

Data Fundamentals

Con Python

24 de Febrero, 2023



red.es



"El FSE invierte en tu futuro"
Fondo Social Europeo

Índice

1. Data Manipulation
 - 1.1.Data Wrangling
 - 1.2.Incompletos
 - 1.3.Atípicos
 - 1.4.Incoherentes
2. Ejercicio de EDA

Aprendiendo a Pensar como un programador

Data Manipulation

Data Wrangling

Data Wrangling

Hemos tratado datos, visto ejemplos y casuísticas...

Data Wrangling

Hemos tratado datos, visto ejemplos y casuísticas...

Hay más contenido teórico-práctico para manejar mejor los datos

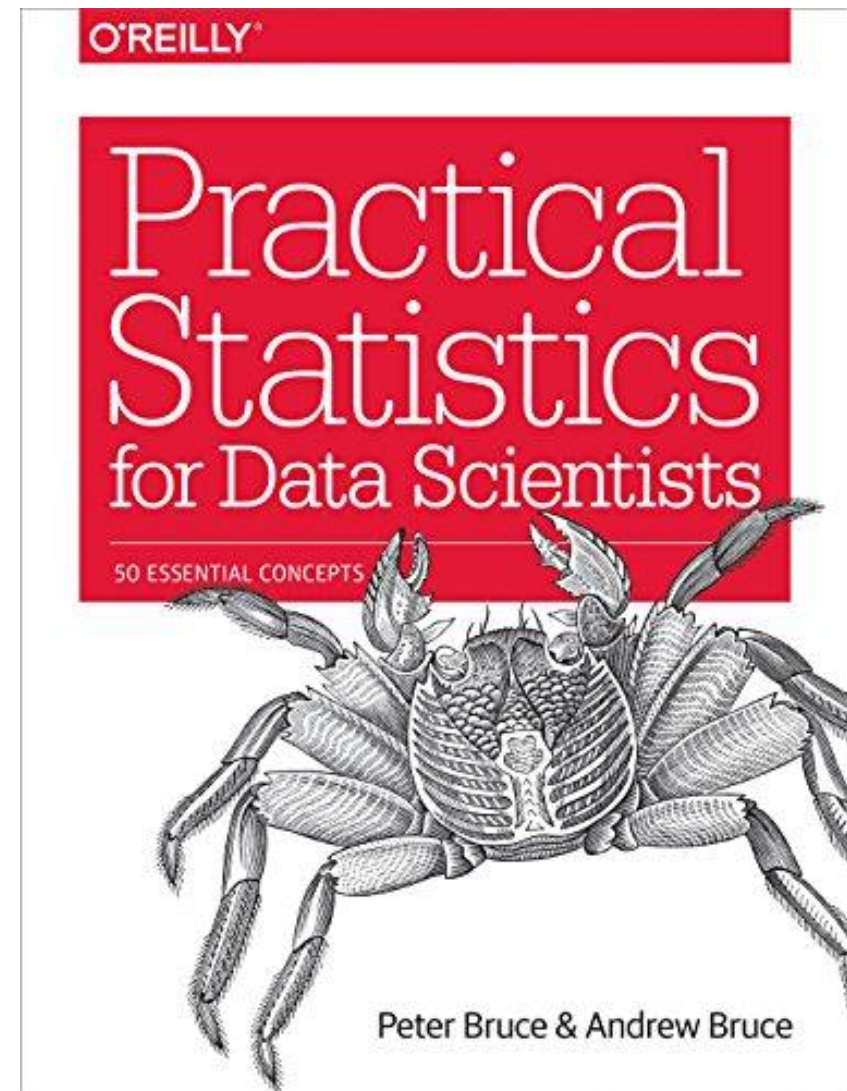
- ❖ Data Cleansing
- ❖ Data Cleaning
- ❖ Data Preparation
- ❖ ...

Data Wrangling

Hemos tratado datos, visto ejemplos y casuísticas...

Hay más contenido teórico-práctico para manejar mejor los datos

- ❖ Data Cleansing
- ❖ Data Cleaning
- ❖ Data Preparation
- ❖ ...

