

A young girl with brown hair, wearing a headband and boxing gloves, is in a boxing stance with her arms raised. She is wearing a dark tank top with a large white letter 'A' on it. The background is a solid blue color.

#3 Data preprocessing

by Saturdays.AI

Saturdays.AI LATAM



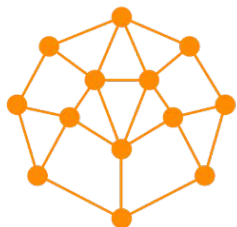
Objetivo

Aprender diferentes técnicas para el preprocesado de los datos



Metodología de Ciencia de Datos.

Week 3



Saturdays.AI
LATAM

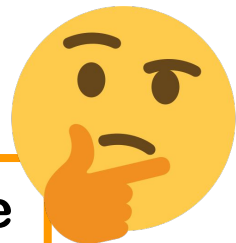
Agenda

- Introducción
- Data cleaning
 - Valores faltantes
 - Valores duplicados
- Data transformation
 - Codificación de los datos
 - Escalamiento de los datos
 - Normalización de los datos
- Data reduction
 - Selección de características
 - Extracción de características

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?



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**.

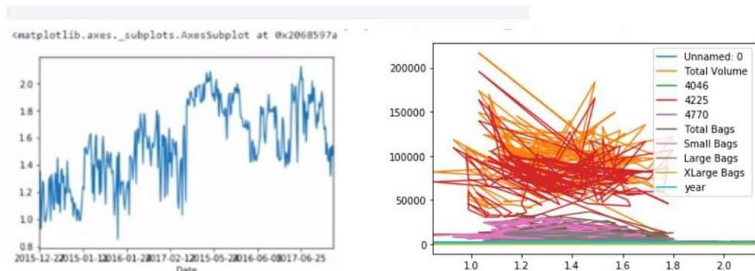
[illegible]

Saturdays.AI
LATAM

Tenemos un problema

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

me everytime:



Expectation

Reality

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



Saturdays.AI
LATAM

Data preprocessing

- Paso muuuy importante en todo proceso de Machine Learning
 - Aplicar transformaciones a los datos para llevarlos a un estado que la máquina ahora puede analizar fácilmente.
- **Requiere muuuucho 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..

