Python for Data Engineering

BY: Myk Ogbinar

Data & Analytics Head @ San Miguel Corp

Founder @ Data Engineering Pilipinas



Workshop Expectations

- You can write basic Python apps
- X You don't need to know what a Data Engineer is yet



Today we'll:

- 1) Demystify what Data Engineers actually do
- 2) Build your first working data pipeline
- 3) Show how Python scripts become reliable data systems



Py-Data Engineering Workshop

From CSV to Dashboard — Building a Mini Data Pipeline in Pure Python pandas • pyarrow • uv • streamlit





Agenda

Segment	Duration

Intro + Setu	30 min	

Extract	10 min
---------	--------





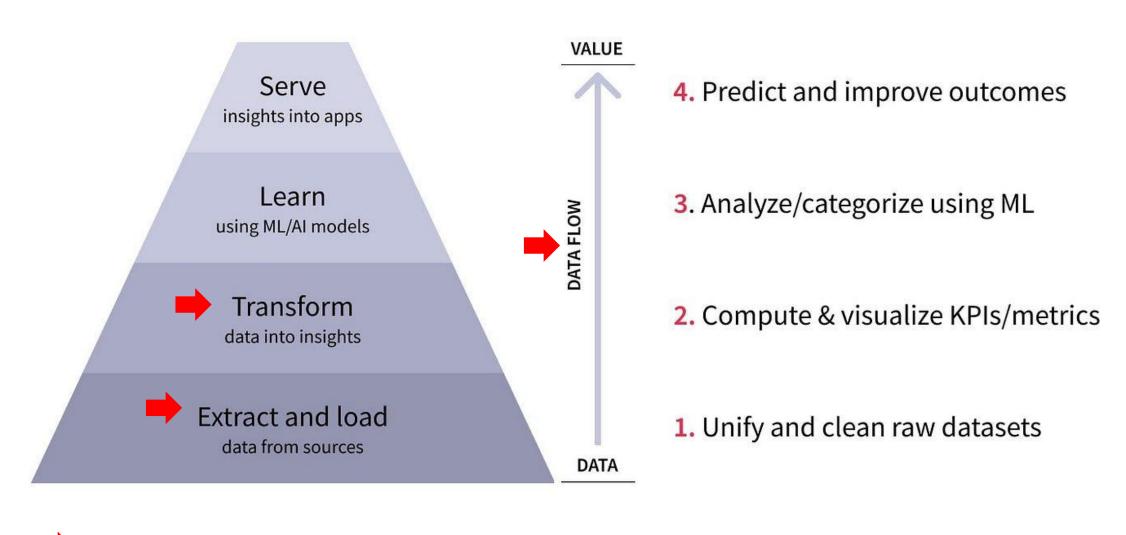
Introductio to Data Engineering

2025 PANAGHIUSA

What is Data Engineering?

