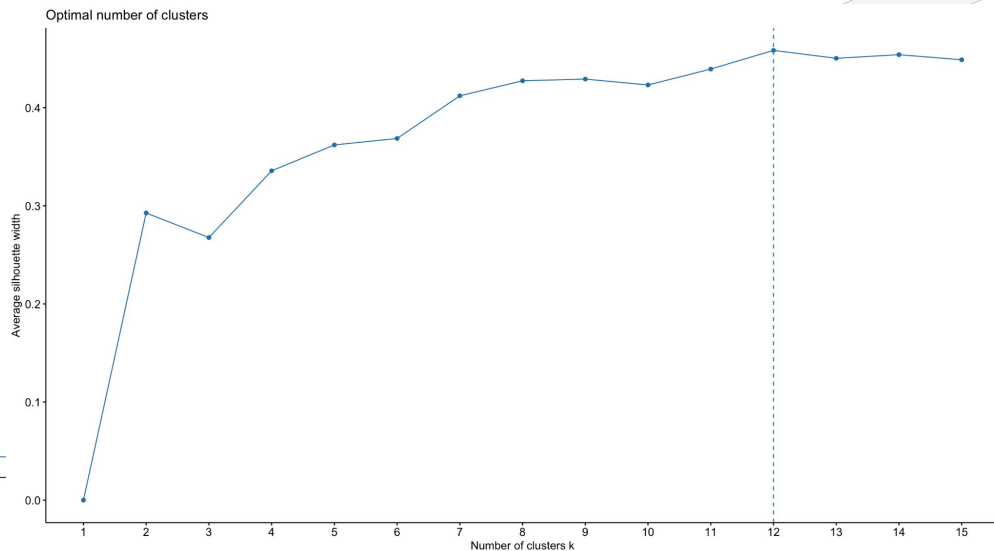
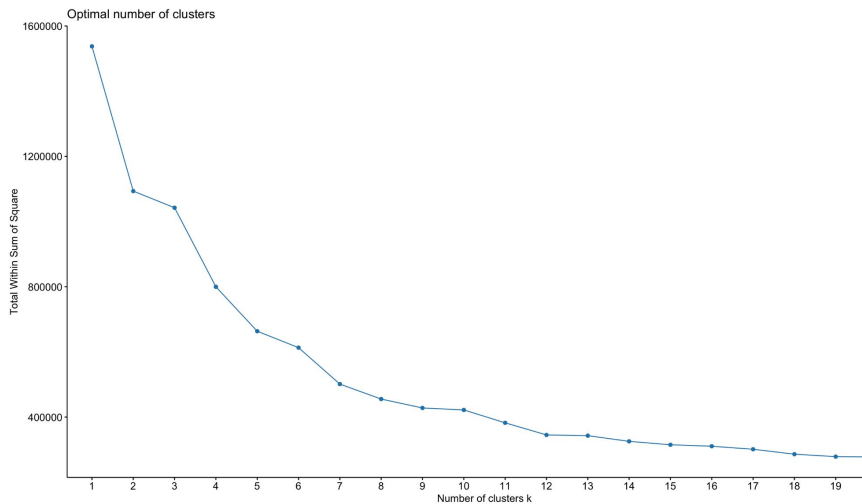


Clustering

Se realizó con el objetivo de comparar y evaluar diferencias entre las diferentes características que puedan ayudar a identificar si el hongo es comestible o no.