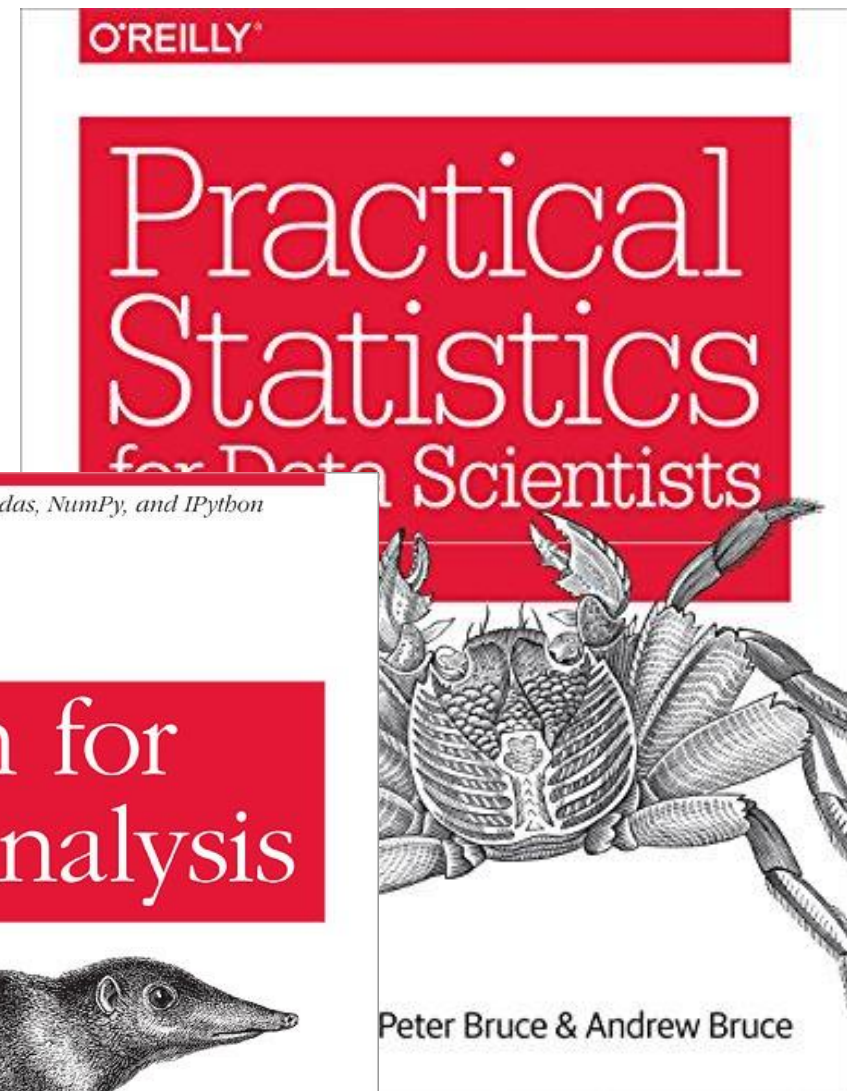
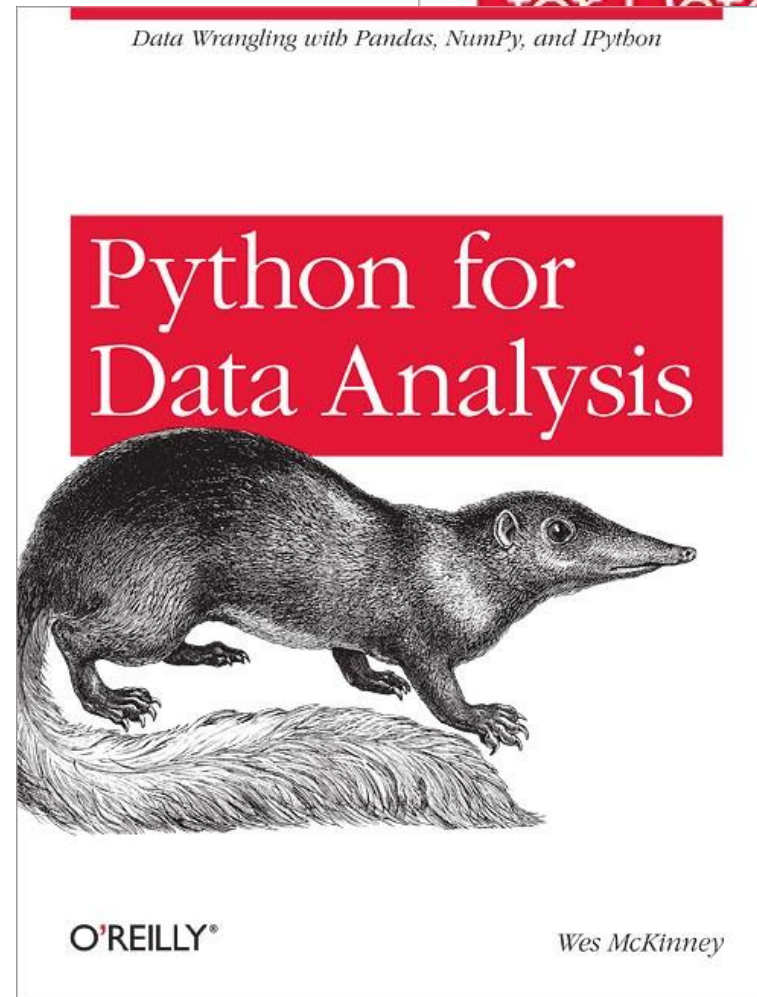


Data Wrangling

Hemos tratado datos, visto ejemplos y casuísticas...

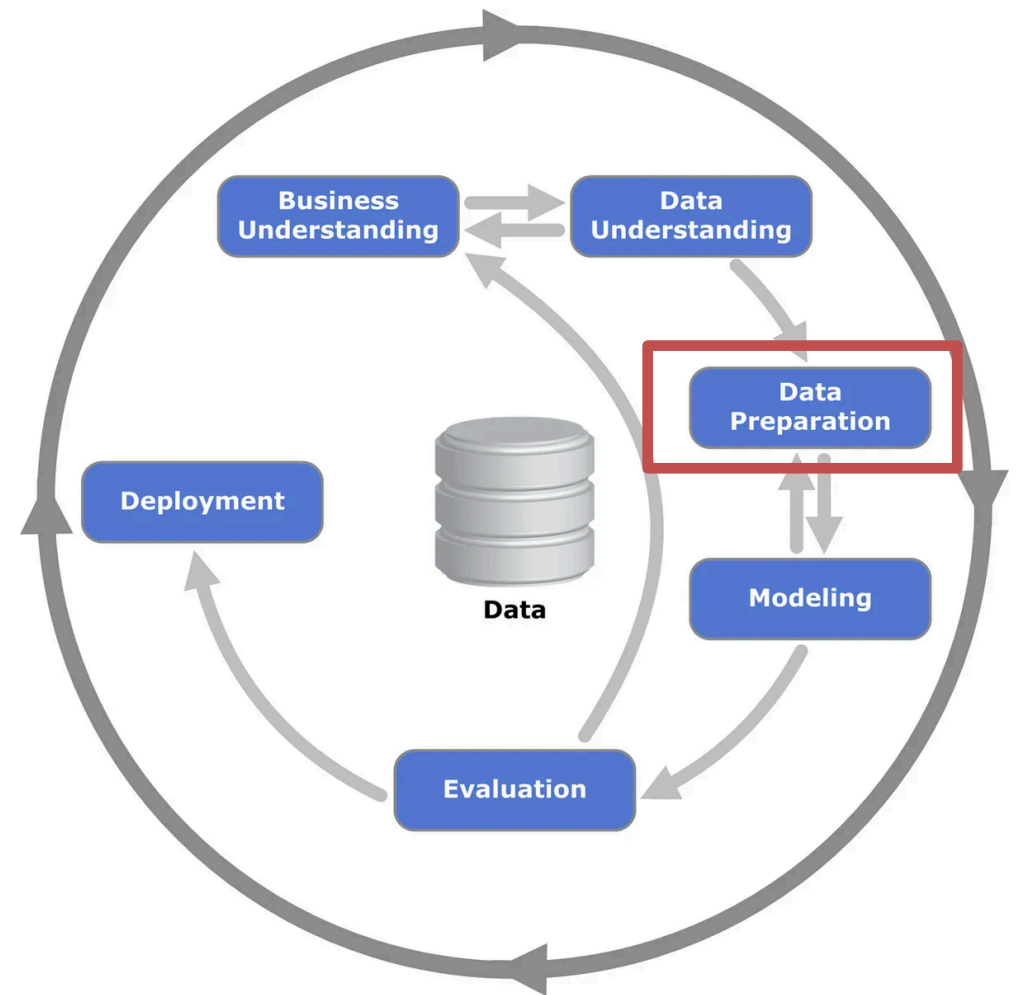
Hay más contenido teórico-práctico para manejar mejor los datos

- ❖ Data Cleansing
- ❖ Data Cleaning
- ❖ Data Preparation
- ❖ ...