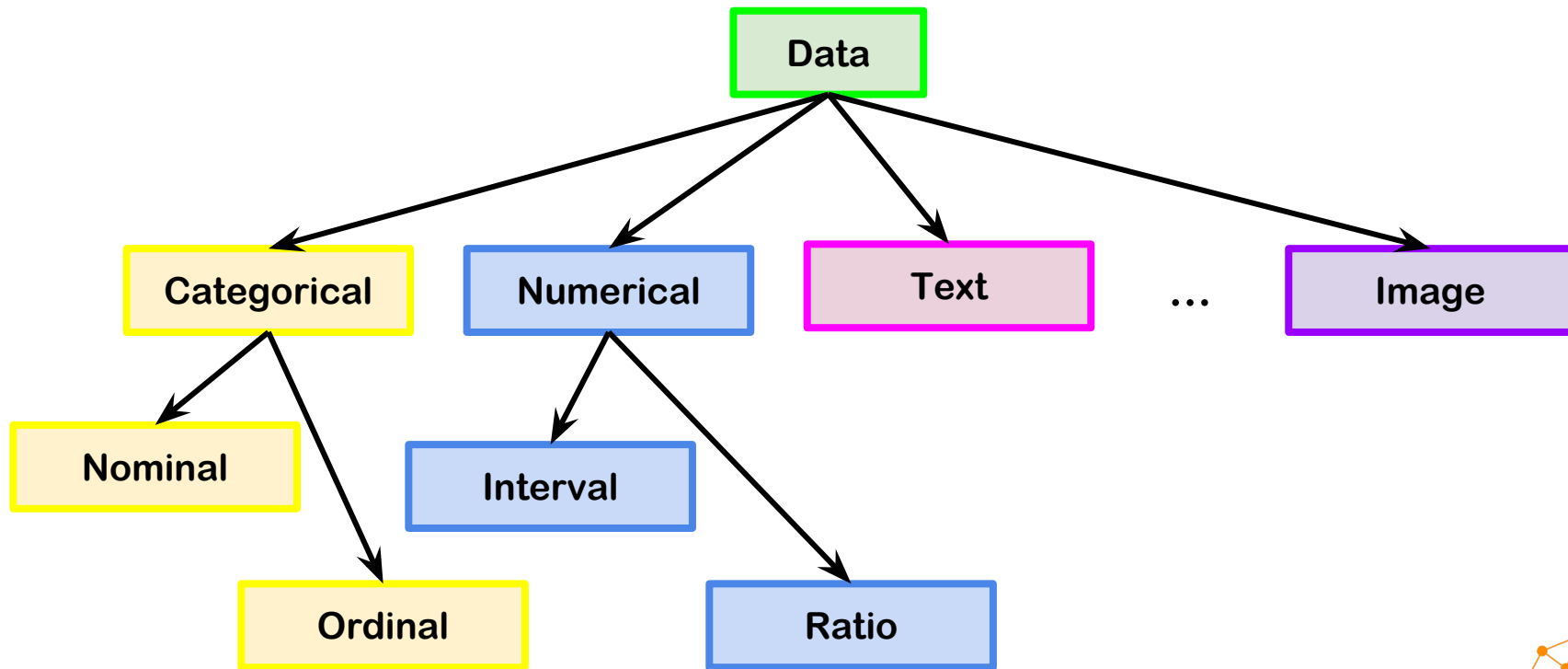
Feature vector

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



Saturdays.AI
LATAM

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

- Variables categóricas **sin un orden implícito**
- **Ejemplo:** Los colores de un carro: negro, morado, rosa

Ordinal

- Variables categóricas **con un orden natural implícito**
- **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

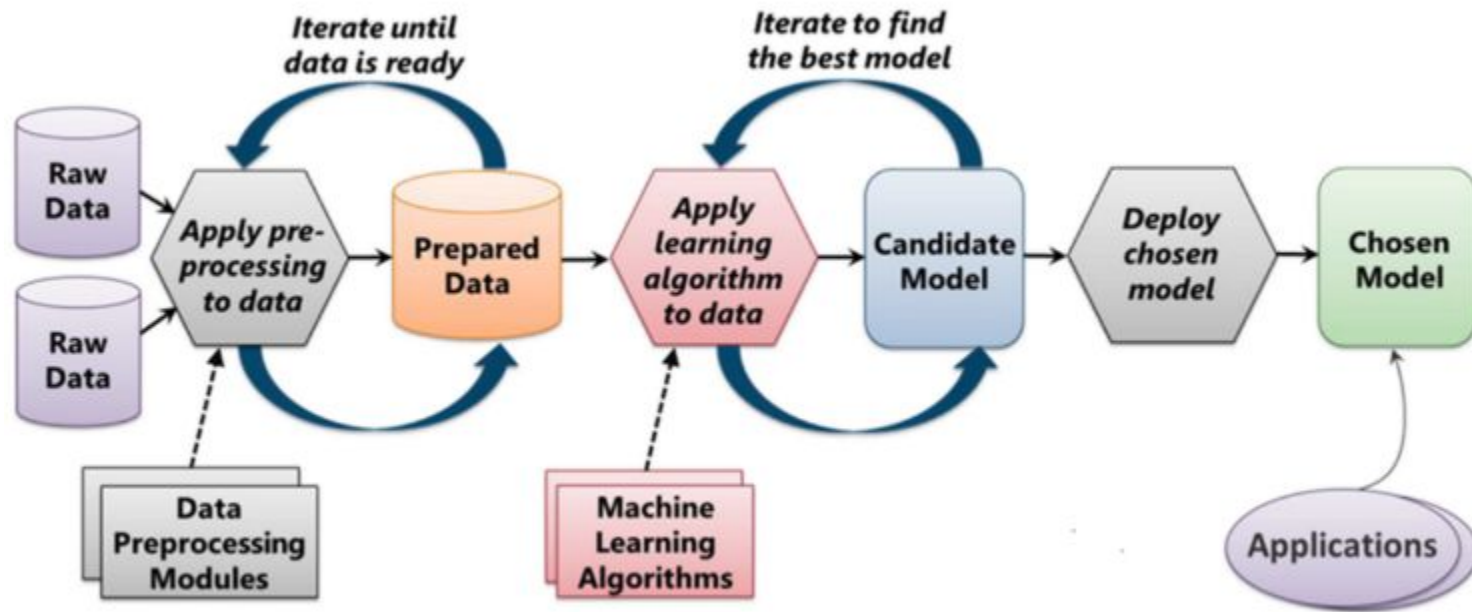
- Con una unidad de medida definida
- Representa valores como 0 y menores que 0
- **Ejemplo:** temperatura en Celsius

Ratio

- Con una unidad de medida definida
- Representa valores de 0 y mayores a 0
- **Ejemplo:** estatura y peso



Metodología de Ciencia de Datos.

