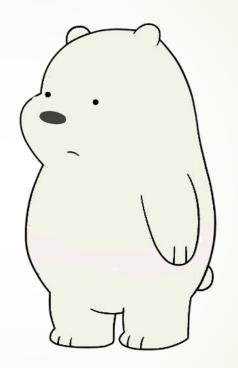
# "Un tú por tú, en el análisis de datos"

#### **Pandas Vs Polars**







#### **Pythonistas GDL**



- ·charlas
- ·comida
- ·convivencia
- sorpresas

Patrocinado por



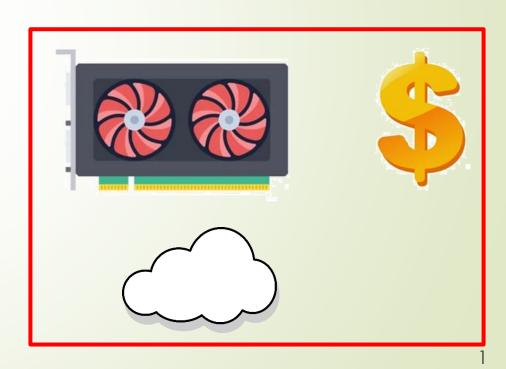
# **Agenda**

- Definición y Contextualización del Problema
- Núcleo Funcional de Pandas
- Comparativa de Librerías
- Comparativa de Rendimiento entre Librerías
- Ejecutando Pandas y Polars
- Egger vs Lazy Frames
- Memoria RAM
- → Conclusión
- Preguntas y Respuestas

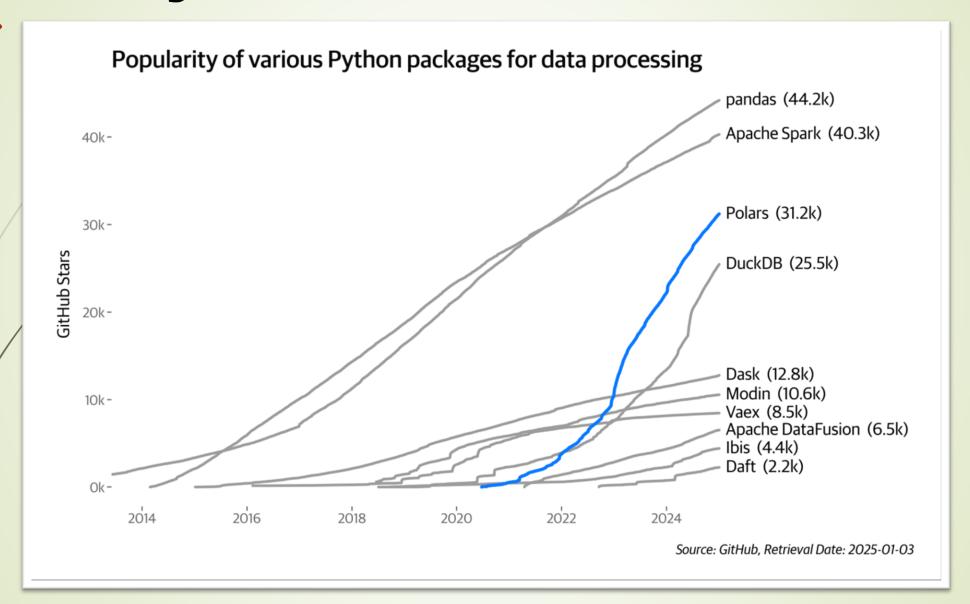
# Definición y Contextualización del Problema