Siguiente Paso

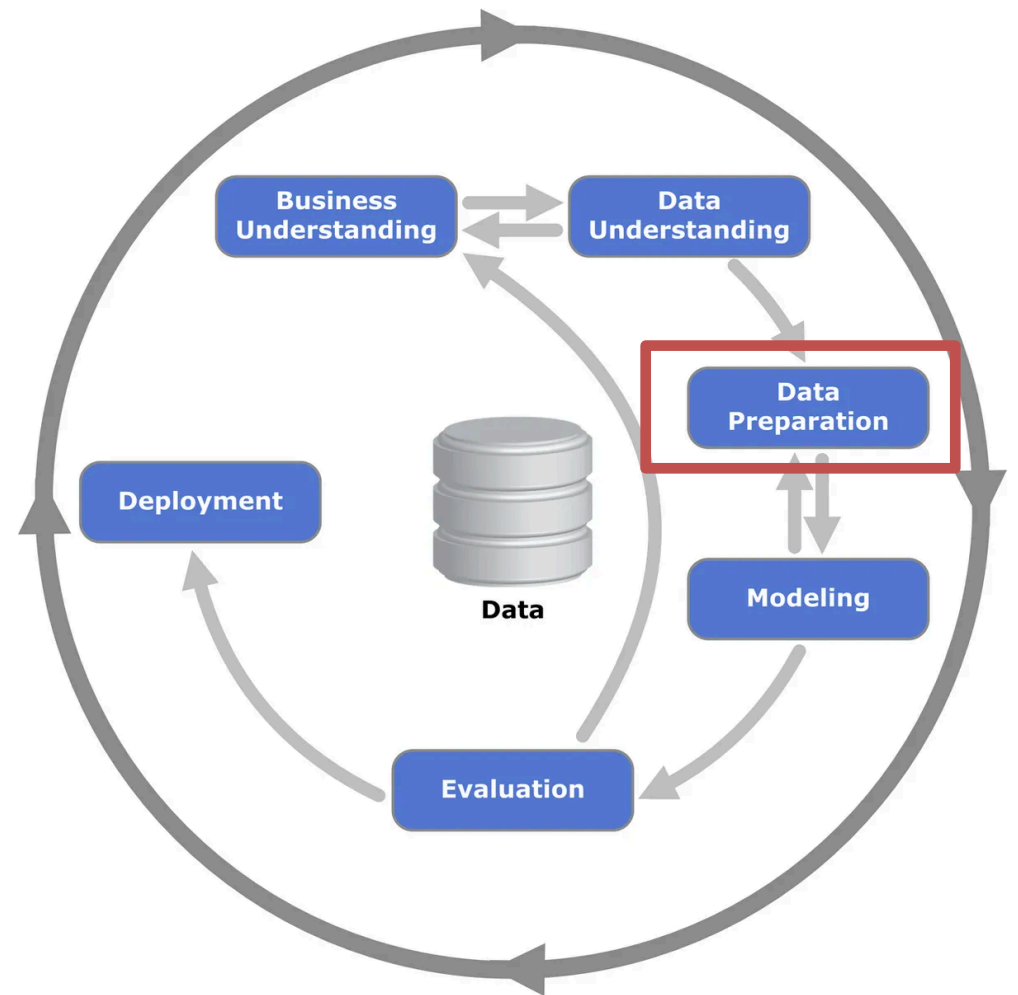
Una vez tenemos bien explorados los datos:



Siguiente Paso

Una vez tenemos bien explorados los datos:

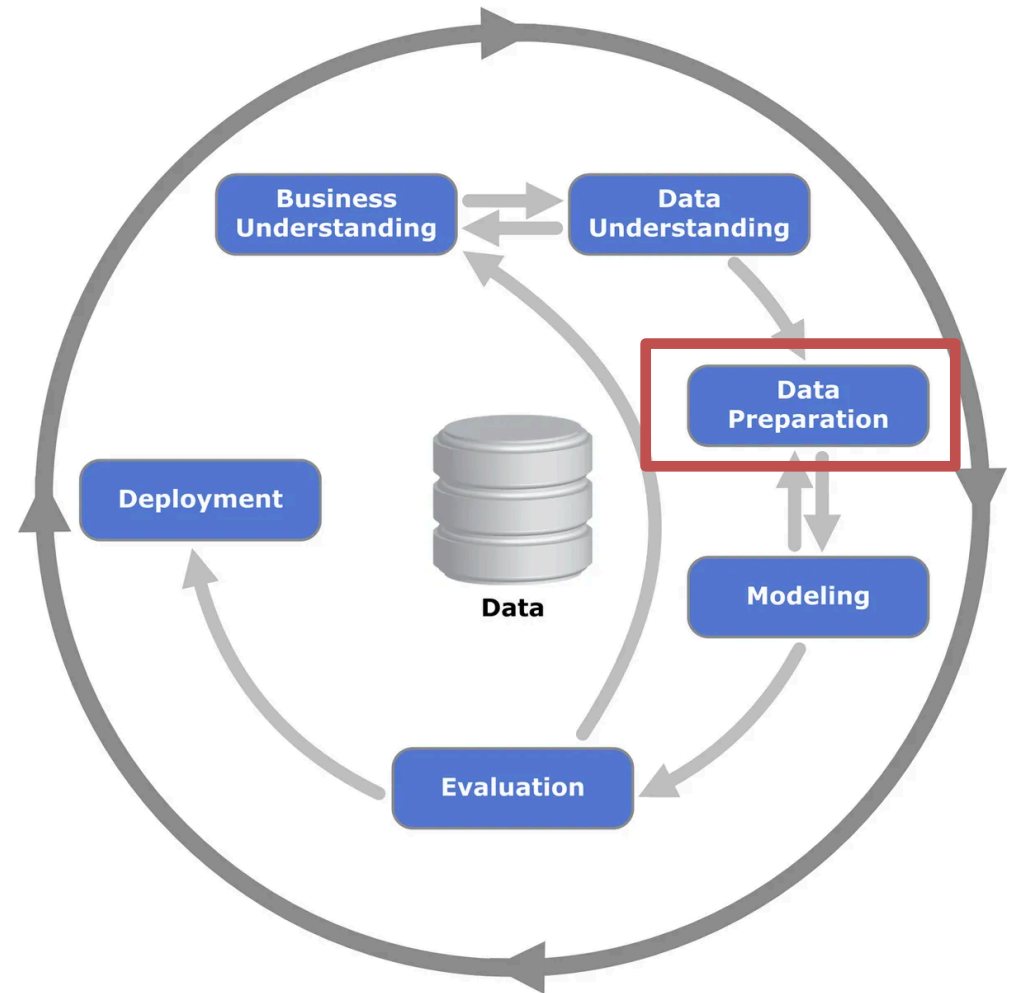
❖ Incompletitudes



Siguiente Paso

Una vez tenemos bien explorados los datos:

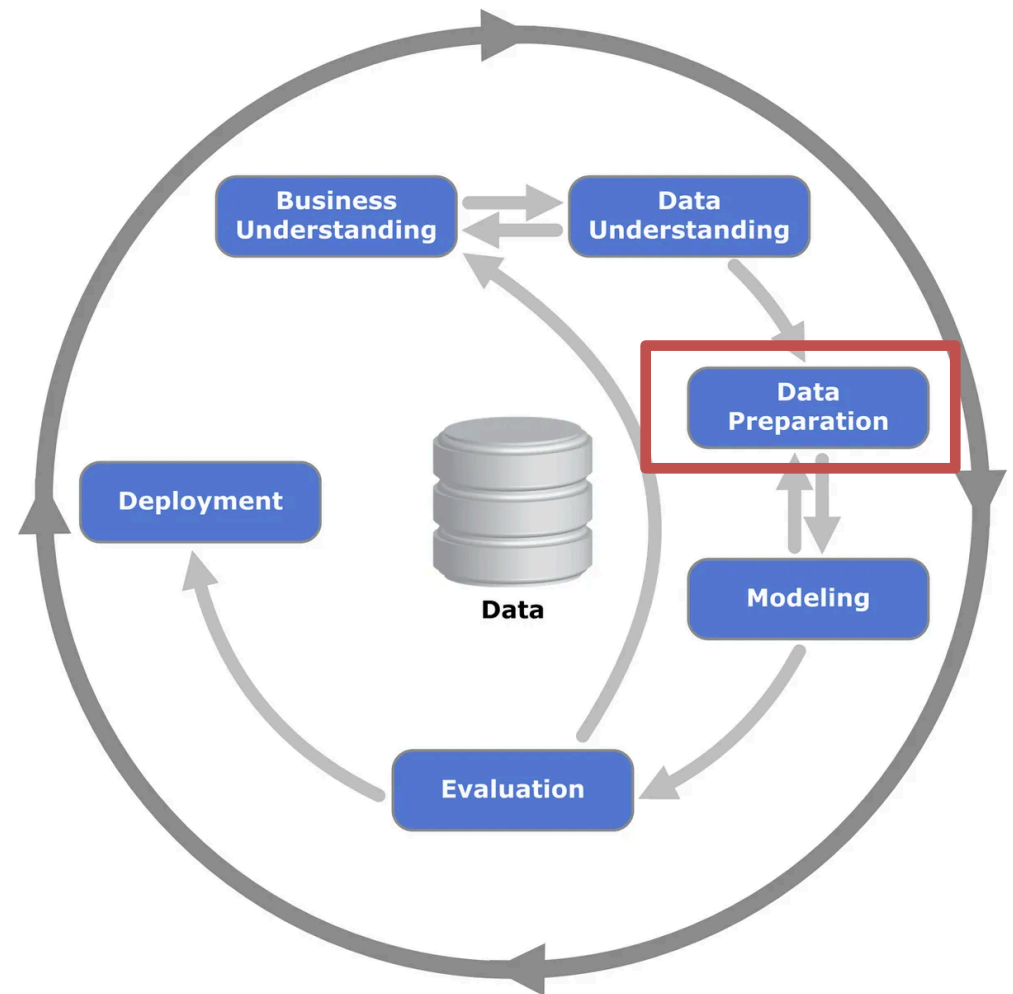
- ❖ Incompletitudes
- ❖ Nulos



Siguiente Paso

Una vez tenemos bien explorados los datos:

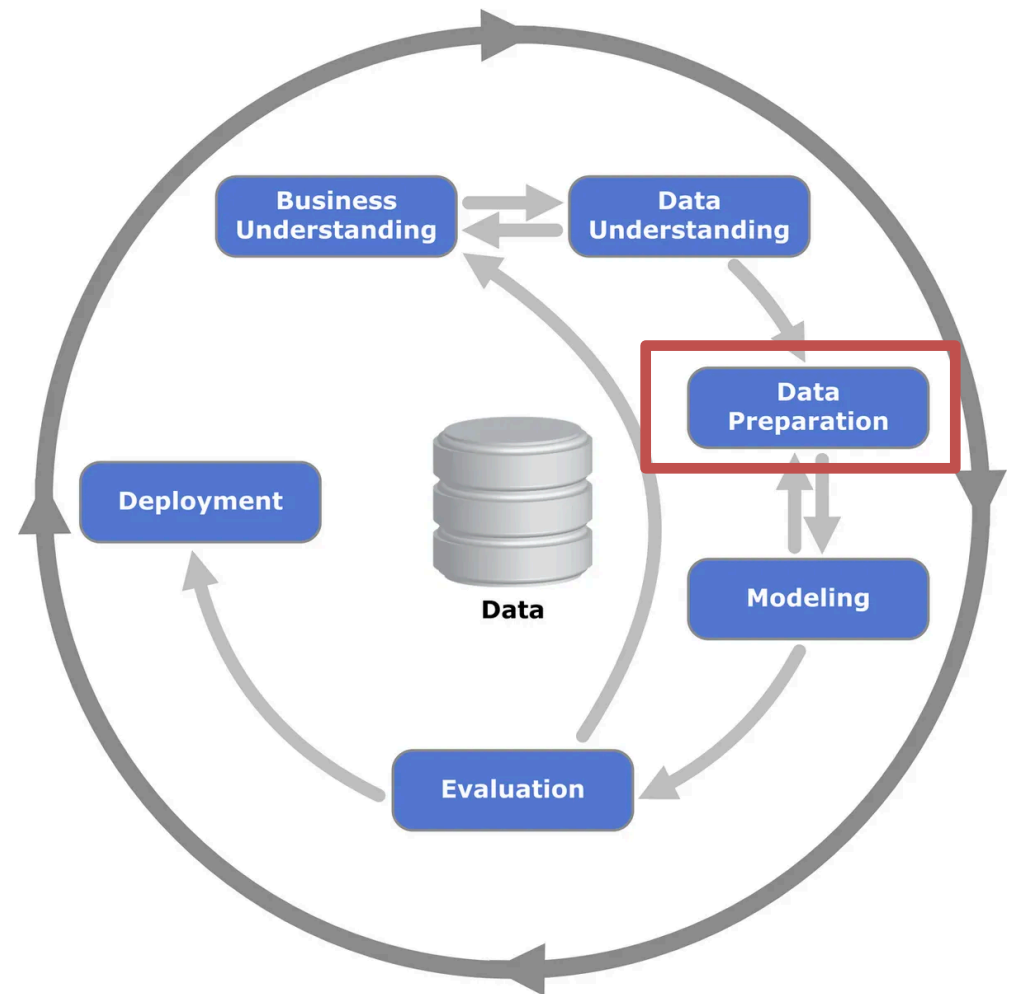
- ❖ Incompletitudes
- ❖ Nulos
- ❖ Incoherentes



