



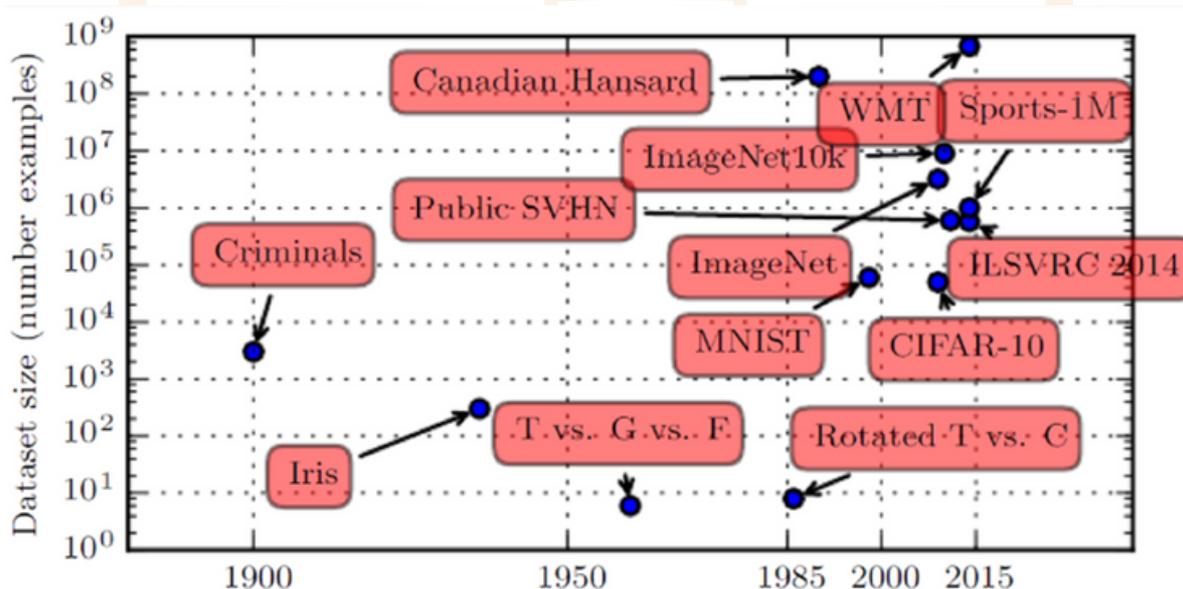
# Aumento en la disponibilidad de datos

# Aumento en la disponibilidad de datos

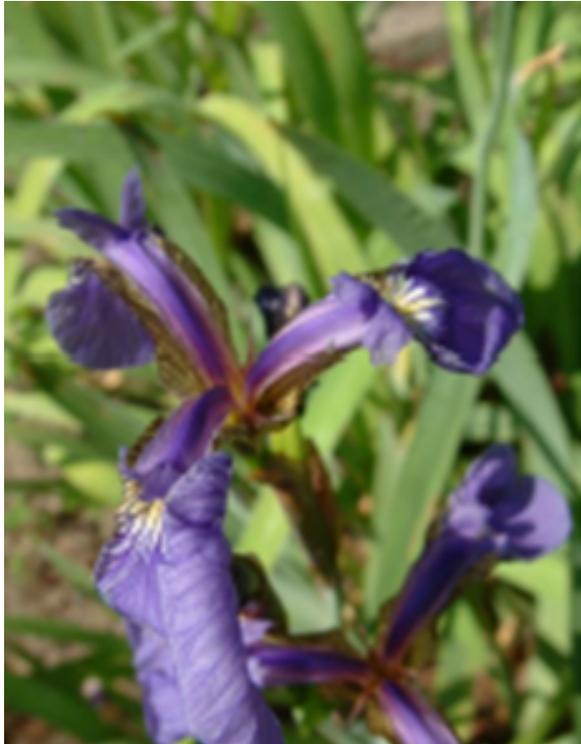
Es una tendencia impulsada por la digitalización de la sociedad. La era del “BigData” ha facilitado el ML, ya que la estadística de estimación generaliza bien los nuevos datos después de observar una pequeña cantidad de datos de entrenamiento.