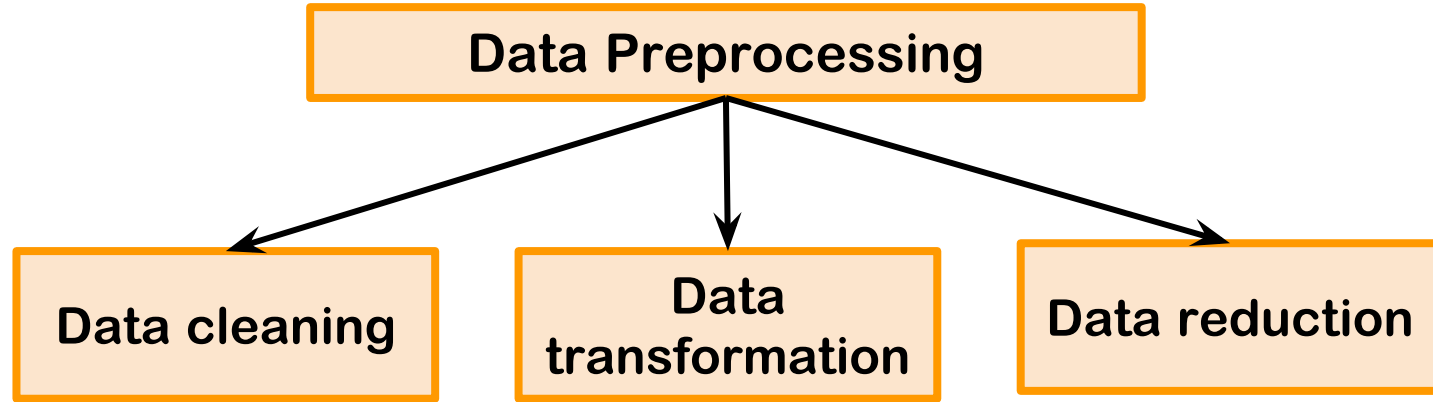


From "Introduction to Microsoft Azure" by David Chappell



Saturdays.AI
LATAM

Data Preprocessing





#3.1 Data cleaning

Data cleaning

- Los datos pueden faltar durante la extracción o recolección de datos.
- Reducen la calidad de cualquiera de nuestras métricas de rendimiento.
- **“Data cleaning”:** Técnica que implica el manejo de datos faltantes, datos ruidosos, etc.

Los valores faltantes pueden aparecer como un signo de interrogación (?) O un cero (0) o menos uno (-1) o un espacio en blanco.

EDA



Saturdays.AI
LATAM

Data cleaning - Missing values

Ignorarlos: Solo cuando el conjunto de datos que tenemos es bastante grande.

Eliminar filas con datos faltantes:

- Estrategia simple y a veces efectiva.
- Falla si muchos objetos tienen valores faltantes.
- Si una *característica tiene valores perdidos en su mayoría, entonces esa característica en sí también se puede eliminar* (+ 75%).

Estimar valores perdidos:

- Si solo falta un porcentaje razonable de valores (20%)
- Métodos más comunes:
 - Manualmente.
 - Promedio.
 - El valor más probable de la característica respectiva.
 - ML algoritmo.



Saturdays.AI
LATAM

¿Que usarias?

Ejemplo 1

ID	Children	Age of youngest child	Did you drink Coca-Cola in the last 24 hours?	How many colas did you drink in the past 24 hours?
1	No		No	
2	Yes	18	Yes	2
3	No		No	
4	Yes	13	No	
5	Yes	8	Yes	1

Ejemplo 2

ID	Gender	Age	Income
1	Male	Under 30	Low
2	Female	Under 30	Low
3	Female	30 or more	High
4	Female	30 or more	
5	Female	30 or more	High



Data cleaning - Duplicate values



- Puede suceder cuando la misma persona envía un formulario más de una vez.
- En la mayoría de los casos, los duplicados **se eliminan** para no dar a ese objeto de datos en particular una ventaja o sesgo.

