

# Hoja de resumen – Preprocesamiento y EDA en ML

(*Fundamentos de Machine Learning – caso Iris*)

*Jean Carlo Londoño Ocampo*

---

## 🎯 Objetivo del laboratorio

Construir una **buena representación de los datos**:

$(X, y)$

antes de entrenar cualquier modelo.

👉 El modelo aprende sobre **X**, no sobre la realidad.

---

## 1 Carga del dataset

Se construye:

- **X** → variables (features)
- **y** → etiqueta (clase)

En Iris:

- X: medidas de sépalo y pétalo
  - y: especie
- 

## 2 Exploración inicial (EDA)

Sirve para entender la forma de la distribución de X.

Incluye:

- `info()`

- `describe()`
- valores únicos

👉 No es estética, es para detectar problemas de calidad.

---

## 3 Calidad de datos (checks básicos)

Se revisa principalmente:

- Completitud → valores nulos
- Validez → rangos y tipos
- Consistencia básica

Esto garantiza que:

$(X,y) \approx \text{feno}'$  meno real observado  $(X, y) \approx \text{feno}'$  real observado

---

## 4 Visualización

Se usan gráficas para detectar:

- ✓ patrones predominantes
- ✓ separaciones naturales
- ✓ estructuras

Ejemplos usados:

- scatter 3D
- pairplot
- mapa de correlación

👉 Las visualizaciones ayudan a ver si el problema es separable con la representación actual.

---

## 5 Detección básica de outliers

Se usa z-score:

$$z = \frac{x - \mu}{\sigma}$$

Sirve para encontrar:

- posibles anomalías
- errores de medición

👉 No siempre se eliminan.

---

## 6 Estandarización

Se aplica:

$$x' = \frac{x - \mu}{\sigma}$$

Porque:

- muchos modelos son sensibles a escala
- mejora la geometría del espacio de X
- acelera el aprendizaje

👉 Es una transformación de representación:

$$\mathbf{X} \rightarrow \mathbf{X}' \quad \mathbf{X}' = \mathbf{X} \mathbf{W}$$

---

## 7 PCA

No es un clasificador.

Sirve para:

- encontrar direcciones de máxima varianza
- visualizar patrones dominantes

Formalmente:

$$\mathbf{Z} = \mathbf{X} \mathbf{W} \mathbf{Z} = \mathbf{X} \mathbf{W} \mathbf{W}^{-1} \mathbf{X} = \mathbf{X}$$

👉 Permite ver si existen estructuras internas en los datos.

---

## 8 Partición train / test

Se separan los datos en:

- entrenamiento
- prueba

Para evitar:

 evaluar el modelo con los mismos datos con los que aprendió.

---

## 9 Exportación

Se guardan los datasets limpios para:

- desacoplar el preprocesamiento
  - reutilizar en distintos modelos
- 



## Idea central del laboratorio

Este laboratorio no mejora un modelo.

Mejora la **calidad de la representación**.

---



## Relación directa con lo visto en clase

Todo el flujo es:

datos crudos  $\rightarrow \phi(X) \rightarrow X' \rightarrow f\theta(X')$   
|text{ datos crudos } \rightarrow \phi(X) \rightarrow X' \rightarrow f\theta(X')

Donde:

- la mayor ganancia suele venir de  **$\phi$  (preprocesamiento)**,
- no de cambiar el modelo.