A partir de 2016, la regla de oro en tareas de clasificación es que un algoritmo de IA supervisado, logrará rendimiento aceptable con 5,000 ejemplos etiquetados por categoría.

Pero el mismo algoritmo de IA superará el rendimiento humano cuando tenga al menos 10 Millones de ejemplos etiquetados. Es importante tener en cuenta los múltiples ejemplos de casos en los que una gran cantidad de datos ha impulsado avances significativos.



# 1. Base de datos Iris (Ronald Fisher, 1936)



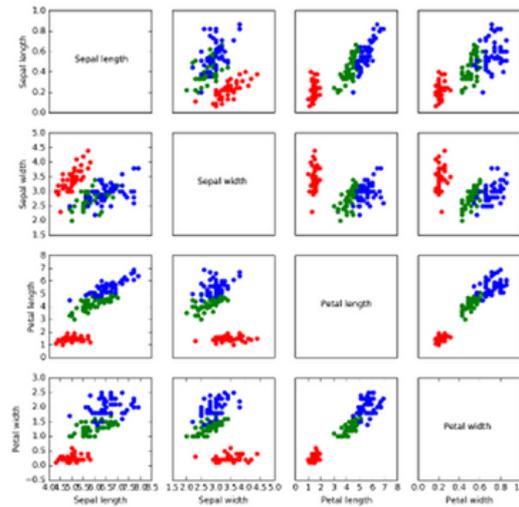
La base de datos Iris, creada por el estadístico y biólogo británico Sir Ronald A. Fisher en 1936, es un conjunto de datos clásico en el ámbito de la estadística y el aprendizaje automático. Fisher utilizó esta base de datos para ilustrar métodos de discriminación lineal.

La base de datos Iris consiste en 150 muestras de iris de tres especies diferentes: Iris setosa, Iris virginica e Iris versicolor. Para cada especie, se midieron cuatro características diferentes en centímetros: longitud del sépalo, ancho del sépalo, longitud del pétalo y ancho del pétalo.



El objetivo de la base de datos Iris es clasificar correctamente las flores de iris en una de las tres especies basándose en estas características. Es un problema de clasificación multiclase donde se pueden aplicar diversas técnicas de aprendizaje supervisado.

La base de datos Iris es muy conocida y se utiliza comúnmente como un conjunto de datos de prueba para evaluar algoritmos y modelos de clasificación. Su simplicidad y clara estructura lo convierten en un excelente punto de partida para aquellos que están aprendiendo conceptos fundamentales de clasificación y análisis de datos.



Iris setosa en Rojo, Iris virginica en Azul e Iris versicolor en Verde

## 2. Base de Datos MNIST

La base de datos Iris es muy conocida y se utiliza comúnmente como un conjunto de datos de prueba para evaluar algoritmos y modelos de clasificación. Su simplicidad y clara estructura lo convierten en un excelente punto de partida para aquellos que están aprendiendo conceptos fundamentales de clasificación y análisis de datos.



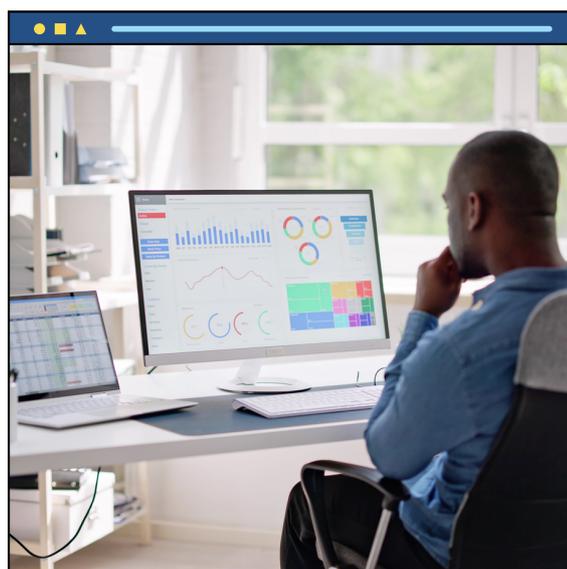
La base de datos MNIST es un conjunto de datos muy conocido y ampliamente utilizado en el campo de la visión por computadora y el aprendizaje automático. MNIST se compone de un conjunto de imágenes de dígitos escritos a mano, del 0 al 9. Estas imágenes son en escala de grises y tienen una resolución de 28x28 píxeles, lo que significa que cada imagen está representada por una matriz de 28 filas y 28 columnas.