Time to Code!!! :D

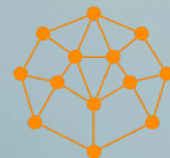


Saturdays.AI
LATAM



#3.3 Data transformation

Feature Encoding



Saturdays.AI
LATAM

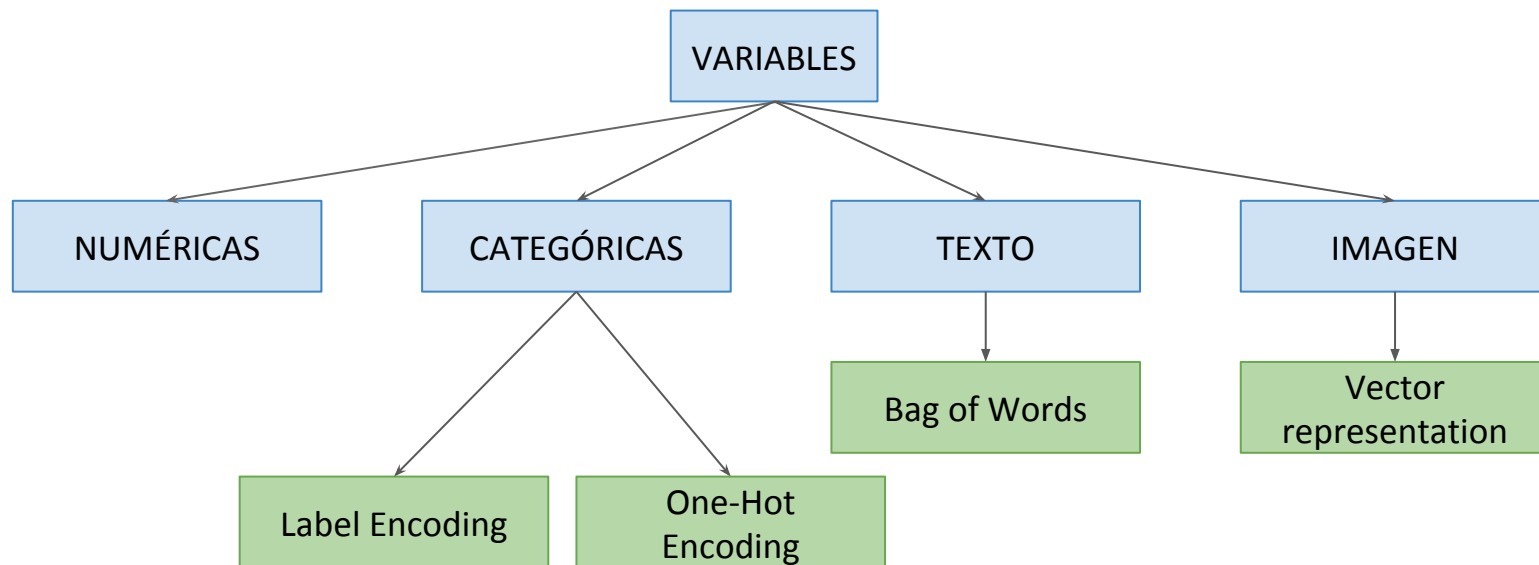
Data transformation

- Proceso en el que toma datos de su estado de origen sin procesar, y los transforma en datos listos para el análisis.

Data transformation: Feature encoding

- Consiste en realizar transformaciones en los datos para aceptarse como entrada para algoritmos de aprendizaje automático, manteniendo su significado original.
- Los algoritmos de aprendizaje automático sólo puede aceptar valores numéricos.

Tipos de características



Variables categóricas

- Una variable **categórica** es una variable que puede tomar uno de un número **limitado** de posibles **valores o categorías**.

ID	TIPO
1	COMPRA
2	VENTA
3	VENTA

Variables categóricas - Label Encoding

- Asigna a cada categoría un valor numérico distinto.
- Útil cuando las categorías **tienen relación**.
 - Ejemplo nivel socioeconómico: bajo = 1, medio = 2, alto = 3.
- **Ligero computacionalmente**, genera solamente una columna.
- Si las categorías **no tiene relación** puede generar sesgos a los algoritmos de ML.

ID	TIPO	TIPO
1	COMPRA	1
2	VENTA	2
3	VENTA	2

Variables categóricas - One-Hot Encoding

- Asigna a cada categoría un vector que contiene 1 y 0 que denota la presencia o ausencia de la característica.
- Útil cuando las categorías **no tienen relación** mejorando la precisión de los algoritmos.
- Puede generar **demasiadas columnas**.

ID	TIPO	COMPRA	VENTA
1	COMPRA	1	0
2	VENTA	0	1
3	VENTA	0	1

Variables de texto

PROD_DESC
Unspecified vomiting of pregnancy, unspecified as to episode of care or not applicable
Retained (old) foreign body following penetrating wound of orbit.
Subjective visual disturbance, unspecified.

Variables de texto - Bag of Words

- **Bag-of-words** es una representación numérica de variables de texto que describe la aparición de palabras dentro de un documento.
- Este utiliza:
 - **Vocabulario** de palabras conocidas.
 - Una **métrica** que evalúa la **presencia** de palabras del vocabulario.