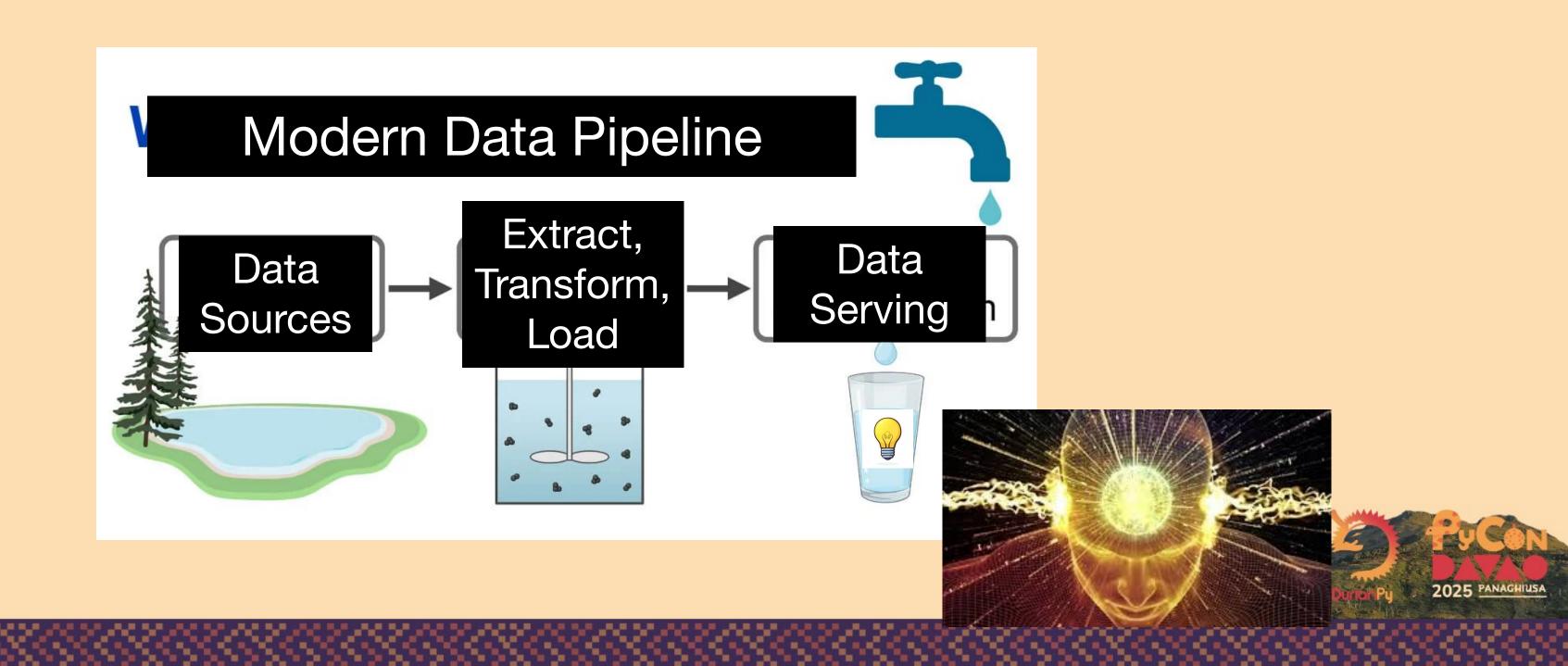


Hierarchy of needs for data driven companies



Data Engineering is tasked with building the Data Stack to support these needs

What is the DE Workflow?