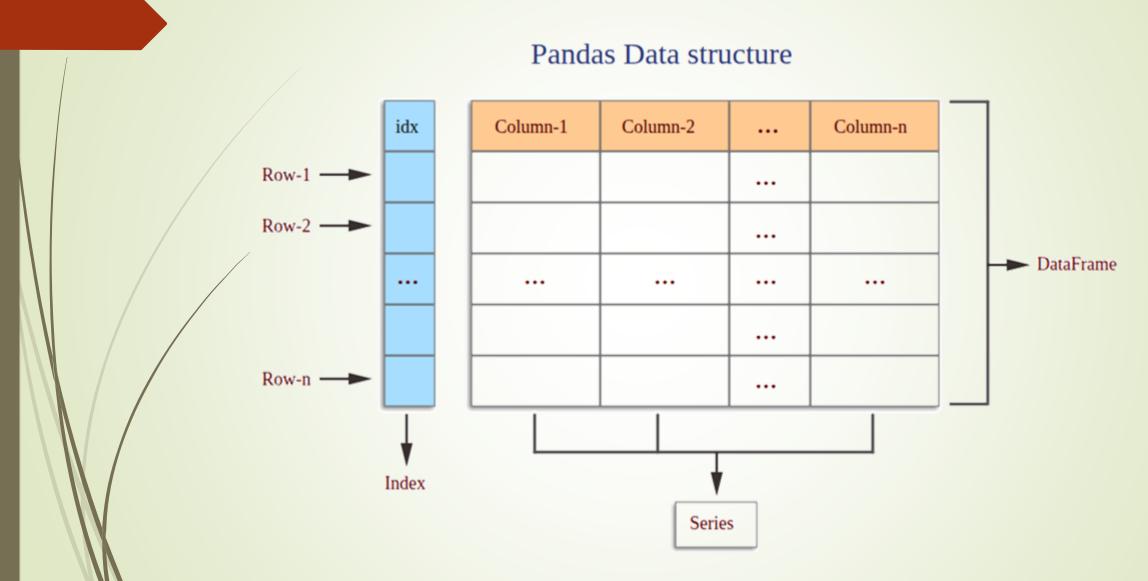




# Definición y Contextualización del Problema



#### **Núcleo Funcional de Pandas**



# Comparativa de Librerías

#### Software



**Wes Mckinney** 2008

Pandas	Polars
Pandas DataFrame	Polars DataFrame
Numpy	Polar Series
Python-dateutil	Polars Datetime
Tzdata	Polars Time Zone
Analisis Temporal	Analisis Temporal
Lectura/Escritura de	Lectura/Escritura de
Archivos	Archivos
csv,excel,parquet	csv,excel,parquet
No disponible	Lazy Frame

Rust + Apache Arrow



**Ritchie Vink** 2020

# **Hardware Pandas**

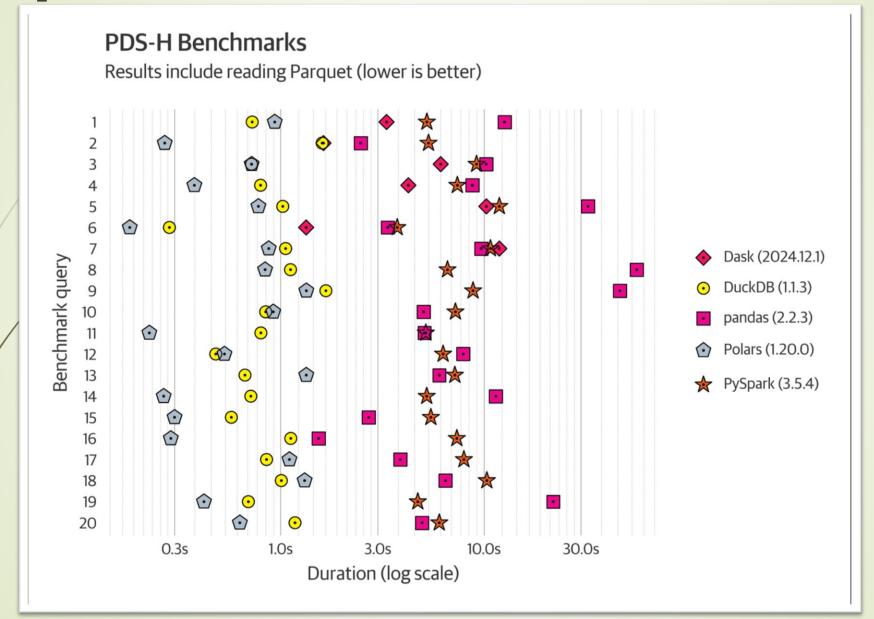
#### **Polars** Operations No usa múltiples núcleos Ejecución Paralela de CPU por defecto Automática No tiene soporte nativo para GPU Diseño Eficiente en CPU No escala bien datasets mayores a 1Gb Escalamiento horizontal Alto Consumo de Bajo Consumo de

Memoria

Backend Python +C

Memoria

#### Comparativa de Rendimiento de Librerías



```
import pandas as pd
import polars as pl
import sys
animals_pd = pd.read_csv("data/animals.csv", sep=",", header=0)
animals_pl = pl.read_csv("data/animals.csv", separator=",", has_header=True)
print(sys.getsizeof(animals_pd)) # Tamaño en bytes del objeto 'x'
print(sys.getsizeof(animals_pl)) # Tamaño en bytes del objeto 'x'
print(f"{type(animals_pd) = }")
print(f"{type(animals_pl) = }")
 3830
 48
 type(animals_pd) = <class 'pandas.core.frame.DataFrame'>
 type(animals_pl) = <class 'polars.dataframe.frame.DataFrame'>
```