Variables de texto - Bag of Words

Dataset

ID	PROD_DESC
1	Algodón, poliéster, sin cordones. Preferible lavar a máquina.
2	Solo algodón, con cordones. No lavar a máquina.
3	Playera azul. No lavar a máquina.

Vocabulario

1. Algodón
2. Poliéster
3. Cordones
4. Preferible
5. Lavar
6. Máquina.
7. Playera
8. Azul

Representación Bag of Words

ID	Alg odón	Poli éster	Cor don es	Pre feri ble	Lav ar	Má qui na	Pla yer a	Azu l
1	1	1	1	1	1	1	0	0
2	1	0	1	0	1	1	0	0
3	0	0	0	0	1	1	1	1



Imágenes - Vector representation



```
0 2 15 0 0 0 11 10 0 0 0 0 0 0 9 9 0 0 0
0 0 0 4 60 157 236 255 255 255 177 95 61 32 0 0 29
0 10 15 115 238 255 244 245 243 250 249 255 222 103 10 0
0 14 170 255 219 244 254 255 253 245 255 249 253 251 124 1
2 98 255 228 255 251 254 211 141 116 122 215 251 218 255 49
13 217 243 255 192 33 120 52 2 0 10 13 232 255 215 36
14 229 252 254 49 12 0 0 7 7 0 70 237 252 235 62
6 141 245 255 212 25 11 9 3 0 115 236 243 255 137 0
0 87 252 250 248 216 60 0 1 177 252 255 249 14 5 0
0 13 113 255 255 245 255 182 181 248 252 242 256 36 0 19
1 0 5 117 251 255 241 255 247 255 241 162 17 0 7 0
0 0 0 4 60 157 236 255 244 254 253 250 11 0 1 0
0 0 4 60 157 236 255 248 252 252 244 250 10 0 4
0 22 206 252 246 251 241 100 240 11 255 245 250 246 9 0
0 111 255 242 255 192 24 0 0 6 39 255 232 230 56 0
0 128 251 255 11 7 11 0 0 0 2 45 255 250 100 0
0 113 255 255 10 4 20 0 15 3 11 182 251 245 61 0
0 89 251 241 255 235 16 15 10 231 248 255 252 62 4
0 164 115 255 247 255 255 255 249 249 249 245 139 0 5
0 0 75 110 255 255 250 248 255 255 248 245 134 12 0
0 0 0 6 1 0 10 115 233 255 255 255 37 0 0 4 1
0 0 5 7 0 0 0 0 0 0 14 1 0 6 6 0 0
```



```
0 2 15 0 0 11 10 0 0 0 0 0 9 9 0 0 0
0 0 0 4 60 157 236 255 255 255 177 95 61 32 0 0 29
0 10 15 115 238 255 244 245 243 250 249 255 222 103 10 0
0 14 170 255 219 244 254 255 253 245 255 249 253 251 124 1
2 98 255 228 255 251 254 211 141 116 122 215 251 218 255 49
13 217 243 255 195 33 120 52 2 0 10 13 232 255 215 36
16 229 252 254 49 12 0 0 7 7 0 70 237 252 235 62
6 141 245 255 212 25 11 9 3 0 115 236 243 255 137 0
0 87 252 250 248 215 60 0 1 121 252 255 248 14 6 0
0 13 113 255 255 245 255 182 181 248 252 242 256 36 0 19
1 0 5 117 251 255 241 255 247 255 241 162 17 0 7 0
0 0 0 4 60 157 236 255 244 254 253 250 11 0 1 0
0 0 4 60 157 236 255 248 252 252 244 250 10 0 4
0 22 206 252 246 251 241 100 240 11 255 245 250 246 9 0
0 111 255 242 255 192 24 0 0 6 39 255 232 230 56 0
0 128 251 255 11 7 11 0 0 0 2 45 255 250 100 0
0 113 255 255 10 4 20 0 15 3 11 182 251 245 61 0
0 89 251 241 255 235 16 15 10 231 248 255 252 62 4
0 164 115 255 247 255 255 255 249 249 249 245 139 0 5
0 0 75 110 255 255 250 248 255 255 248 245 134 12 0
0 0 0 6 1 0 10 115 233 255 255 255 37 0 0 4 1
0 0 5 7 0 0 0 0 0 0 14 1 0 6 6 0 0
```



```
[ 0 2 15 0 0 11 10 0 0 0 0 0 9 9 0 0 0 0 0 4 60 157 236 255 255 177 95 61 32 0 0 29 ...
... 0 0 6 1 0 52 153 233 255 252 147 37 0 0 4 1 0 0 5 5 0 0 0 0 14 1 0 6 6 0 0 ]
```