Siguiente Paso

Una vez tenemos bien explorados los datos:

- ❖ Incompletitudes
- ❖ Nulos
- ❖ Incoherentes
- ❖ Inexistentes



Incompletos

Unidades	Edad	Sexo	PAS	PAD	Peso	Talla
1	×	×	?	?	×	×
⋮	⋮	⋮	⋮	⋮	⋮	⋮
100	×	×	?	?	×	×
101	×	×	×	×	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
200	×	×	×	×	?	?
201	×	×	?	?	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
300	×	×	?	?	?	?
301	×	×	×	×	×	×
⋮	⋮	⋮	⋮	⋮	⋮	⋮
1.000	×	×	×	×	×	×

Nota: El símbolo ? representa los valores ausentes
y ×, los observados

Incompletos

Tenemos que identificarlos y “marcarlos”

Unidades	Edad	Sexo	PAS	PAD	Peso	Talla
1	x	x	?	?	x	x
⋮	⋮	⋮	⋮	⋮	⋮	⋮
100	x	x	?	?	x	x
101	x	x	x	x	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
200	x	x	x	x	?	?
201	x	x	?	?	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
300	x	x	?	?	?	?
301	x	x	x	x	x	x
⋮	⋮	⋮	⋮	⋮	⋮	⋮
1.000	x	x	x	x	x	x

Nota: El símbolo ? representa los valores ausentes
y x, los observados

Incompletos

Tenemos que identificarlos y “marcarlos”

-1, Null, 0, Símbolos... Tenemos que saber dónde están y colocar una señal.

Unidades	Edad	Sexo	PAS	PAD	Peso	Talla
1	×	×	?	?	×	×
⋮	⋮	⋮	⋮	⋮	⋮	⋮
100	×	×	?	?	×	×
101	×	×	×	×	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
200	×	×	×	×	?	?
201	×	×	?	?	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
300	×	×	?	?	?	?
301	×	×	×	×	×	×
⋮	⋮	⋮	⋮	⋮	⋮	⋮
1.000	×	×	×	×	×	×

Nota: El símbolo ? representa los valores ausentes y ×, los observados

Incompletos

Tenemos que identificarlos y “marcarlos”

-1, Null, 0, Símbolos... Tenemos que saber dónde están y colocar una señal.

Unidades	Edad	Sexo	PAS	PAD	Peso	Talla
1	×	×	?	?	×	×
⋮	⋮	⋮	⋮	⋮	⋮	⋮
100	×	×	?	?	×	×
101	×	×	×	×	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
200	×	×	×	×	?	?
201	×	×	?	?	?	?
⋮	⋮	⋮	⋮	⋮	⋮	⋮
300	×	×	?	?	?	?
301	×	×	×	×	×	×
⋮	⋮	⋮	⋮	⋮	⋮	⋮
1.000	×	×	×	×	×	×

Nota: El símbolo ? representa los valores ausentes y ×, los observados

Incompletos

Incompletos

Habr  veces que no tendremos buena distribuci n de los datos, cuando esto ocurra diremos que est n "imbalanceados" o que no tienen dispersi n suficiente.

Incompletos

Habr  veces que no tendremos buena distribuci n de los datos, cuando esto ocurra diremos que est n "imbalanceados" o que no tienen dispersi n suficiente.

 Qu  hacemos en este caso?

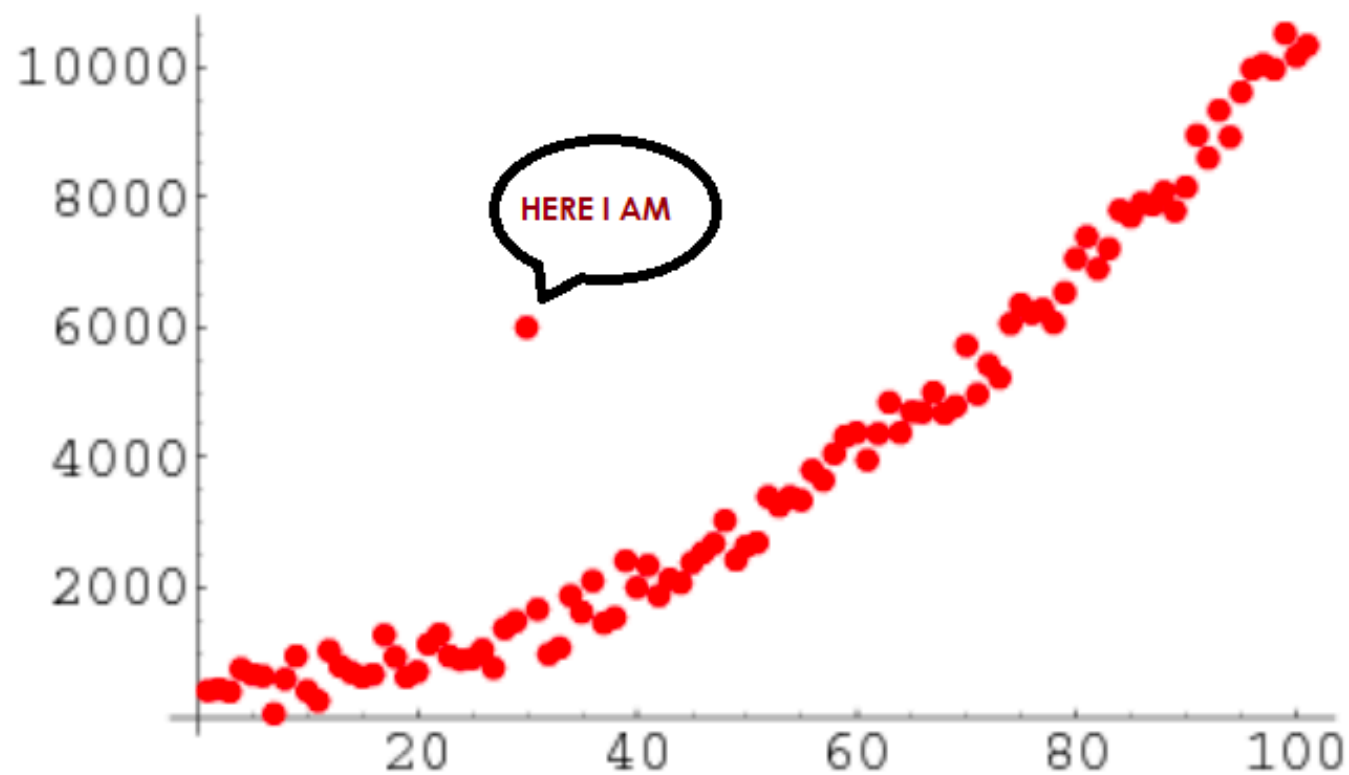
Valores Atípicos

Valores Atípicos

Valores fuera de la normal general...