En total, la base de datos MNIST contiene 60,000 imágenes de entrenamiento y 10,000 imágenes de prueba. Cada imagen está etiquetada con el dígito que representa, lo que convierte a MNIST en un conjunto de datos etiquetado.

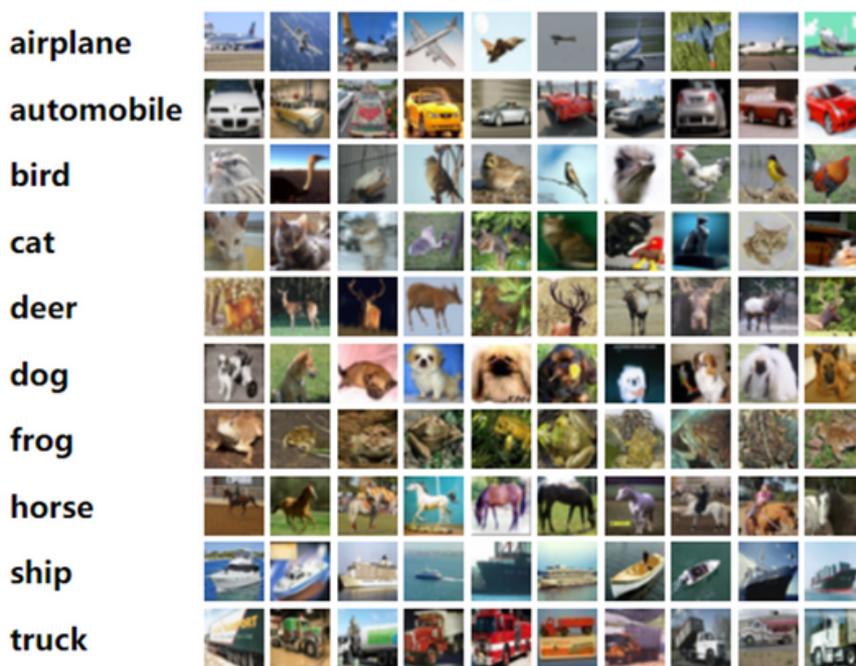


El propósito principal de MNIST es servir como un conjunto de datos estándar para evaluar y comparar algoritmos de clasificación, especialmente en el contexto de reconocimiento de dígitos escritos a mano. Muchos principiantes en el aprendizaje automático comienzan con MNIST debido a su simplicidad y facilidad de comprensión.

A lo largo de los años, MNIST ha sido un punto de referencia común para probar y demostrar la efectividad de diversos enfoques y modelos, incluidas redes neuronales convolucionales. Sin embargo, dado su carácter relativamente simple, algunos investigadores han abogado por la transición a conjuntos de datos más desafiantes para abordar mejor la complejidad de las tareas de visión por computadora.



### 3. Base de datos CIFAR-10



La base de datos CIFAR-10 es otro conjunto de datos ampliamente utilizado en el campo de la visión por computadora y el aprendizaje automático. CIFAR-10 consiste en 60,000 imágenes en color de 32x32 píxeles, divididas en 10 clases distintas. Cada clase representa un objeto o escena diferente.

