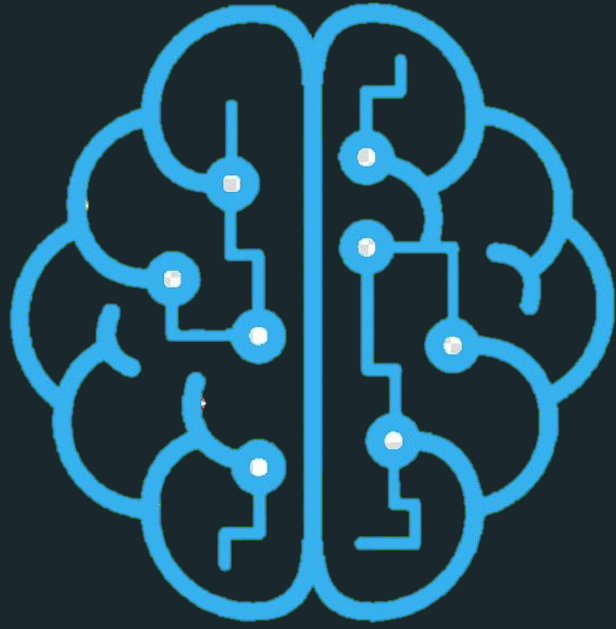# APRENDIZAJE DE MÁQUINA

Preprocesamiento

ISTA

# Preprocesamiento

**Introducción**, **Data cleaning:** Valores faltantes y duplicados, **Data transformation:** Codificación de los datos  Escalamiento de los datos Normalización de los datos, **Data reduction:** Selección de características Extracción de característica, **Conclusiones**
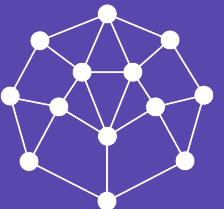
# #1. Introducción

# Datooos!!! Muchos datos!!

- Recurso más valioso hoy en el mundo actual.

- Según el Foro Económico Mundial, para el 2025 estaremos generando alrededor de **463 exabytes de datos a nivel mundial por día**.

- Pero:

¿Todos estos **datos son lo suficientemente adecuados** para ser utilizados por algoritmos de aprendizaje automático?

WORLD ECONOMIC FORUM

Saturdays.AI

# Hablando de datos pensamos …

- Grandes conjuntos de datos con una gran cantidad de filas y columnas.

- No siempre es el caso: los **datos pueden estar en muchas formas diferentes:** tablas estructuradas, imágenes, archivos de audio, videos, etc.

- Las **máquinas no entienden** el texto libre, las imágenes o los datos de video tal como están, **entienden los 1 y 0.**
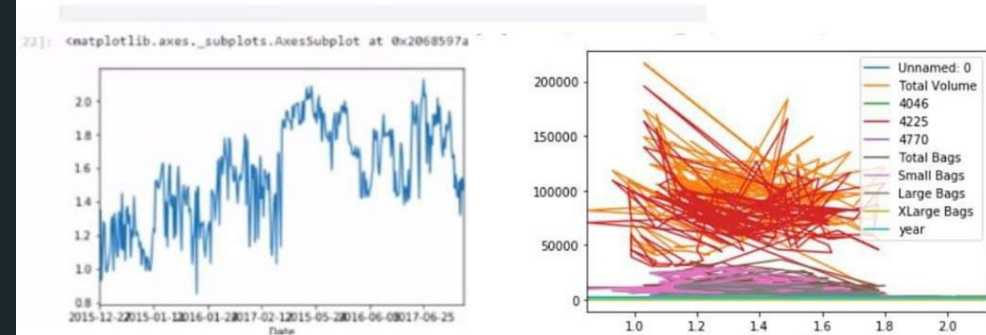
# Tenemos un problema

- Los datos del mundo real a menudo son:
  - incompletos,
  - inconsistentes,
  - contienen muchos errores.



**La calidad de los datos afecta directamente la capacidad de nuestro modelo para aprender.**

# Y una necesidad

- Contar con datos que sean
  - completos,
  - precisos,
  - válidos



Las características universalmente reconocidas de los datos de buena calidad, que incluyen: accesibilidad, exactitud, exhaustividad, consistenci (o coherencia), definición a vigencia, , institucional e interpretabilidad (Teslow, 2016). granularidad, precisión, relevancia, entorno oport

# Data preprocessing

Es el proceso de **aplicar transformaciones a los datos** para llevarlos a un estado que la máquina pueda **analizarlos fácilmente.**



**Paso muuuy importante en todo proceso de Machine Learning!!!**

Su objetivo es **disponer de datos de calidad** previo al modelado utilizando algoritmos.