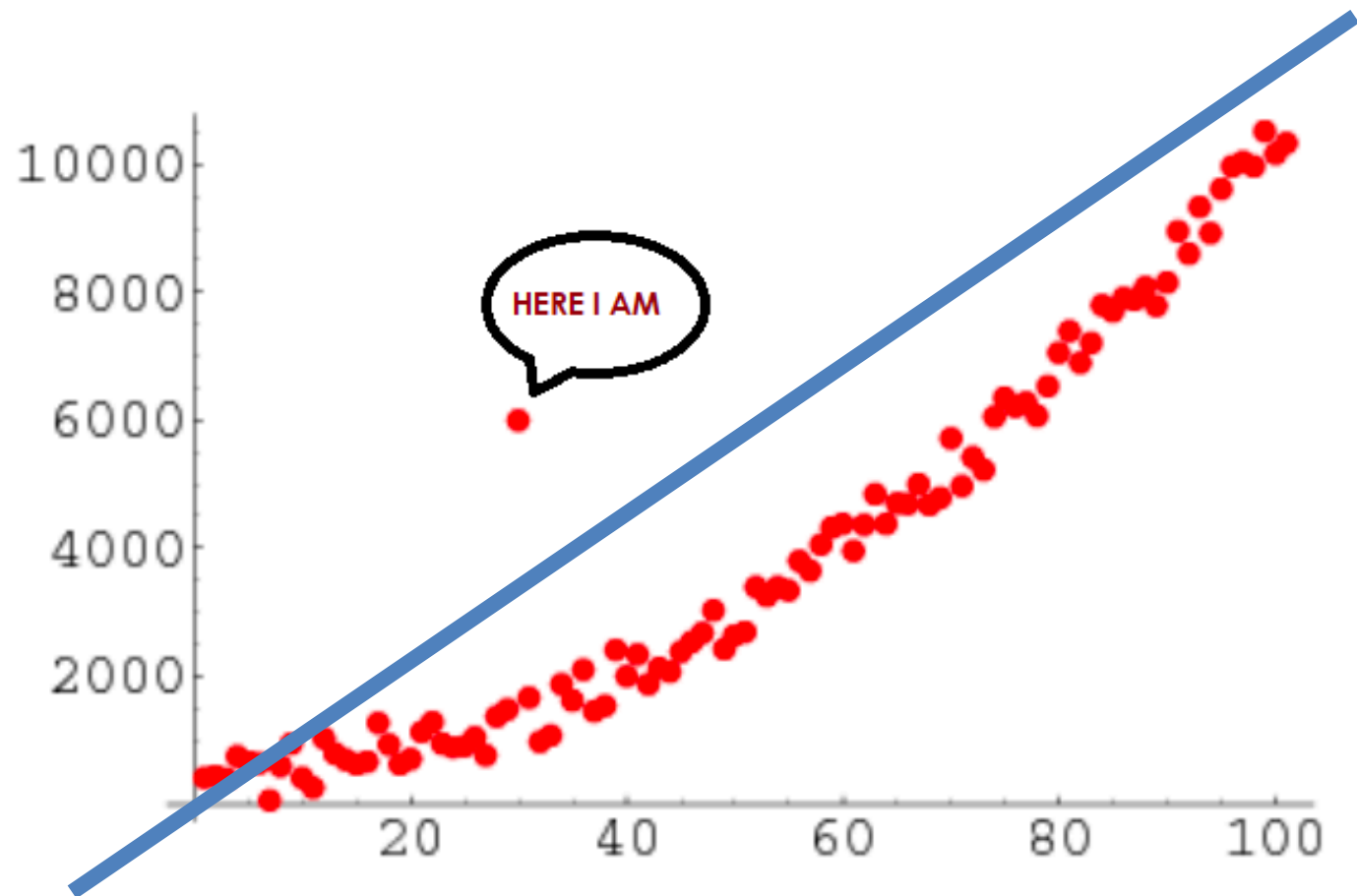
Valores Atípicos

Valores fuera de la normal general...