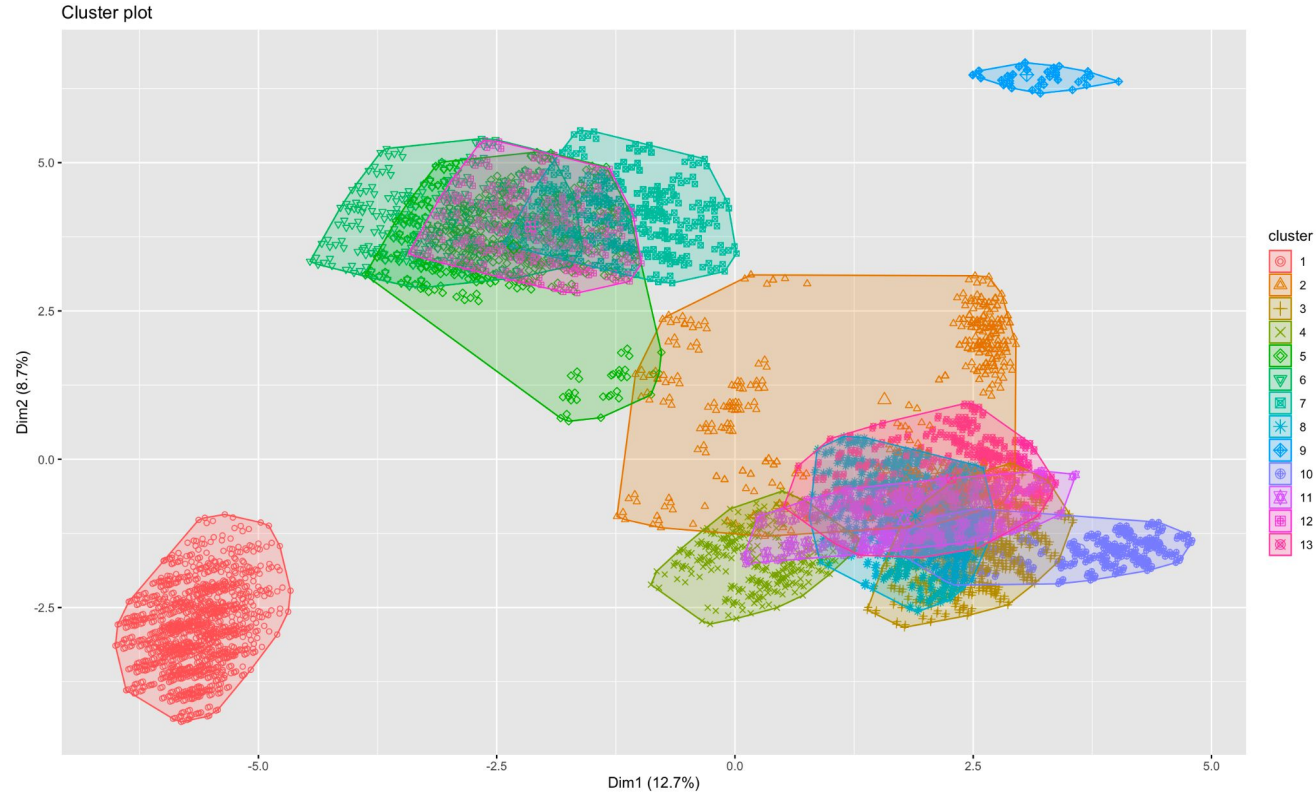
Medimos el valor del **estadístico de Hopkins** para evaluar si es conveniente la clusterización, obteniendo un valor cercano al **81,4%**.

