

# Análisis y visualización con pandas, nivel intermedio

Departamento de Física.

---

Corodinadora: C Loyola

Profesores C Femenías / F Bugini / D Basantes

Primer Semestre 2025

Universidad Andrés Bello

Departamento de Física y Astronomía



# Índice

- Fuentes de datos
- Carga y limpieza
- Diagnóstico de faltantes
- Agrupaciones avanzadas
- Reshape melt y pivot
- Series temporales
- Merge ilustrativo
- Visualización integrada
- Estilos
- Exportar resultados
- Rendimiento
- Mini proyecto en parejas
- Guía rápida de solución
- Recursos recomendados
- Preguntas frecuentes
- Próximos temas

## Fuentes de datos

---

- Meteoritos NASA (limitado a 5 000 filas)  
`https://data.nasa.gov/resource/y77d-th95.csv?`  
`\protect\TU\textdollar\protect\TU\`  
`textdollarlimit=5000`
- Exoplanetas Seaborn  
`https://raw.githubusercontent.com/mwaskom/`  
`seaborn-data/master/planets.csv`
- Masas PDG (extracto)  
`https://raw.githubusercontent.com/`  
`particle-physics-book/data/master/pdg_mass.csv`

Copia estas direcciones directamente en tu notebook.

## Carga y limpieza

---

## Carga y limpieza ∈ Lectura de los tres conjuntos

```
1 met =  
  ↪ pd.read_csv("https://data.nasa.gov/resource/y77d-th95.csv?$$limit=  
2 pln =  
  ↪ pd.read_csv("https://raw.githubusercontent.com/mwaskom/seaborn-data/  
3 pdg =  
  ↪ pd.read_csv("https://raw.githubusercontent.com/particle-physics-b
```

‘read\_csv’ descarga y parsea automáticamente; no requiere bibliotecas externas.

```
1 met = met[met.mass.notna()].copy()
2 met["massKg"] = met.mass.astype(float) / 1000      # gramos a kg
3 met = met[met.massKg < 1e4]                        # quitar >10
   ↪ t
```

Tres líneas bastan para dejar el conjunto listo para análisis.

## Diagnóstico de faltantes

---



## Diagnóstico de faltantes ∈ Porcentaje de valores nulos

```
1 faltantes = met.isna().mean() * 100
2 faltantes.sort_values(ascending=False).head()
```

Así priorizas las columnas que necesitan limpieza.

## Agrupaciones avanzadas

---

```
1 met["year"] = pd.to_datetime(met.year,  
    ↪ errors="coerce").dt.year  
2 decade = (met.groupby(met.year // 10 * 10)["massKg"]  
3             .agg(conteo='size', promedio='mean',  
    ↪ total='sum'))  
4 decade.head()
```

Agrupar por transformación matemática del año genera la tabla por décadas.

## Agrupaciones avanzadas ∈ Pivot table visual

```
1 pln["decada"] = pln.year // 10 * 10
2 tabla = pln.pivot_table(index="method", columns="decada",
3                           values="distance", aggfunc="count")
4 plt.imshow(tabla.fillna(0), aspect='auto')
5 plt.colorbar(label="número de exoplanetas");
   ↪ plt.yticks(range(len(tabla)), tabla.index)
6 plt.xlabel("decada"); plt.show()
```

La matriz calor ayuda a descubrir métodos dominantes en cada período.

## Reshape melt y pivot

---

## Reshape melt y pivot ∈ Convertir tabla ancha a larga

```
1 wide = decade.reset_index().rename(columns={"year":"decada"})
2 long = wide.melt(id_vars="decada",
3                 var_name="metrica", value_name="valor")
4 long.head()
```

La forma “larga” es conveniente para bibliotecas de visualización como Seaborn.

# Series temporales

---

## Series temporales ∈ Resample anual de masa caída

```
1 metTs = (met.set_index(pd.to_datetime(met.year, format='%Y'))  
2           .massKg.resample("1Y").sum())  
3 metTs.plot(title="Masa total de meteoritos por año");  
   ↪ plt.ylabel("kg"); plt.show()
```

‘resample’ es la versión temporal de ‘groupby’.



## Merge ilustrativo

---

## Merge ilustrativo ∈ Normalizar masas de planetas

```
1 proton = pdg.loc[pdg.symbol == "p", "mass_GeV"].iat[0]
2 pln["massOverMp"] = pln.mass / proton
3 pln[["name", "mass", "massOverMp"]].head()
```

Combinamos física de partículas con astronomía en un solo DataFrame.

# Visualización integrada

---

```
1 pln.plot.scatter(x="mass", y="orbital_period",  
2                 alpha=0.4, logx=True, logy=True,  
3                 title="Masa vs periodo orbital")  
4 pln["distance"].plot.kde(title="Densidad de distancias");  
   ↪ plt.show()
```

La API integrada evita llamadas adicionales a Matplotlib.

# Estilos

---

```
1 plt.style.use("science")
2 decade.total.plot(marker='o')
3 plt.title("Toneladas de meteoritos por década");
  ↪ plt.ylabel("kg"); plt.show()
```

El resultado queda listo para copiar en un informe o paper.

Exportar resultados

---

```
1 decade.to_csv("meteoritos_por_decada.csv", index=True)  
2 pln.to_parquet("planetsCompact.parquet", index=False)
```

Parquet brinda alta compresión y lectura más veloz que CSV.



# Rendimiento



- Columnas con pocos valores → `'astype("category")'` para reducir memoria.
- `'read_csv'` con `'chunksize'` permite procesar archivos gigantes en porciones.
- `'eval'` y `'query'` compilan expresiones y evitan copias intermedias.

## Mini proyecto en parejas

---

1. Clasifica cada meteorito en continente según latitud y longitud.
2. Histograma logarítmico de masa por continente.
3. Pivot table de cantidad de impactos por década y continente.
4. Exporta la pivot table a un archivo Excel para compartir.

## Guía rápida de solución

---

```
1 bins = np.logspace(-2, 4, 30)
2 for cont in met.continent.unique():
3     met.loc[met.continent == cont, "massKg"].plot.hist(
4         bins=bins, alpha=0.5, label=cont)
5 plt.xscale("log"); plt.legend(); plt.show()
```

## Recursos recomendados

---

- Cheat-sheet oficial: [https://pandas.pydata.org/Pandas\\_Cheat\\_Sheet.pdf](https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf)
- Libro “Python Data Science Handbook”, capítulos ocho a once.
- Lista Awesome Public Datasets en GitHub.



## Preguntas frecuentes

---

- **Warning SettingWithCopy** → usa 'loc' y 'copy'.
- CSV muy grande → convierte a Parquet y gana 5-10 x en velocidad.
- Problemas en merge → usa 'validate="one\_to\_one"' para asegurar unicidad.

## Próximos temas

---

Estadística básica, bootstrap y visualización con Seaborn aplicadas a datos físicos.

¡Fin de la sesión!