Fuente: [Link](#)

Time to Code!!! :D

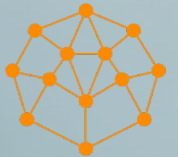


Saturdays.AI
LATAM

#3.2 Transformación de datos

Escalamiento / Feature scaling

Normalización de datos / and Feature normalization



Saturdays.AI
LATAM

Feature scaling vs Feature normalization

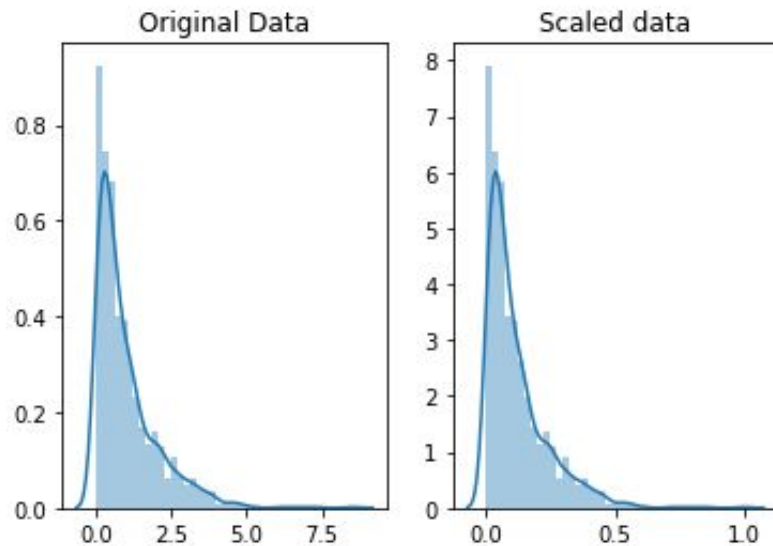
- Cambiar los valores de las columnas numéricas en el conjunto de datos a una escala común

Feature scaling, cambia el **rango** de sus datos

Feature normalization
cambiando la forma de la **distribución** de sus datos.

Feature scaling

- Transformación para que se **ajusten a un rango específico**.
 - Ejemplo: 0–100 o 0–1.
- Útil cuando se usa en métodos basados en medidas de qué tan separados están los puntos de datos



Observe que la forma de los datos no cambia, pero que en lugar de variar de 0 a 8, ahora varía de 0 a 1.

Feature scaling

Con Feature Scaling generalmente nos referimos a generar un cambio en el rango de los valores sin cambiar la forma de la distribución de los mismos, de esta forma es posible mejorar el performance de un algoritmo de machine learning o su tiempo de ejecución cuando sus features se encuentran en una escala similar o cercana a ser normalmente distribuida.

Algunos algoritmos que pueden mejorar su performance mediante estas técnicas son:

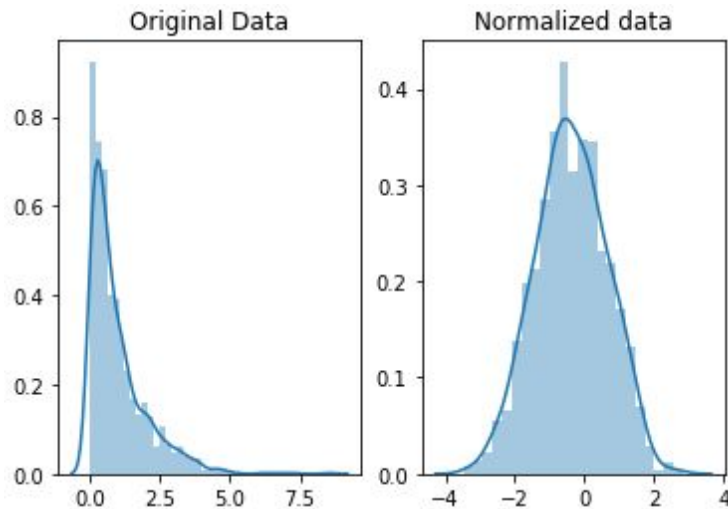
- * Regresión logística y lineal
- * Nearest Neighbors
- * Redes Neurales
- * Análisis de componentes principales

Feature normalization

- Pretendemos cambiar los valores de columnas numéricas en el dataset a una escala en común, sin distorsionar las diferencias en los rangos de valores. En machine learning no siempre los dataset requieren normalización, sólo es requerido cuando los features tienen rangos diferentes.
- La normalización es útil cuando se desconoce la distribución de los datos o cuando la distribución no es Gaussiana (en campana).