**Pandas** 

Polars

		animal	<b>‡</b>	class	÷	habitat	<b>‡</b>	diet	<b>‡</b>	lifespan	¢	status	<b>‡</b>	features	÷	weight	÷
	0	dolphin		mammal		oceans/rivers		carnivore			40	least concern		high intelligence			150.0
ш	1	duck		bird		wetlands		omnivore			8	least concern		waterproof feathers			3.0
ш	2	elephant		mammal		savannah		herbivore			60	endangered		large ears and trunk			8000.0
ш	3	ibis		bird		wetlands		omnivore			16	least concern		long, curved bill			1.0
ш	4	impala		mammal		savannah		herbivore			12	least concern		long, curved horns			70.0
ш	5	kudu		mammal		savannah		herbivore			15	least concern		spiral horns			250.0
ш	6	narwhal		mammal		arctic ocean		carnivore			40	near threatened		long, spiral tusk			NaN
ш	7	panda		mammal		forests		herbivore			20	vulnerable		black and white coloration			100.0
	8	polar bear		mammal		arctic		carnivore			25	vulnerable		thick fur and blubber			720.0
1	9	ray		fish		oceans		carnivore			20	NaN		flat, disc-shaped body			90.0

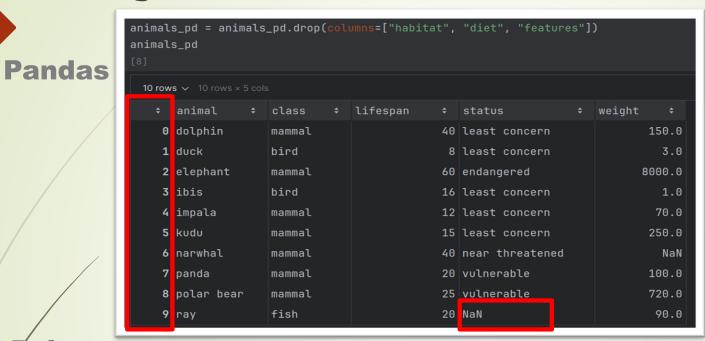
ш	animal	<b>‡</b>	class	<b>‡</b>	habitat	<b>‡</b>	diet ÷	lifespan	<b>‡</b>	status	<b>‡</b>	features	÷ 1	weight	<b>‡</b>
Ш	"dolphin"		"mammal"		"oceans/rivers"		"carnivore"		40	"least concern"		"high intelligence"	1	150	
Ш	"duck"		"bird"		"wetlands"		"omnivore"		8	"least concern"		"waterproof feathers"	3	3	
	"elephant"		"mammal"		"savannah"		"herbivore"		60	"endangered"		"large ears and trunk"	8	8000	
Ш	"ibis"		"bird"		"wetlands"		"omnivore"		16	"least concern"		"long, curved bill"	1	1	
Ш	"impala"		"mammal"		"savannah"		"herbivore"		12	"least concern"		"long, curved horns"	5	70	
Ш	"kudu"		"mammal"		"savannah"		"herbivore"		15	"least concern"		"spiral horns"	2	250	
Ш	"narwhal"		"mammal"		"arctic ocean"		"carnivore"		40	"near threatened	"	"long, spiral tusk"	r	null	
Ш	"panda"		"mammal"		"forests"		"herbivore"		20	"vulnerable"		"black and white coloration"	1	100	
П	"polar bear"		"mammal"		"arctic"		"carnivore"		25	"vulnerable"		"thick fur and blubber"	7	720	
Ш	"ray"		"fish"		"oceans"		"carnivore"		20	1111		"flat, disc-shaped body"	ç	90	

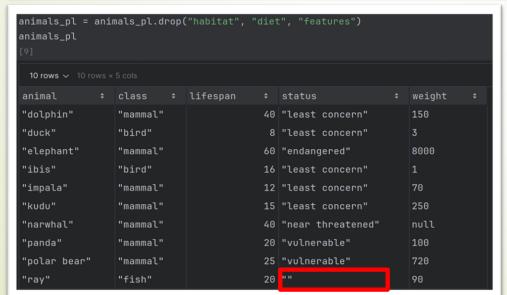
#### **Pandas**

```
animals_pd["animal"]
         dolphin
            duck
        elephant
 3
            ibis
          impala
            kudu
         narwhal
           panda
      polar bear
             ray
 Name: animal, dtype: object
```

#### **Polars**

```
animals_pl.get_column("animal")
  10 rows ∨ 10 rows × 1 cols
animal
"dolphin"
"duck"
"elephant"
"ibis"
"impala"
"kudu"
"narwhal"
"panda"
"polar bear"
"ray"
```