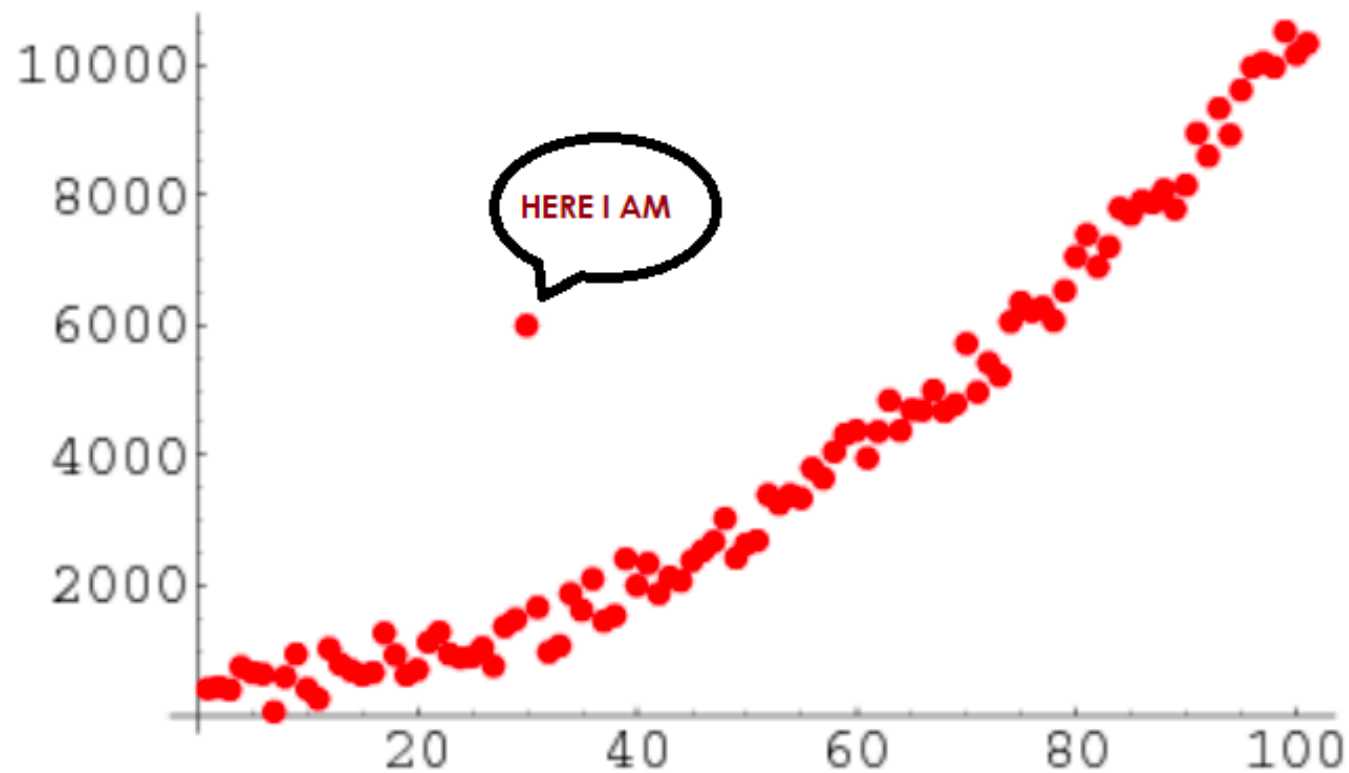
Valores Atípicos

Valores fuera de la normal general...