Las clases en CIFAR-10 son las siguientes:



Cada clase tiene 6,000 imágenes, con 5,000 destinadas a entrenamiento y 1,000 para pruebas. Las imágenes de CIFAR-10 presentan desafíos adicionales en comparación con MNIST debido a su naturaleza en color y la mayor complejidad de los objetos representados.

CIFAR-10 se utiliza comúnmente como un conjunto de datos de referencia para tareas de clasificación de imágenes más complejas y realistas en comparación con MNIST.



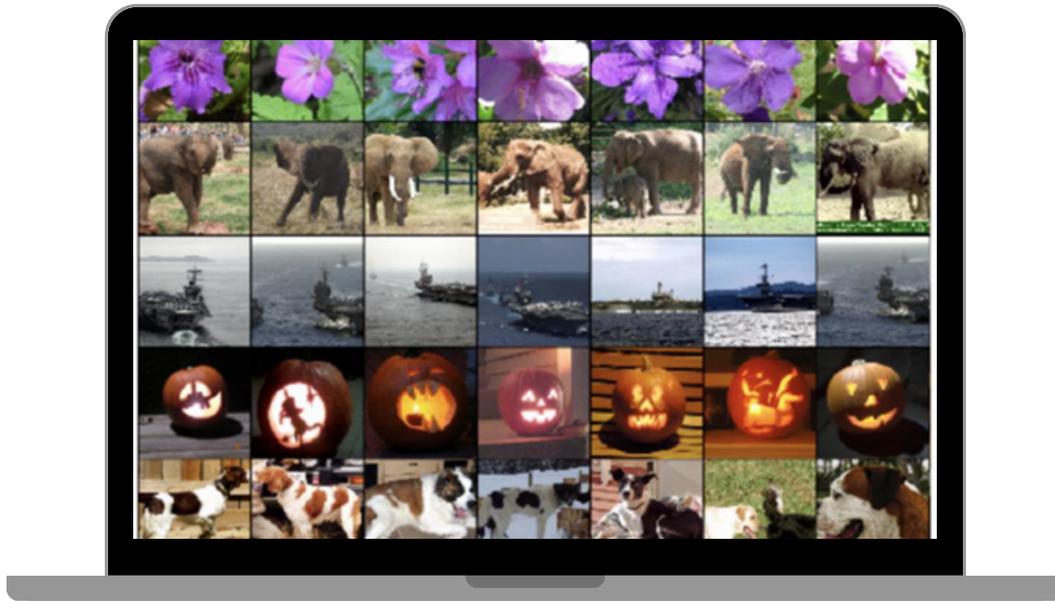
Al igual que MNIST, CIFAR-10 ha sido fundamental para evaluar y comparar diversos algoritmos y modelos en el contexto de la visión por computadora. Su uso sigue siendo prominente en la investigación y la educación en aprendizaje profundo y visión por computadora.

## 4. ImageNet

Es una base de datos de imágenes utilizada para entrenar y evaluar algoritmos de reconocimiento de objetos en imágenes.

Fue creada por el equipo de investigación liderado por Fei-Fei Li en la Universidad de Stanford y lanzada en el año 2009. ImageNet contiene millones de imágenes etiquetadas manualmente en más de 20,000 categorías distintas, abarcando una amplia variedad de objetos, animales y escenas.





El conjunto de datos ImageNet Challenge, asociado con ImageNet, ha sido un componente clave para evaluar y comparar el rendimiento de algoritmos de visión por computadora y aprendizaje profundo. En particular, la competición anual ImageNet Large Scale Visual Recognition Challenge (ILSVRC) se ha convertido en una referencia importante en la comunidad de investigación de inteligencia artificial.

La tarea principal en el ILSVRC es clasificar correctamente las imágenes en las categorías específicas de la base de datos ImageNet. El conjunto de datos ha sido crucial para avanzar en la capacidad de las computadoras para **comprender y reconocer objetos en imágenes**, especialmente a través del desarrollo de arquitecturas de redes neuronales profundas como las **redes neuronales convolucionales (CNN)**, que han demostrado ser altamente efectivas en tareas de clasificación de imágenes.