#### **Eager vs LazyFrames**

```
lazy_query = (
    pl.scan_csv("data/animals.csv")
        .group_by("class")
        .agg(pl.col("weight").mean())
        .filter(pl.col("class") == "mammal")
)
```

```
lazy_query.show_graph(optimized=False)

/ [3] 67ms

FILTER BY [(col("class")) == ("mammal")]

AGG [col("weight").mean()]
BY
[col("class")]

Csv SCAN [data/animals.csv]
π */8;
```

```
Lazy_query.show_graph()

/[4] 36ms

AGG [col("weight").mean()]
BY
[col("class")]

Csv SCAN [data/animals.csv]
π 2/8;
σ [(col("class")) == ("mammal")]
```

#### **Eager vs LazyFrames**

```
%%time
trips = pl.read_parquet("data/taxi/yellow_tripdata_*.parquet")
sum_per_vendor = trips.group_by("VendorID").sum()
income_per_distance_per_vendor = sum_per_vendor.select(
    "VendorID",
    income_per_distance=pl.col("total_amount") / pl.col("trip_distance")
top_three = income_per_distance_per_vendor.sort(
   by="income_per_distance", descending=True
).head(3)
top_three
 CPU times: user 10.2 s, sys: 2.89 s, total: 13 s
 Wall time: 1.27 s
```

**Eager vs LazyFrames** 

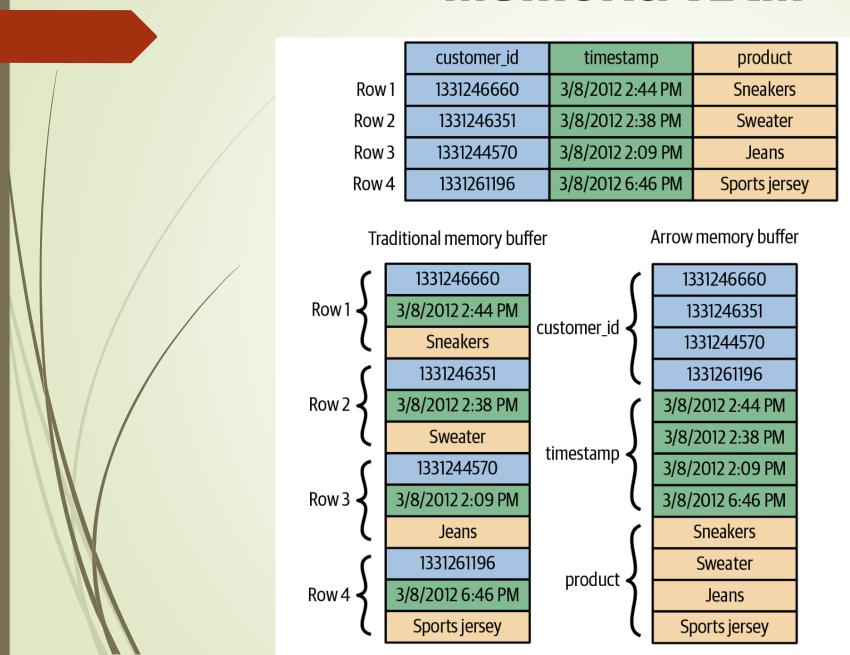
```
%%time
trips = pl.scan_parquet("data/taxi/yellow_tripdata_*.parquet")
sum_per_vendor = trips.group_by("VendorID").sum()

income_per_distance_per_vendor = sum_per_vendor.select(
    "VendorID",
    income_per_distance=pl.col("total_amount") / pl.col("trip_distance"),
)

top_three = income_per_distance_per_vendor.sort(
    by="income_per_distance", descending=True
).head(3)
print(top_three.describe())
top_three.collect()
```

str	f64	f64
count	3.0	3.0
null_count	0.0	0.0
mean	4.0	5.487613
std	2.645751	0.867551
min	1.0	4.731557
25%	5.0	5.296493
50%	5.0	5.296493
75%	6.0	6.434789
max	6.0	6.434789
CPU times: use	er 3.94 s,	sys: 46.6 ms, total 3.99 s
Wall time: 432	2 ms	

#### **Memoria RAM**



#### Conclusión

- Pandas es ideal para tareas educativas, prototipos rápidos, análisis exploratorio con data sets no mayores a 1Gb
- Polars es mejor en rendimiento para grandes volúmenes de datos con mejoras en operaciones de filtrado y manipulación de datos
- Sin embargo se puede adoptar una estrategia hibrida.

# **Preguntas y Respuestas**

# Redes Sociales



in migueltlapa



