

# 1. El Teorema del Límite Central (Central Limit Theorem)

## Definición:

El Teorema del Límite Central (TLC) establece que, si se extraen muestras independientes de una población (sin importar la distribución original), la distribución de las medias muestrales tiende a aproximarse a una distribución normal conforme el tamaño de la muestra aumenta. En otras palabras, aunque los datos originales no sean normales, la distribución de la media de muchas muestras lo será.

## Importancia:

Este teorema es fundamental en la inferencia estadística porque permite aplicar métodos paramétricos basados en la normalidad, como la construcción de intervalos de confianza y la realización de pruebas de hipótesis, aun cuando la población de origen tenga una distribución arbitraria. Su aplicabilidad en análisis muestrales lo convierte en una herramienta indispensable para evaluar incertidumbres y estimar parámetros.

## Aplicación Práctica:

En Machine Learning, el TLC se utiliza, por ejemplo, en la validación de modelos mediante el muestreo de datos para estimar métricas de error. Al analizar la media de los errores obtenidos en diversas particiones de los datos (como en el cross-validation), se puede asumir, bajo condiciones adecuadas, que la distribución de estos errores es aproximadamente normal, lo cual simplifica el análisis y la comparación de modelos.

# 2. Sampling (Muestreo)

## Definición:

El muestreo es el proceso de seleccionar un subconjunto representativo de una población con el fin de inferir características o comportamientos del conjunto total. Este procedimiento es esencial para obtener conclusiones generalizables sin tener que trabajar con la totalidad de los datos, lo que muchas veces resulta inviable por cuestiones de tiempo o recursos.

## Tipos de Muestreo:

- **Muestreo aleatorio simple:** Cada elemento de la población tiene la misma probabilidad de ser seleccionado, lo que minimiza el sesgo en la elección.
- **Muestreo estratificado:** La población se divide en subgrupos (estratos) homogéneos y se seleccionan muestras de cada estrato, garantizando que se reflejen adecuadamente las diferentes características.
- **Muestreo sistemático:** Se elige un punto de partida aleatorio y, a partir de ahí, se selecciona cada  $k$ -ésimo elemento, lo que puede simplificar la logística de la recolección.

## Relevancia en Machine Learning:

Un muestreo adecuado es crucial para entrenar y validar modelos de Machine Learning, ya que asegura que la variabilidad y la representatividad de los datos se mantengan. Esto

reduce el riesgo de sobreajuste y mejora la capacidad del modelo para generalizar a nuevos datos, permitiendo una evaluación más realista de su desempeño.

### 3. Diferencia entre Error Tipo I y Error Tipo II

#### **Error Tipo I (Falso Positivo):**

Ocurre cuando se rechaza incorrectamente la hipótesis nula siendo verdadera. Por ejemplo, en un sistema de detección de fraudes, clasificar una transacción legítima como fraudulenta generaría un falso positivo, afectando la experiencia del usuario.

#### **Error Tipo II (Falso Negativo):**

Se da cuando se acepta la hipótesis nula siendo falsa. Siguiendo el ejemplo anterior, no identificar una transacción realmente fraudulenta constituye un falso negativo, lo cual podría tener implicaciones financieras graves.

#### **Comparación:**

La elección de un umbral en pruebas estadísticas o en modelos predictivos implica un balance entre ambos tipos de error. Mientras que un falso positivo puede generar alarmas y costos operativos innecesarios, un falso negativo puede resultar en la omisión de eventos críticos o riesgos significativos. La decisión de priorizar la reducción de uno u otro error dependerá del contexto y de las consecuencias asociadas.

### 4. Regresión Lineal y sus Métricas

#### **Definición de Regresión Lineal:**

La regresión lineal es un método estadístico que modela la relación lineal entre una variable dependiente y una o más variables independientes. En el contexto del Machine Learning, se utiliza para realizar predicciones y entender cómo cada predictor influye en el resultado.