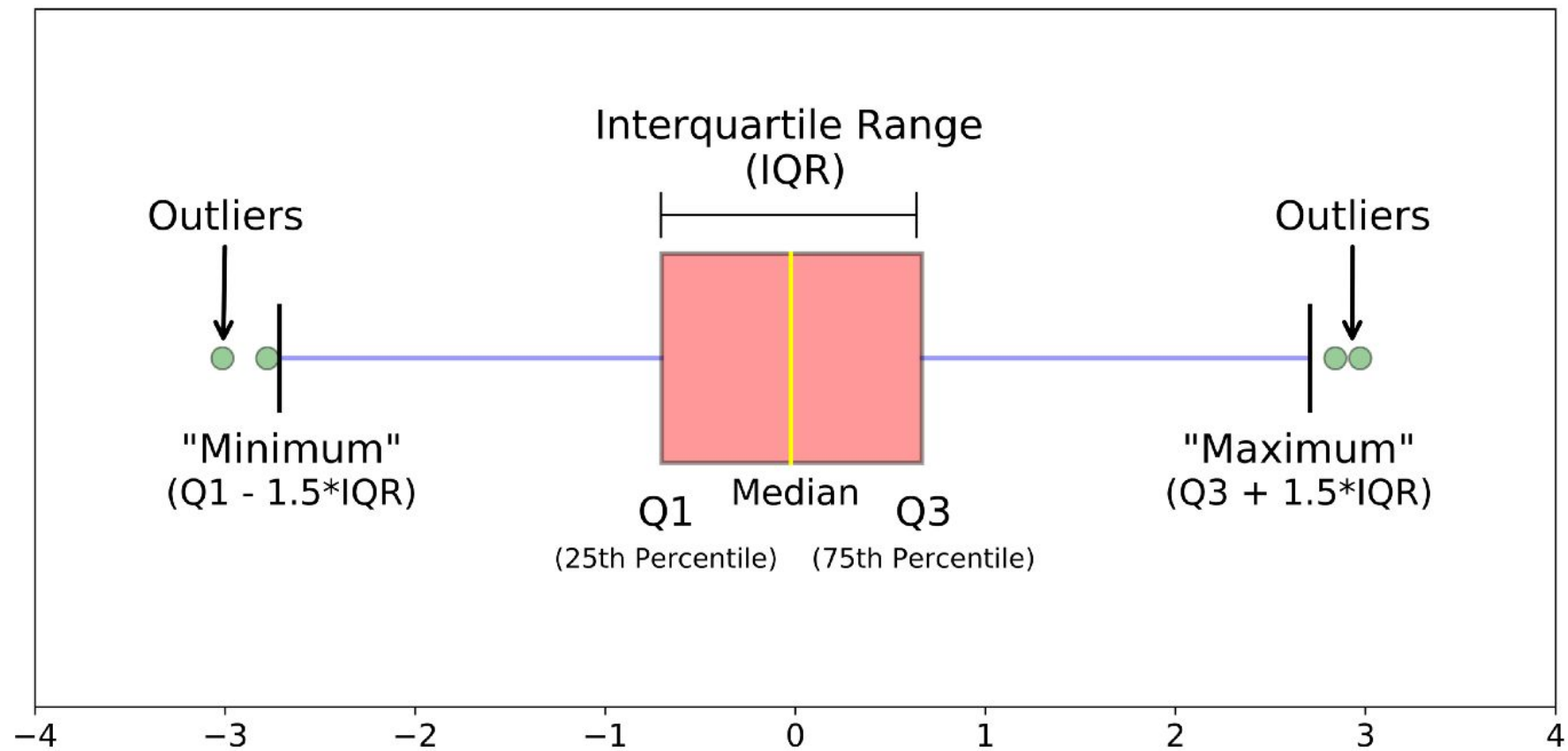


Valores Atípicos

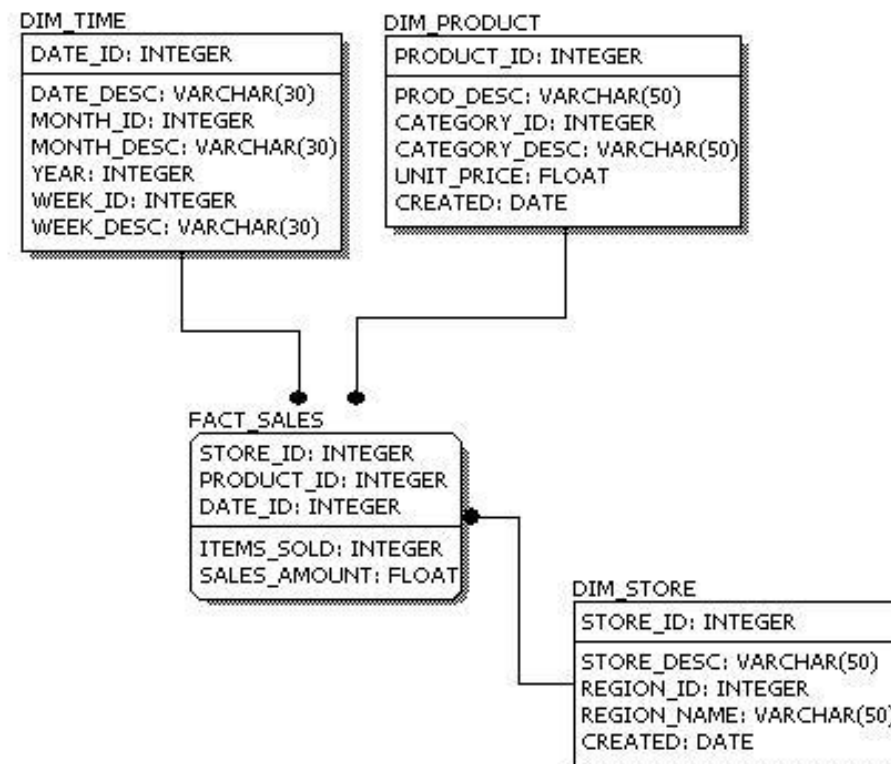
OUTLIERS!!!!



Algunas Métricas

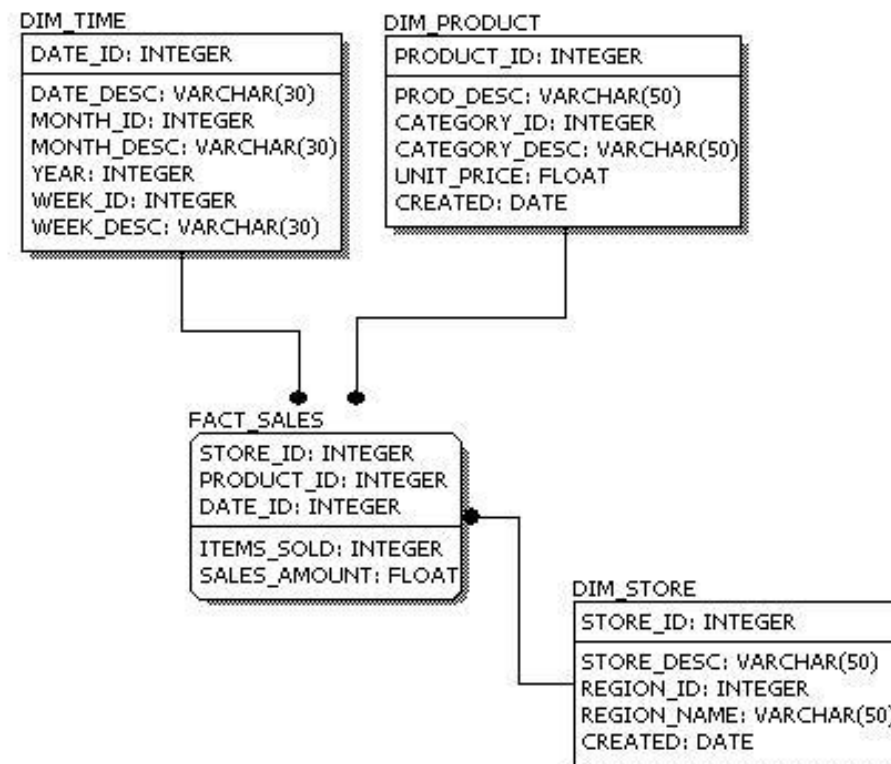


Incoherentes



Incoherentes

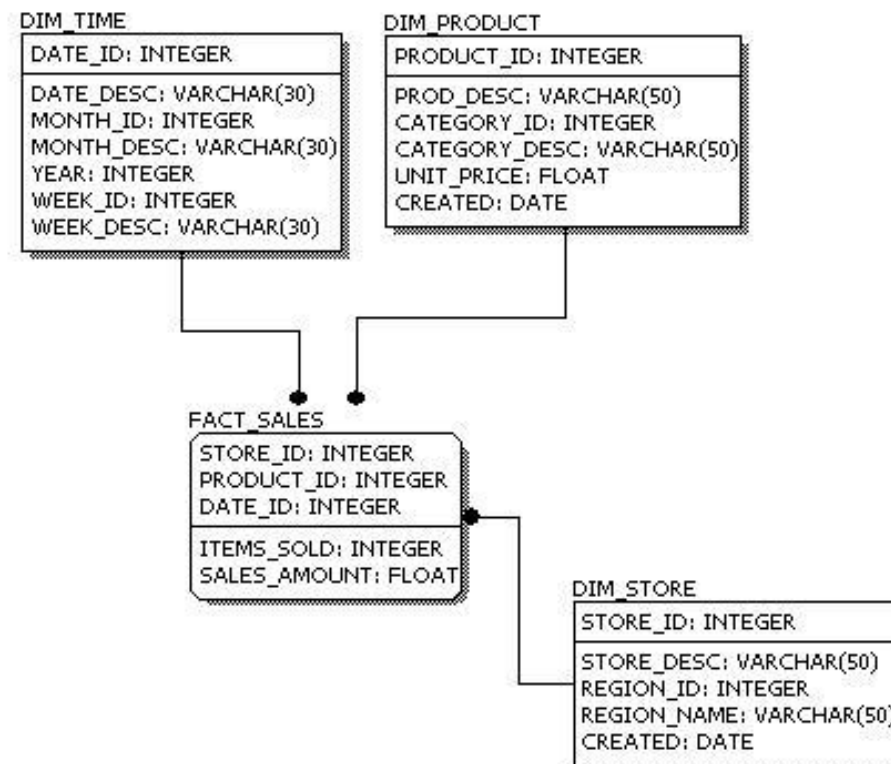
¿Cómo son los datos? ¿Cuál es el rango de valores posible? ¿Existe un orden?



Incoherentes

¿Cómo son los datos? ¿Cuál es el rango de valores posible? ¿Existe un orden?

¡Un Diccionario de Datos es la respuesta!



Aprendiendo a Pensar como un programador

Ejercicios

Ejercicio 1

Ejercicio guiado de EDA:

Buscamos Outliers

Identificamos datos incompletos

Creamos un diccionario de Datos

...

Generamos datos Sintéticos para los imbalanceados





red.es



"El FSE invierte en tu futuro"

Fondo Social Europeo