# Data preprocessing

Alrededor del **60% del tiempo** de los científicos de datos es empleado aquí, **preparando los datos para el modelado.**





**WHAT DATA SCIENTISTS SPEND THE MOST TIME DOING**

- 60% Cleaning and organising data
- 19% Collecting data sets
- 3% Building training sets
- 4% Refining algorithms
- 5% Other
- 9% Mining data for patterns

Source: CrowdFlower 2016

Fuente: Forbes (2016).

**Paso que requiere mucho tiempo!!!**

# Definiciones

## Dataset:

- Colección de datos.

  - **Ejemplos:** registros de interacciones, eventos, observaciones.

- Descritos mediante una serie de características o features.

  - **Ejemplos**: la masa de un objeto físico o el momento en que ocurrió un evento, etc.

**Feature vector**

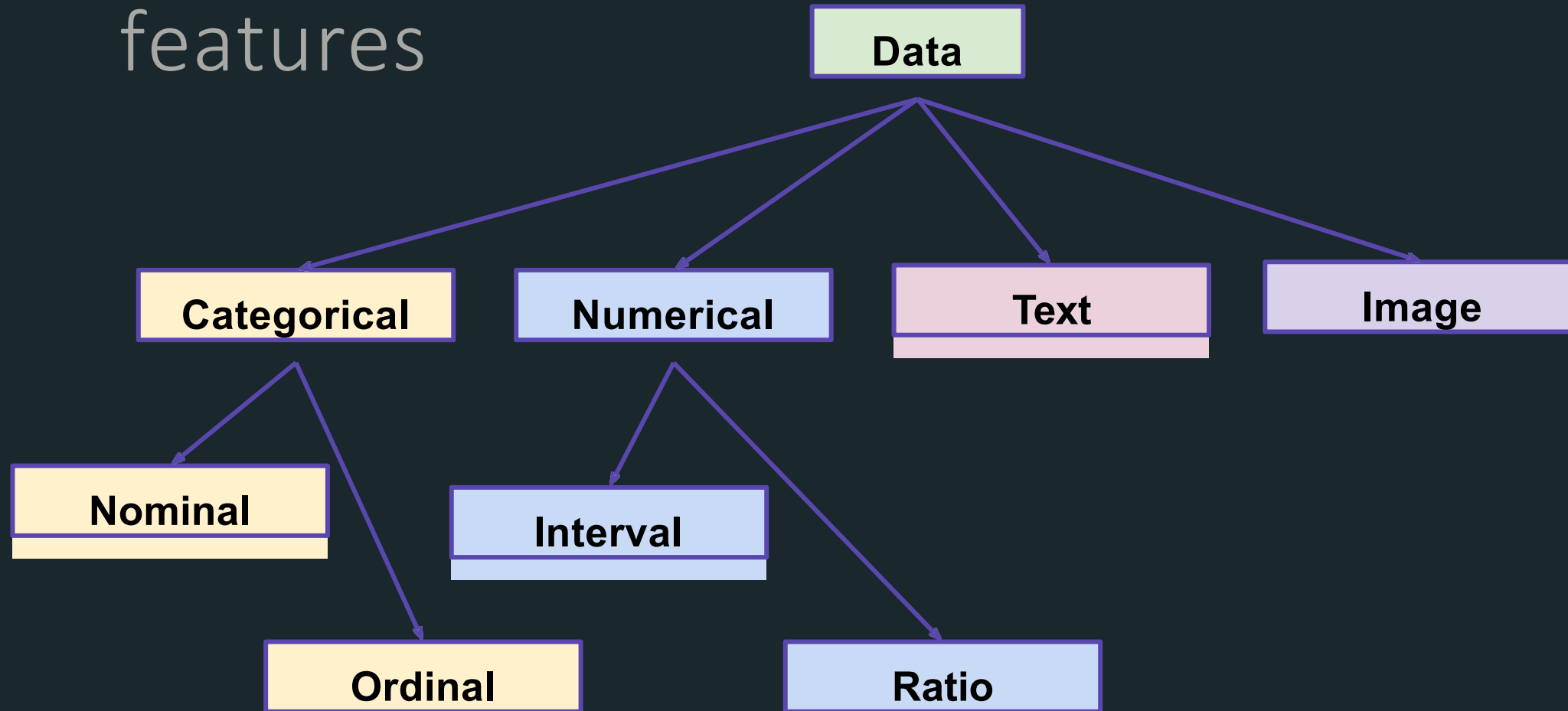| | Features | | | | Label |
|---|---|---|---|---|---|
| **Position** | **Experience** | **Skill** | **Country** | **City** | **Salary ($)** |
| Developer | 0 | 1 | USA | New York | 103100 |
| Developer | 1 | 1 | USA | New York | 104900 |
| Developer | 2 | 1 | USA | New York | 106800 |
| Developer | 3 | 1 | USA | New York | 108700 |
| Developer | 4 | 1 | USA | New York | 110400 |
| Developer | 5 | 1 | USA | New York | 112300 |
| Developer | 6 | 1 | USA | New York | 114200 |
| Developer | 7 | 1 | USA | New York | 116100 |
| Developer | 8 | 1 | USA | New York | 117800 |
| Developer | 9 | 1 | USA | New York | 119700 |
| Developer | 10 | 1 | USA | New York | 121600 |