Clustering - Agrupaciones



Número óptimo de clusters: 13

Visualización de Clusters



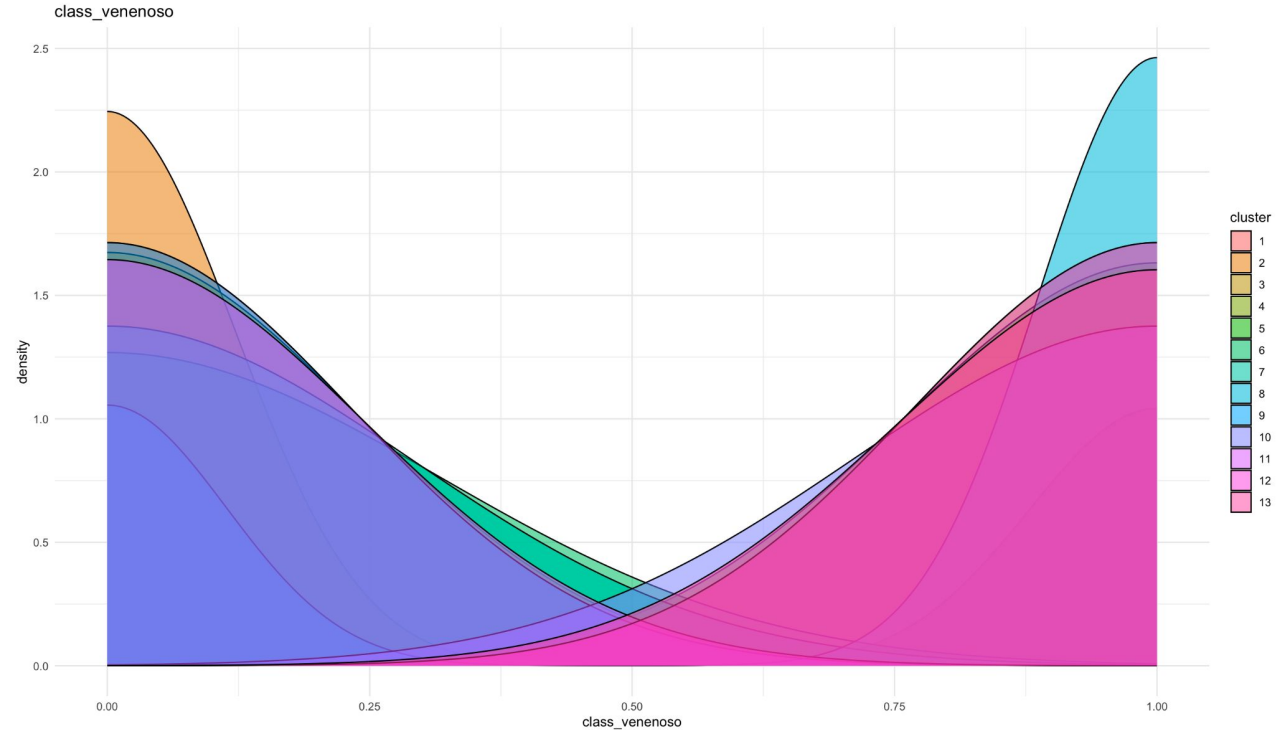
Análisis de clusters

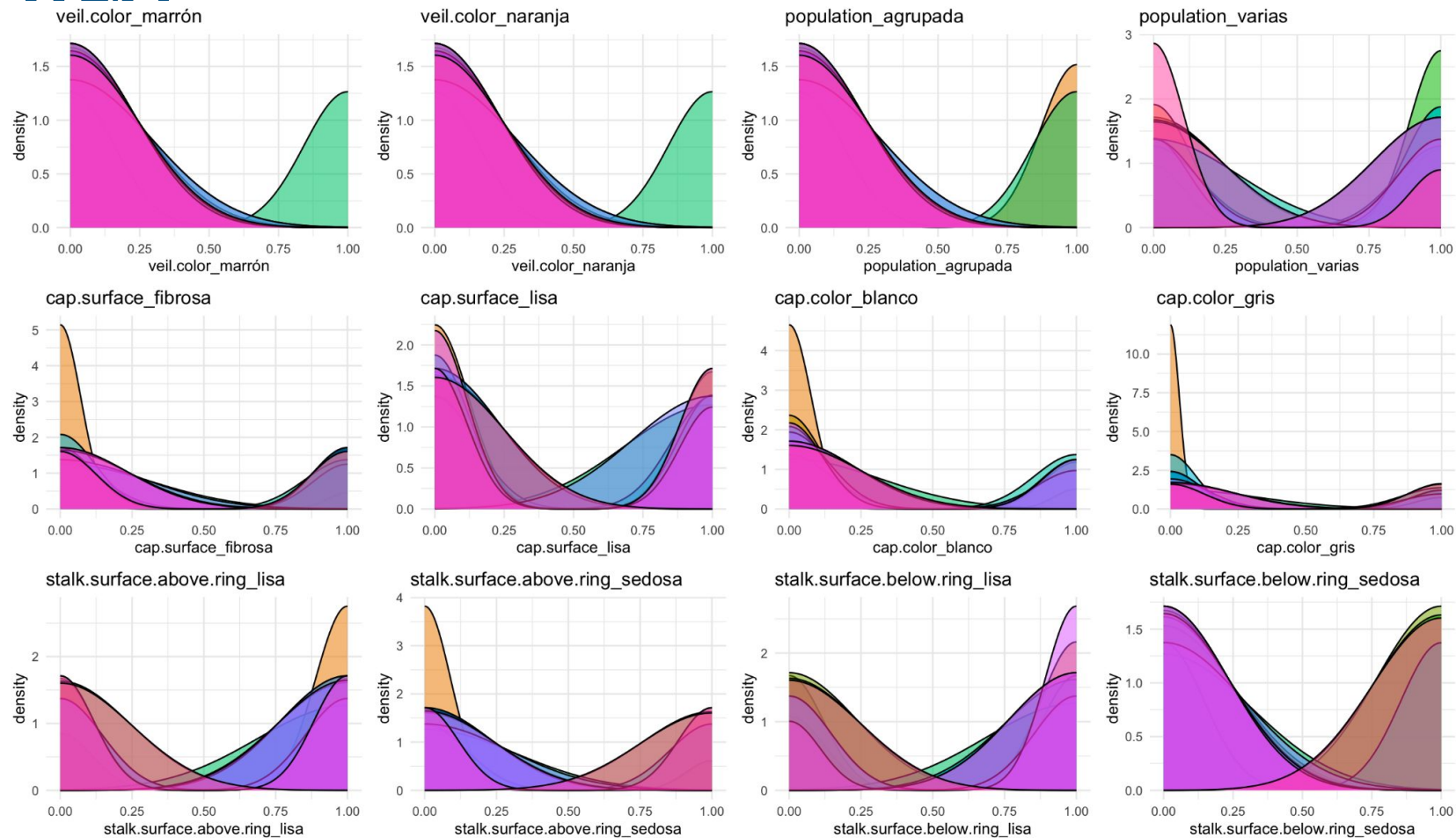
Cluster 2:

- Comestibles
- Se encuentran agrupados
- Tallo liso

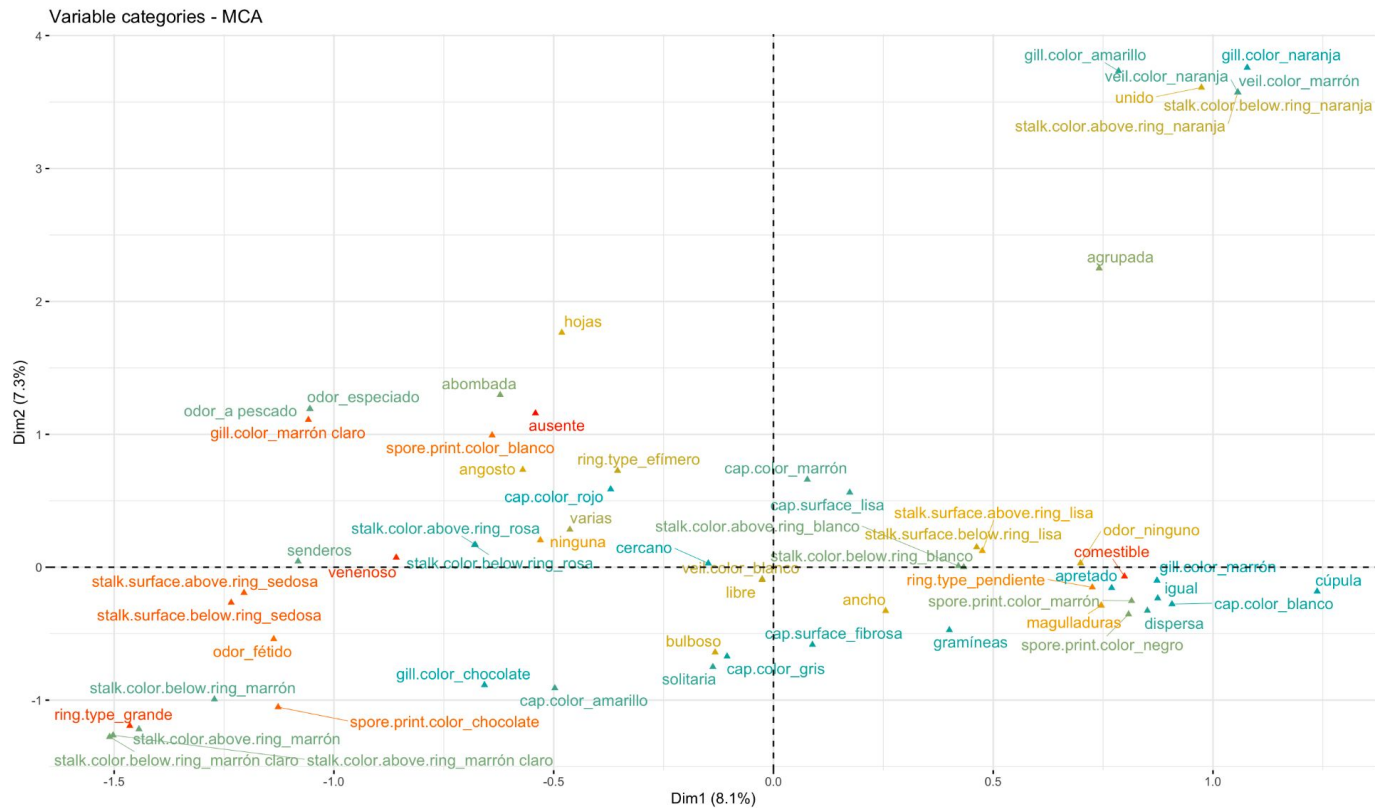
Cluster 8:

- Venenosos
- Cabeza lisa
- No se encuentra en grupos





MCA

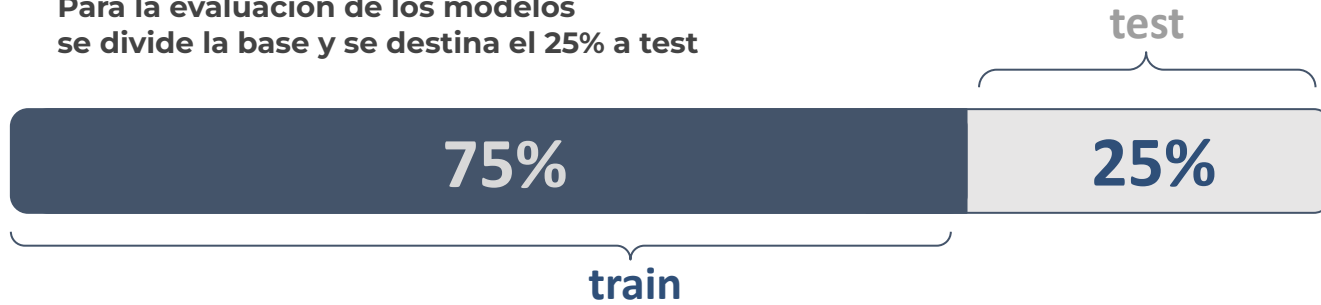


Evaluación



Partición de datos

Para la evaluación de los modelos
se divide la base y se destina el 25% a test



Modelo Baseline



¿Qué modelo use como baseline?

Features:

- Colores
- Forma de la cabeza

Modelo utilizado: LogisticRegression

Rendimiento alcanzado: 0.673

0	656	384
1	280	711
	0	1

Selección de modelos



Modelos evaluados

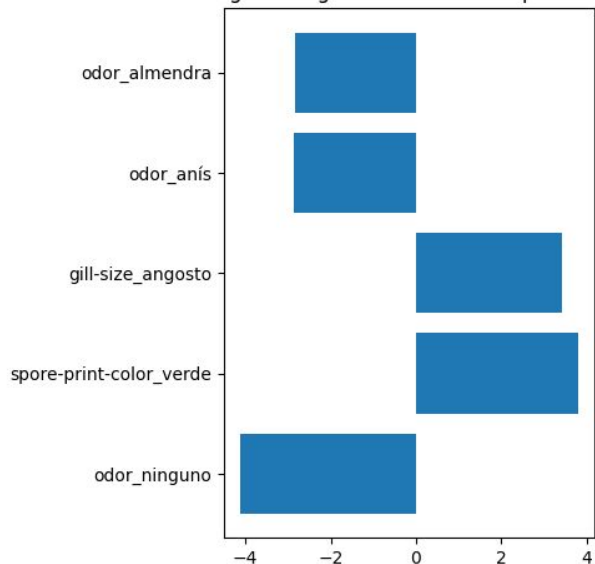
- Árboles
 - ◆ XGBoost
 - ◆ Decision Tree Classifier
 - ◆ Random Forest
- Support Vector Classifier
- KNN
- Logistic Regression