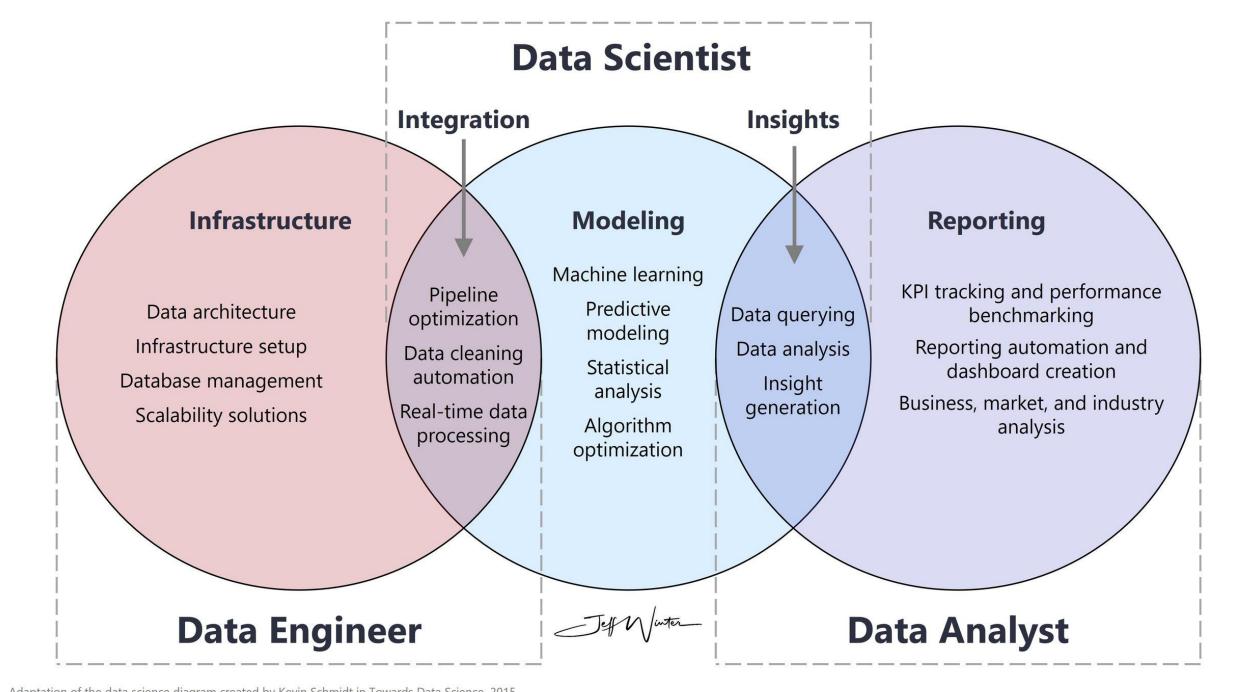


Why is Data Engineering Important?





Data Engineer vs Data Scientist vs Data Analyst





We're currently looking for a Data Migration Engineer based in the Philippines to join our growing team. Hiring! Data Engineer / Business Intelligence Developer (Power BI)

Key Requiremen

- Lead data onl
- Expand and ir
- Support creat
- Reconcile inventor
- Analyze data

Head of Data Engineering (Remote work)

- Metro Manila (Remote)
- Architects (Information & Communication Technology)
- (Full time
- @ P300,000 P390,000 per month

Posted 4d ago . Medium application volume

Compensation: \$1500 to \$1800 per month

Data Analyst / Data Engineer (cloud) (wfh) (day shift) 150-200K

Jobstreet

Data Analyst / Data Engineer (cloud) (wfh) (day shift) Job in Makat...

We're looking for both Data Engineer and Data Analyst that has years of experience and has experience in Cloud

- 🜠 Hiring: Data Engineers (Multiple positions available) for a global technology firm in Australia 🜠 Fully-Remote | Philippines
 - Base Salary: PHP 150-170k/month + sign-on bonus + statutory benefits

Data Modeling Techniques and Methodologies in Mandaluyong

Location:Quezon City

Salary: 80,000-110,000/month

Job Level: Experienced Hire

Senior Data Engineer in Taguig Salary: 85,000-108,000/month

Job Level: Experienced Hire

Salary:PHP80000 - PHP90000 per month + To be discussed

We're looking for a Lead Data Engineer to lead the design, development, and implementation of scalable data pipelines and infrastructure! If you're passionate about big data, cloud platforms, and leadership, this role is for you! 💥

🔑 Location: Ortigas, Pasig (Hybrid Work Setup 2x : 🎇 WE'RE HIRING: SENIOR DATA ENGINEER

Salary Range: Php 200,000 - 220,000 monthly

💅 Hiring: Lead Data Engineer 🚀

Location: QC Hybrid/ GY Hours Mon-Fri

Start Date: ASAP

Engagement Type: Long-Term Contract

Budget: Up to PHP 100K/month

🌌 We're Hiring! Data Engineer (Junior to Senior Level) 🜠



Location: UP Ayala, Quezon City (Hybrid)

Schedule: Day Shift

Salary: PHP 60,000 - 115,000

Hi Engineers! We are hiring for experienced Data Engineer.

Position: Director Data Engineering

WFH set-up

Salary: Php 300, 000



What is



What is Data?

- Originally, data is plural for "datum", a Latin word
- Datums are most often called "data points"
- Data represents a collection of data points (also called datasets)
- Can be used as both singular and plural form

Data are characteristics or <u>information</u>, usually numerical, that are collected through observation.