# Definiciones: Tipos de features

# Tipos de features : Categorical

- Características cuyos valores se toman de un **conjunto definido de valores.**
  - Ejemplos?…. En el chat :D

| Nominal | Ordinal |
|---|---|
| - Variables categóricas **sin un orden implícito.** | - Variables categóricas **con un orden natural implícito.** |
| - **Ejemplo:** Los colores de un carro: negro, morado, rosa. | - **Ejemplo:** Los tamaños de la ropa: chico, mediano, grande. |

# Tipos de features : Numerical

- Características representadas por **números cuyos valores son continuos o discretos.**
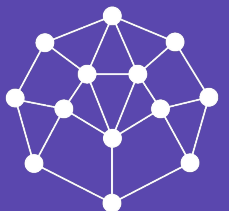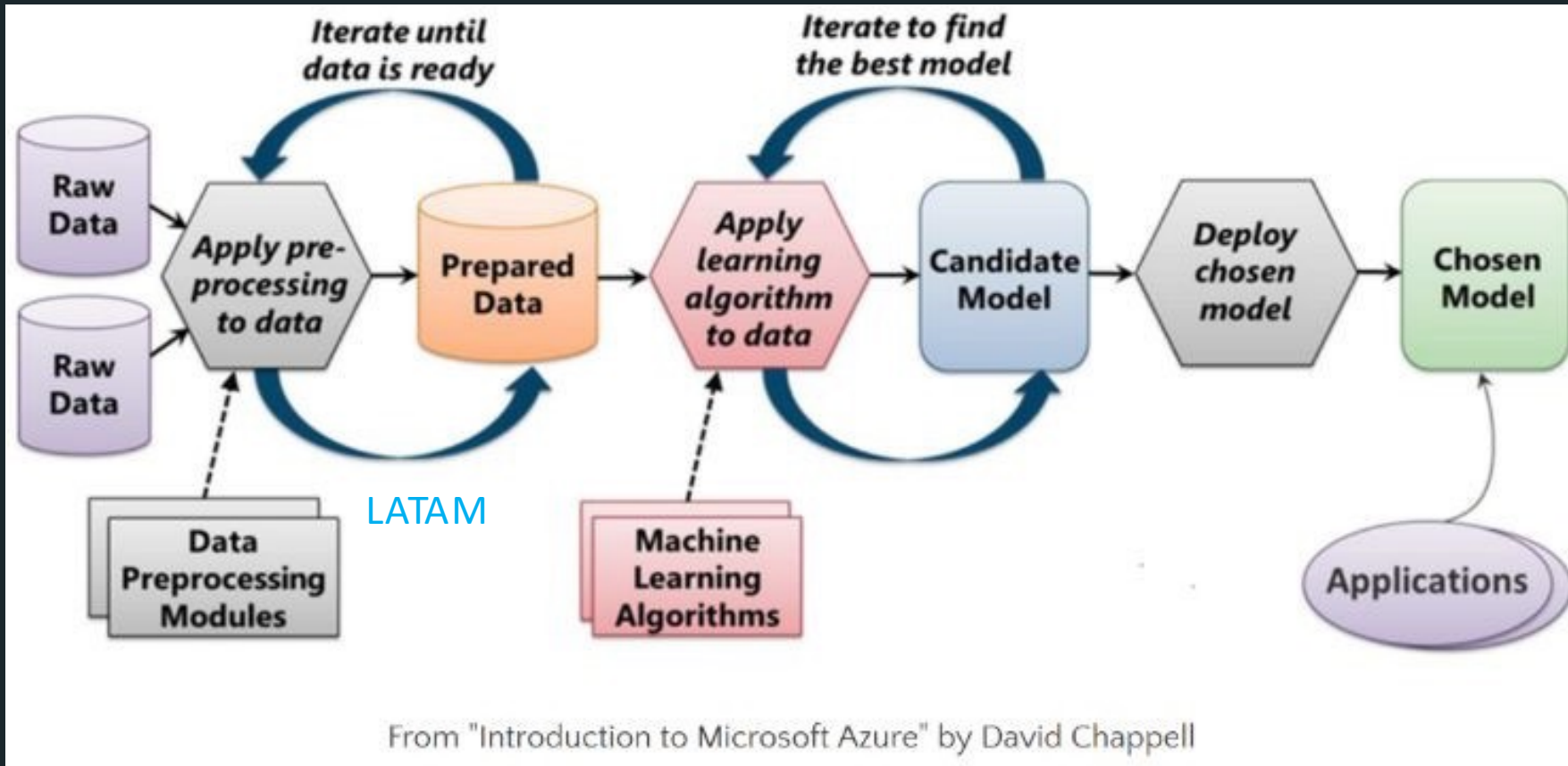  - Ejemplos?…. En el chat :D