Tratamiento de variables categóricas

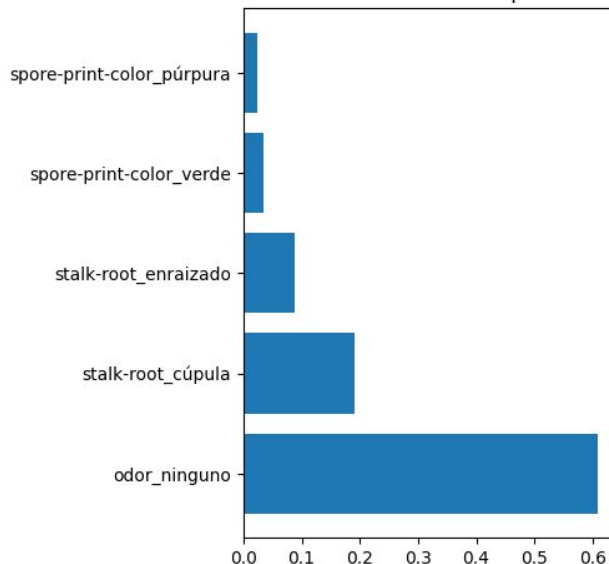
- One hot encoding
 - Se eliminó la primer columna
- Flags de diferencia en colores

Feature importance

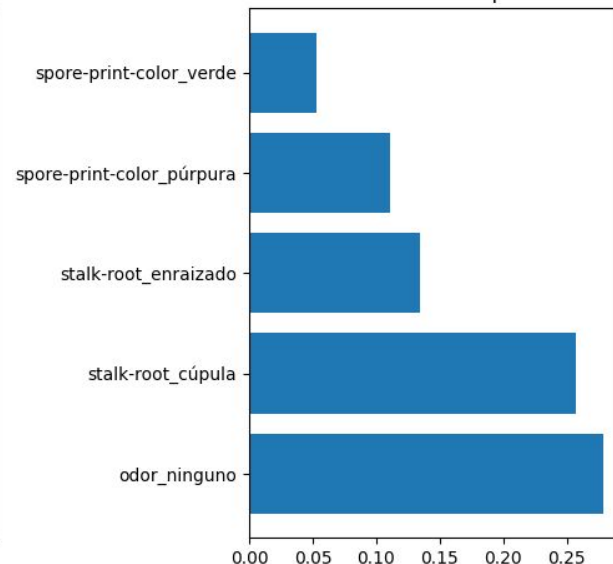
Logistic Regression Feature Importance



Decision Tree Feature Importance



XGBoost Feature Importance



Modelo Final



Modelo Final

XGBoost

Score: 1.0

Hiper Parámetros

- Booster Type: 'gbtree'
- Learning Rate: 0.3
- n_estimators: 100
- max_depth: 6
- min_child_weight: 1

Limitaciones y posibles mejoras

Limitaciones

- Inexperiencia en modelos predictivos
- Dataset limitado

Posibles mejoras

- Aumentar el tamaño del dataset para encontrar edge cases

Conclusiones

Conclusiones

Objetivo: Predecir si un hongo es **comestible** o no.

→ EDA

- ◆ Es muy fácil confundir hongos ya que **comparten** muchas **características**.
- ◆ Los comestibles **no** suelen tener olores.
- ◆ La forma que tienen los hongos de **producir esporas impacta** bastante en si es comestible o no.

Conclusiones

→ Modelo predictivo

◆ Herramientas utilizadas

- OneHotEncoder
- Flags en variables importantes

◆ Modelo final:

- El **olor** es una variable relevante.

→ Metricas evaluadas

- ◆ Score
- ◆ Matriz de confusión



Instituto Tecnológico
de Buenos Aires

Gracias!