- Wikipedia



Shapes of Data

Structure Example Common Format

Tabular Sales table CSV / SQL

Hierarchical Product catalog JSON / XML

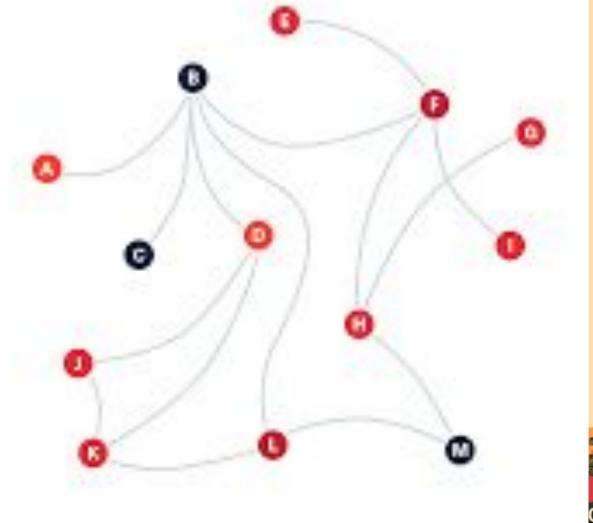
Graph Social network Neo4j, RDF

Time Series Sensor readings Parquet / Influx

Geospatial Map GeoJSON, coordinates Shapefile

D	F	G	H.	- 1
Sales resentative	Item	Price Per Unit (\$)	Units Sold	Total Sales (\$)
celia Blake	HJK_10	\$25	50	\$1,250
nie Freeman	HJK_12	\$20	100	\$2,000
vieve Evans	HJK_33	\$10	15	\$150
don Cohen	HJK_10	\$25	25	\$625
an Davis	HJK_33	\$10	110	\$1,100
an Grant	HJK_27	\$12	55	\$660
ry Herrera	HJK_12	\$20	120	\$2,400
egan Hicks	HJK_64	\$30	115	\$3,450
er Fleming	HJK_27	512	60	5720
ra Rogers	HJK 64	39		







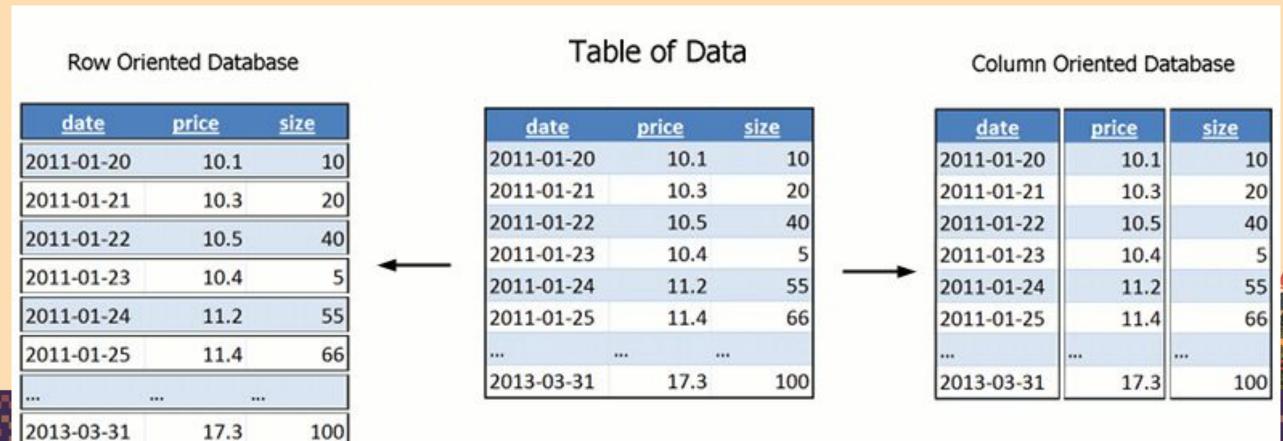
Column vs Row Storage

Row-Oriented DB

- > Stores/retrieves data by row → may read unnecessary columns
- > Fast for inserts/updates (OLTP)
- > Aggregations are expensive
- > Lower compression efficiency

Column-Oriented DB

- > Stores/retrieves data by column → reads only needed columns
- > Slower writes, optimized for reads (OLAP)
- > Aggregations are efficient
- > High compression due to similar column values





Special Values in Data & Data Quality

Value	Meaning	Example Sources	Data Cleansing Impact
NULL	Missing / unknown	Unfilled field	Standardize & decide: impute, ignore, or drop
NaN	"Not a Number" (invalid numeric result)	0/0, bad calculations	Convert to NULL or fix upstream logic
Infinity / - Infinity	Overflow or division by zero	1/0.0	Usually errors; filter, cap, or convert
Empty String ('')	String with length 0 (≠ NULL)	User form input	Standardize to NULL if it means missing
Sentinel Values	Fake placeholders	-999, 'NA'	Detect & replace with NULL

- Identify & standardize special values.
- Make missingness explicit (prefer NULL).
- Decide: Impute, Exclude, or Flag depending on context.