Feature normalization - Standardization

- **Z-Normalization (Standardization):**
 - Transformación para que se **ajusten a una distribución normal**.

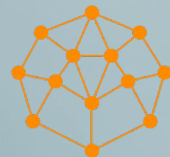


Time to Code!!! :D



Saturdays.AI
LATAM

#3.3 Data reduction



Saturdays.AI
LATAM

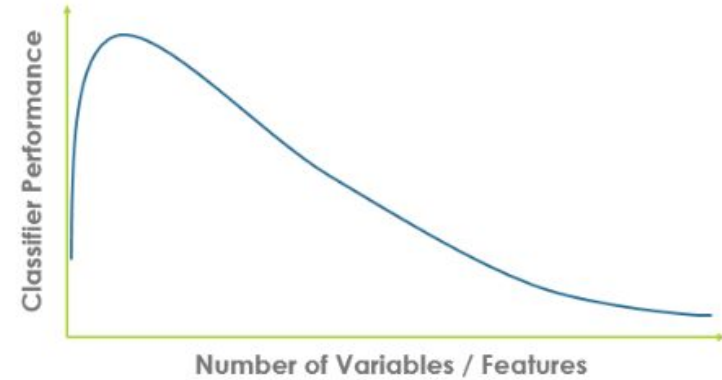
Data reduction

- Manejamos una gran cantidad de datos.
 - Problema:
 - **Gran volumen de datos =
Análisis más difícil.**
 - Solución: Técnicas de reducción de datos.
- Objetivo:
 - Reducir el almacenamiento de datos.
 - Eliminar características innecesarias.
 - Reducir los costos computacionales durante análisis.



The curse of dimensionality

- "Dimensionalidad" simplemente se refiere al **número de características** en su conjunto de datos.
- Cuando el **número de características es muy grande en relación con el número de observaciones** en su conjunto de datos, ciertos algoritmos luchan por entrenar modelos efectivos.



Data reduction - Feature selection

- Filtrar características irrelevantes o redundantes de su conjunto de datos.

Feature selection != Feature extraction

1. Feature selection **mantiene un subconjunto** de las características originales
2. Feature extraction **crea nuevas**.

Feature selection - Variance thresholds

- Eliminan características cuyos **valores no cambian mucho de una observación a otra** (es decir, su varianza cae por debajo de un umbral).
- Estas características proporcionan poco valor.
- Debido a que la **variación depende de la escala, siempre debes normalizar** sus características primero.