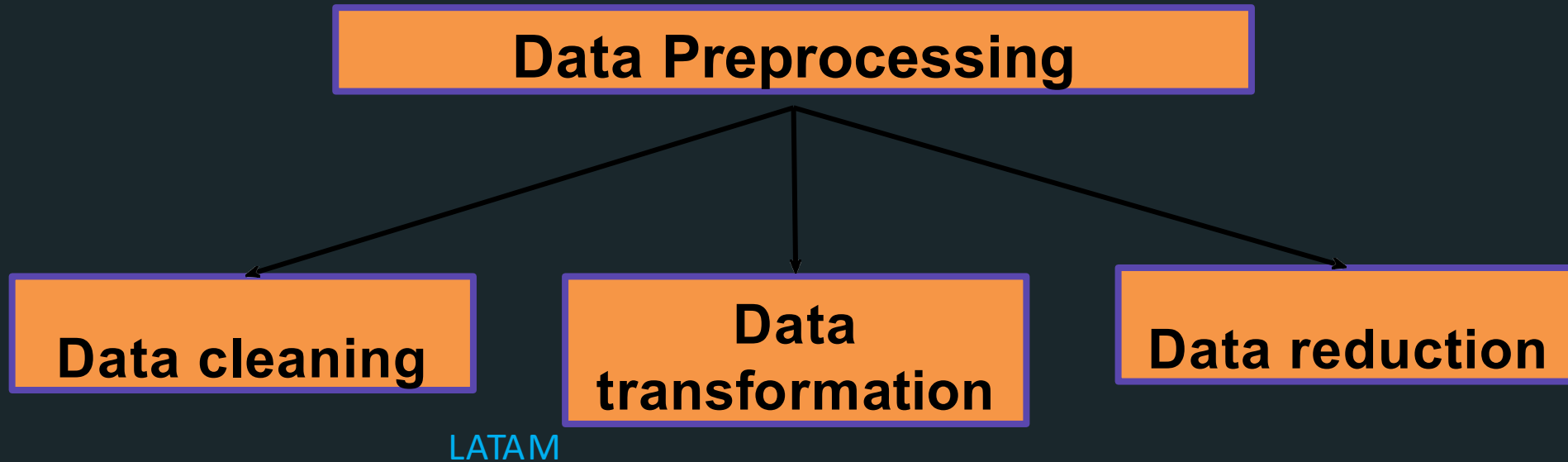| Interval | Ratio |
|----------|-------|
| - Con una unidad de medida definida.<br>- Representa valores como 0 y menores que 0.<br><br>- **Ejemplo:** temperatura en Celsius | - Con una unidad de medida definida.<br>- Representa valores de 0 y mayores a 0.<br><br>- **Ejemplo:** estatura y peso. |

# Metodología de Ciencia de Datos.



From "Introduction to Microsoft Azure" by David Chappell

# Data Preprocesing

# CONCLUSIONES!!!