Data Modeling



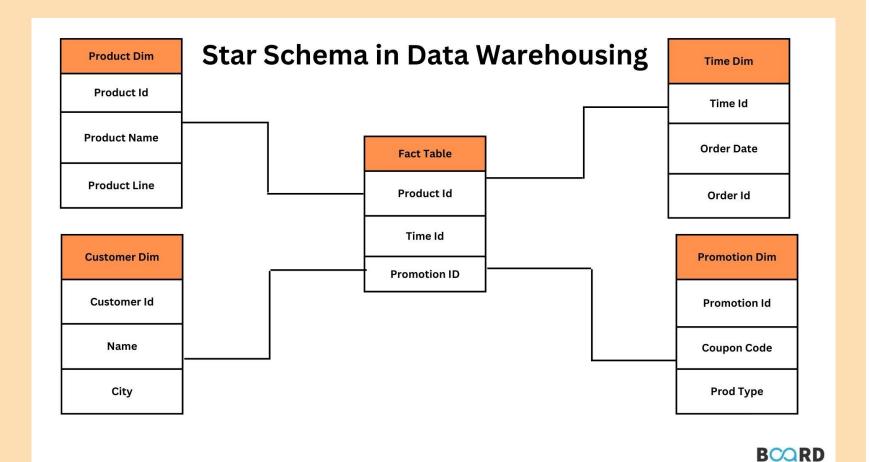
er_id	customer_name	product_name	price	
1	Mike Smith	Laptop	800	
2	Sara Jones	Mouse	20	,
3	Mike Smith	Laptop	800	
4	Sara Jones	Laptop	800	

Original unnormalized table

product_id	product_name	price
501	Laptop	800
502	Mouse	20

customer_id	customer_name	email
151	Mike Smith	mike@example.com
152	Sara Jones	sara@example.com

order_id	customer_id	order_date
1001	151	2024-08-16
1002	152	2024-08-19
1003	151	2024-08-22
1004	152	2024-08-25





Why Data Cleaning & Data Quality Matters

- Raw data = messy, inconsistent, incomplete
- Poor quality data → bad decisions
- Cleaning ensures accuracy, reliability, trust
- DE Mindset: Cleaning & Standardization happens in stages.

Getting the data "right"

1. Completeness

How complete is your dataset?

4. Validity

Does the data fit the defined range and definition?

2. Accuracy

Is the data an accurate representation of the element it describes?

5. Uniqueness

Are there any duplicates within your data?

3. Consistency

Is your data synchronized across your organization?

6. Integrity

Can your data be traced and connected across the organization?



Python for Data Engineering

Li	b	ra	ry

pandas

pyarrow

pathlib

streamlit

UV

Role

data manipulation

Parquet I/O

file system abstraction

lightweight visualization

dependency management

Concept

SQL + Excel hybrid

columnar efficiency

portability

serving layer

reproducibility



Let's get started!



Challenges of Data Engineering

- Inconsistent data formats
- Missing or invalid values
- Schema drift
- Scaling and performance
- Governance issues





How to Expand the Project

- Add more tables (Employees, Suppliers, Categories)
- Add new DQ rules
- Replace Parquet with DuckDB
- Schedule with Prefect
- Add new charts to Streamlit





Key Takeaways

- Data Engineering = making data usable
- Python can handle full workflow
- Focus on structure, quality, reproducibility





Closing Notes

- You didn't just learn Python you learned how to engineer data.
- Visit: dataengineering.ph
- Add: linkedin.com/in/ogbinar