POBLACIÓN

$$\sigma^2 = \frac{\sum (x - \bar{x})^2}{N}$$

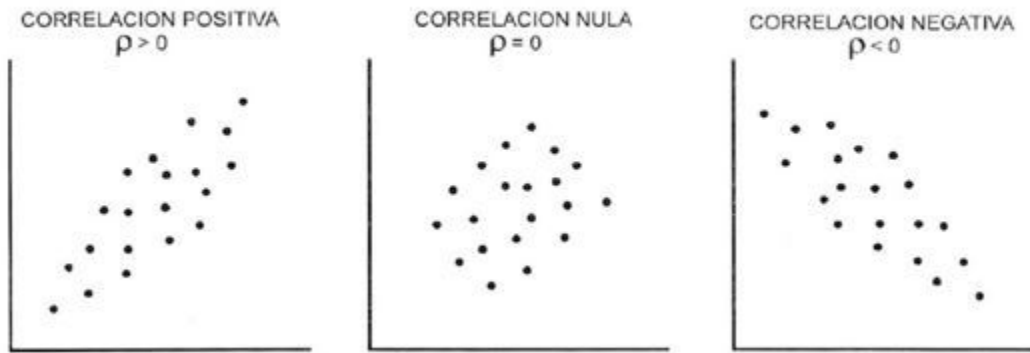
MUESTRA

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1}$$

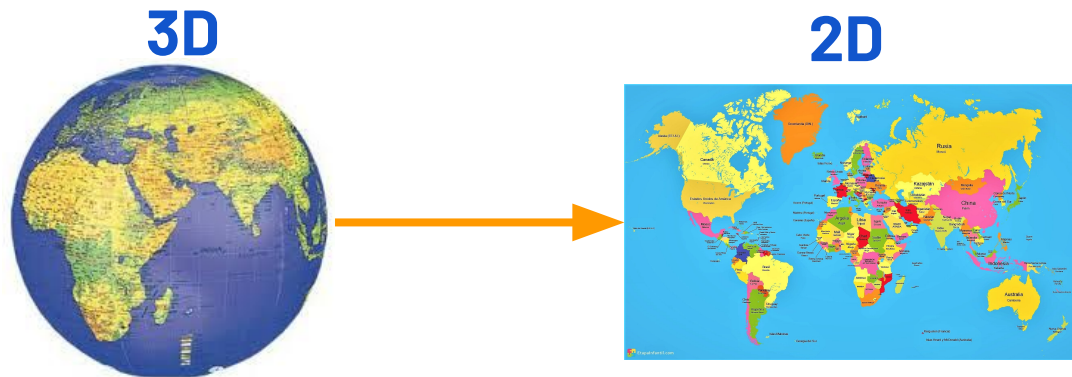


Feature selection - Correlation thresholds

- Eliminan características que están **altamente correlacionadas con otras** (es decir, sus valores cambian de manera muy similar a la de los demás).
- Estas características proporcionan información redundante.



Data reduction - Feature extraction



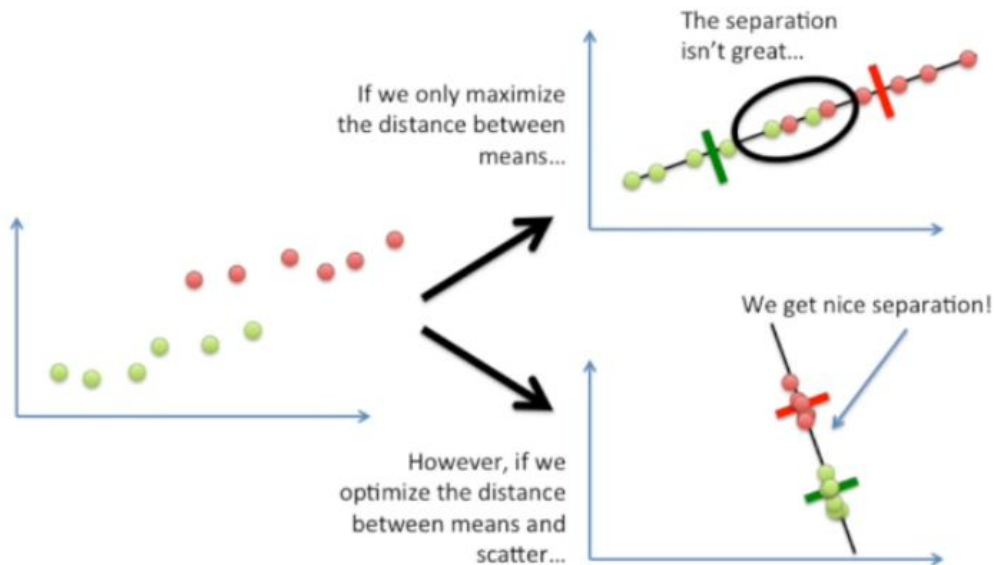
- Creación de un **nuevo conjunto de características más pequeñas** que aún captura la mayor parte de la información útil.

Feature extraction - Principal Component Analysis (PCA)

- Crea combinaciones lineales de las características originales.
- Las nuevas características son ortogonales = que no están correlacionadas.
- No supervisado.
- Los nuevos componentes se ordenan según la máxima varianza: PC1 explica la mayor variación en su conjunto de datos, PC2 explica la segunda variación más, ...

Feature extraction - Linear Discriminant Analysis (LDA)

- Crea combinaciones lineales de sus características originales.
- PCA vs LDA
 - LDA es supervisado
 - LDA maximiza la separabilidad entre clases.



Fuente: <https://www.youtube.com/watch?v=azXCzI57Yfc>

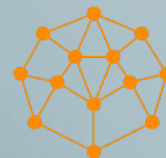
Time to Code!!! :D



Saturdays.AI
LATAM



#3.4 Conclusión



Saturdays.AI
LATAM

¿Que técnica usar? ¿Cuál es mejor?

"No Free Lunch" theorem

En resumen Ningún algoritmo funciona mejor para cada problema.

Solución: Prueba muchos algoritmos diferentes para tu problema, usa un "test set" de datos para evaluar el rendimiento y seleccionar al ganador.



Contacto



María Inés Calderón Zetter

ines@saturdays.ai

<https://www.linkedin.com/in/ineszetter/>



Favio Vázquez

favio@closter.net

<https://www.linkedin.com/in/faviovazquez/>



Lea Vega Romero

lea@saturdays.ai

<https://www.linkedin.com/in/lea-vega-66a18011b/>



Saturdays.AI
LATAM