#### **Métricas de Evaluación:**

- **p-value:**  
Es una medida que indica la probabilidad de obtener resultados al menos tan extremos como los observados, asumiendo que la hipótesis nula es cierta. Un p-value bajo (generalmente menor a 0.05) sugiere que la relación entre la variable independiente y la dependiente es estadísticamente significativa.
- **Coeficientes (Pendientes e Intercepto):**  
Los coeficientes determinan el impacto de cada variable independiente sobre la variable dependiente. La pendiente indica el cambio esperado en la variable dependiente por cada unidad de cambio en la independiente, mientras que el intercepto representa el valor de la dependiente cuando las independientes son cero.
- **R-squared ( $R^2$ ):**  
Este coeficiente mide la proporción de la variabilidad de la variable dependiente que es explicada por el modelo. Un valor cercano a 1 indica que el modelo explica bien los datos, mientras que valores bajos sugieren que existen otros factores no considerados.

**Ejemplo Práctico:**

Imaginemos un modelo que predice el precio de las viviendas en función de su tamaño. Aquí, el coeficiente asociado al tamaño nos indica cuánto se incrementa el precio por metro cuadrado adicional; un p-value bajo en este coeficiente confirmaría la relevancia del tamaño, y un  $R^2$  alto indicaría que el modelo captura gran parte de la variabilidad en los precios.

## 5. Estadística por Iteración

**Definición:**

La estadística por iteración involucra la aplicación repetida de métodos estadísticos para estimar parámetros o evaluar la estabilidad de las estimaciones. Técnicas como el bootstrapping consisten en remuestrear los datos originales para construir una distribución empírica del estimador.

**Aplicación en Machine Learning:**

En el desarrollo de modelos complejos, métodos iterativos permiten evaluar la robustez de los parámetros y estimar intervalos de confianza más realistas. Por ejemplo, al aplicar bootstrapping para estimar la varianza de los coeficientes en un modelo de regresión, se obtiene una visión más precisa de la incertidumbre asociada a cada estimación.

## 6. Selección de Bias (Sesgo de Selección)

**Definición:**

El sesgo de selección se produce cuando la muestra de datos no es representativa de la población objetivo, lo que puede derivar en conclusiones erróneas y en modelos que no generalizan correctamente.

**Impacto:**

Utilizar datos sesgados puede afectar gravemente la validez de un estudio o de un modelo predictivo. Por ejemplo, si un modelo de reconocimiento facial se entrena mayoritariamente con imágenes de un solo grupo demográfico, su desempeño en otros grupos será pobre, lo que podría tener consecuencias éticas y prácticas.

**Estrategias de Mitigación:**

- **Diseño de Muestreo Apropiado:** Emplear técnicas como el muestreo estratificado para garantizar que todas las subpoblaciones estén representadas adecuadamente.
- **Validación Cruzada y Recolección de Datos Diversa:** Implementar métodos de validación que permitan detectar y corregir sesgos, complementados con esfuerzos para recolectar datos de fuentes variadas y representativas.

## 7. Probabilidad Binomial

**Definición:**

La distribución binomial es una distribución discreta que describe el número de éxitos en un

conjunto de ensayos independientes, cada uno con dos posibles resultados (éxito o fracaso), y en los que la probabilidad de éxito es constante.

#### **Aplicación Práctica:**

Un ejemplo práctico en Machine Learning es la evaluación de un clasificador binario. Supongamos que queremos estimar la probabilidad de obtener un número determinado de clasificaciones correctas (éxitos) en un conjunto de pruebas. Utilizando la distribución binomial, podemos calcular la probabilidad de que el desempeño observado se deba al azar o refleje verdaderamente la capacidad predictiva del modelo.

## **8. Referencias**

Mendenhall, W., Beaver, R. J., & Beaver, B. M. (2010). Introducción a la probabilidad y estadística (13ª ed.). CENGAGE Learning.

Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (2012). Probabilidad y estadística para ingenierías y ciencias (9ª ed.). Pearson.