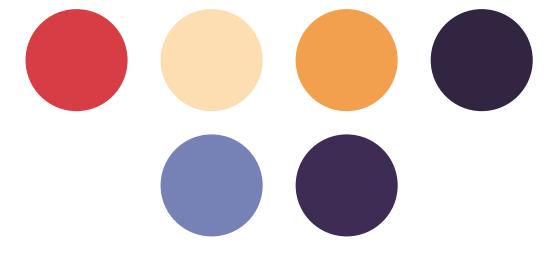




BRANDING

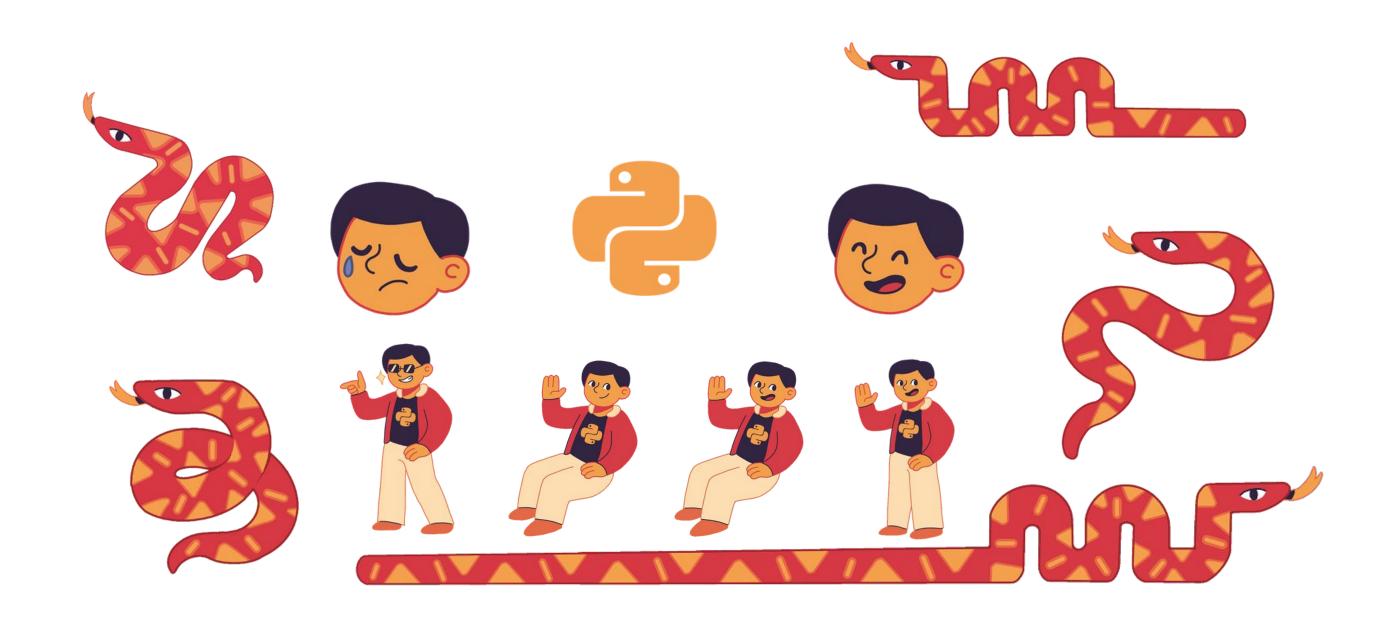






PIPANGANAN	H1	
Nunito	H1	
Haite Sans (H2	
(ဗြောင်လ)Sans (Bold) Inter (Medium)		Body
Nunito Sans (Semi-Bold)	Body

ASSETS



TITLE HERE

Insert Text

```
fflush(stdin);
scanf("%[^\n]",R.note);
if(fwrite(&R,sizeof(R),1,fp)){
226
227
228
              gotoxy(5,12);
229
              puts("Note is saved sucessfully");
230
              fclose(fp);
231
232
               gotoxy(5,12);
233
               SetColor(12);
              puts("\aFail to save!!\a");
234
              ClearColor();
235
236
         gotoxy(5,15);
         printf("Press any key....");
                                                   getch();
237
238
         getch();
239
         fclose(fp);
240
```

2025 PANAGHIUSA