

Claro, aquí tienes un informe ejecutivo basado en el EDA que has realizado en el notebook:

Informe Ejecutivo: Exploración y Análisis de Datos del Dataset de Hongos

Introducción

Este informe presenta un resumen del análisis exploratorio de datos (EDA) realizado sobre el dataset de hongos *Agaricus-Lepiota*. El objetivo principal de este EDA es comprender las características y la estructura del dataset, identificar posibles problemas de calidad de los datos y preparar los datos para futuros análisis o modelos de machine learning.

1. Descripción del Dataset

- El dataset contiene 8124 instancias descritas por 23 características, incluyendo la clase (venenoso o comestible) y diversas características morfológicas de los hongos.
- Las características son nominales (categóricas), representando atributos como la forma del sombrero, el color, el olor, etc.
- La variable objetivo es 'class', que indica si el hongo es venenoso ('p') o comestible ('e').

2. Calidad de los Datos

- Valores Nulos: Se identificaron valores faltantes en la columna 'stalk-root', representados por el símbolo '?'. Estos valores fueron imputados utilizando la moda de la columna.
- Tipos de Datos: Todas las columnas se identificaron inicialmente como de tipo 'object'. Para un procesamiento adecuado, se convirtieron las columnas categóricas al tipo 'category'.
- Valores Únicos: Se exploró la cantidad de valores únicos por columna para entender la variabilidad de los datos. La columna 'veil-type' tiene un único valor, lo que sugiere que podría no ser informativa para los modelos predictivos.

3. Transformación de Datos

- Codificación One-Hot: Las variables categóricas se codificaron utilizando one-hot encoding para convertirlas en un formato numérico adecuado para los algoritmos de machine learning. Se manejaron las categorías desconocidas para evitar errores durante la predicción.
- Reducción de Dimensionalidad (PCA): Se aplicó el Análisis de Componentes Principales (PCA) para reducir la dimensionalidad del dataset a dos componentes principales. Esto permitió visualizar la distribución de las clases

en un espacio bidimensional, lo que reveló una separación relativamente clara entre las clases de hongos venenosos y comestibles.

4. Modelado Predictivo

- **Modelo de Clasificación:** Se entrenó un modelo de Bosque Aleatorio (Random Forest) para predecir la clase de los hongos. El modelo alcanzó una precisión del 100% en el conjunto de prueba, lo que indica un rendimiento excelente en la clasificación de hongos venenosos y comestibles.
- **Evaluación del Modelo:** Se generó un informe de clasificación que muestra métricas de precisión, recall y F1-score para cada clase. Todas las métricas fueron del 100%, lo que confirma la alta precisión del modelo.

5. Conclusiones

- El dataset de hongos es de alta calidad, con pocos valores faltantes que se manejaron adecuadamente.
- Las características morfológicas proporcionan información suficiente para distinguir entre hongos venenosos y comestibles, como lo demuestra el alto rendimiento del modelo de Bosque Aleatorio.
- La visualización mediante PCA sugiere que las clases son inherentemente separables, lo que facilita la tarea de clasificación.

Recomendaciones

- Considerar la eliminación de la columna 'veil-type' debido a su falta de variabilidad.
- Explorar otros modelos de clasificación para comparar su rendimiento con el Bosque Aleatorio.
- Realizar una validación cruzada más exhaustiva para asegurar la generalización del modelo.
- Investigar la importancia de las características para entender cuáles son los atributos morfológicos más determinantes para la clasificación de los hongos.

Este informe proporciona una visión general del EDA realizado y sienta las bases para futuros análisis y modelado predictivo más profundos.