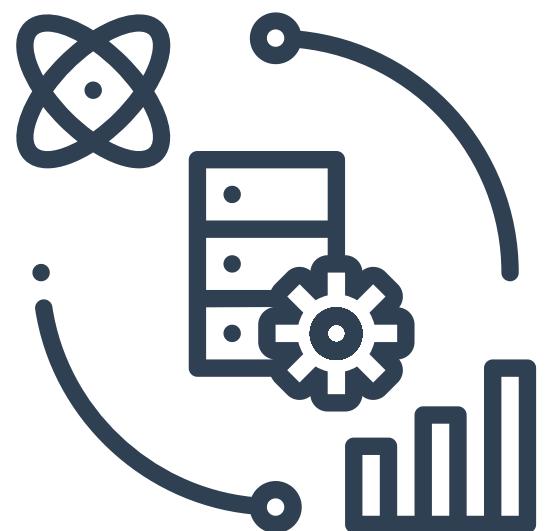


# Data Science Review

لا تنسونا من دعواتكم

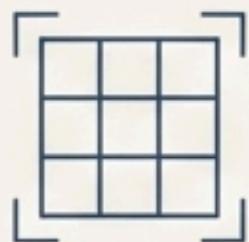
الله يوفقنا جميعاً ونتقابل في الصيفية 🙏  
هذا حسابنا لو عندكم اي اقتراحات او ملفات ممكن ننشرها :  
[Ice Sara Muath](#)



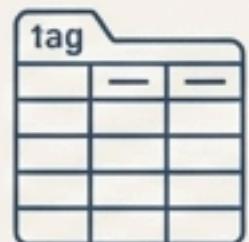
# The Data Science Workflow is a Two-Act Play

## Act I: Forging the Toolkit

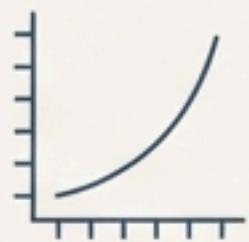
We'll first stock our conceptual toolbox with Python's most powerful data science libraries. The goal is to understand what each tool does and why it matters.



**NumPy:** The foundation for high-performance numerical computing.



**Pandas:** The indispensable tool for manipulating and analyzing structured data.



**Matplotlib:** The essential library for bringing data to life through visualization.

## Act II: The Workshop

Next, we'll apply these tools to a real challenge: cleaning and preparing the famous Titanic dataset for a machine learning model. This is where theory meets practice.

### Key Takeaway

The goal isn't to memorize every function. It's to understand the purpose of the tools, so you know which one to grab when you face a problem.

# NumPy: The Bedrock of Numerical Computing in Python

At its heart, NumPy provides the `ndarray` (n-dimensional array) object. It's like a Python list but supercharged: faster, more memory-efficient, and equipped with a vast library of mathematical functions. Pandas and many other libraries are built directly on top of it.

```
import numpy as np

# From a Python list
arr_1d = np.array([1, 2, 3, 4, 5])

# A 2x3 matrix of zeros
zeros_matrix = np.zeros((2, 3))

# A range of values, like Python's range()
even_numbers = np.arange(0, 10, 2)

# 5 evenly spaced numbers between 0 and 1
linear_space = np.linspace(0, 1, 5)
```

- `shape` is (rows, columns)
- Works like `np.arange(start, stop, step)`
- Extremely useful for plotting functions

## Why it Matters

Nearly every numerical operation in the data science ecosystem relies on NumPy's speed and efficiency. Understanding its arrays is non-negotiable.

# Pandas: Your Programmable Spreadsheet for Data Analysis

Pandas introduces two essential data structures:

1. **Series**: A one-dimensional labeled array (like a single column in a spreadsheet).
2. **DataFrame**: A two-dimensional labeled structure with columns of potentially different types (the entire spreadsheet).



pd.read\_csv()



CSV File

	Name	Age	City
0	Alice	25	New York
1	Bob	30	San Francisco
2	Charlie	22	London

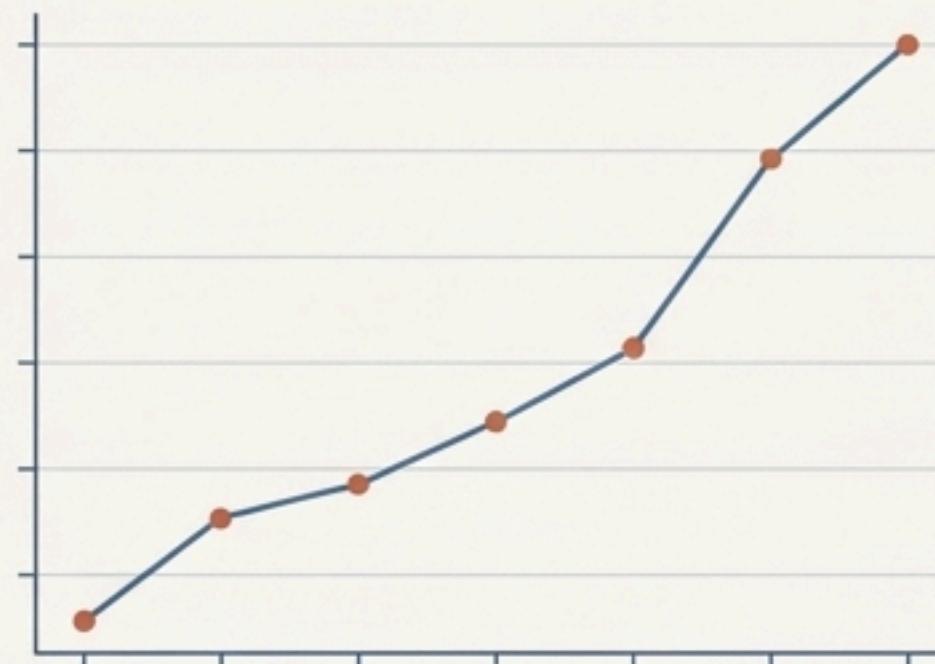
```
import pandas as pd

# Reads a CSV file from a path into a DataFrame
df = pd.read_csv('path/to/your/data.csv')

# Display the first 5 rows to inspect the data
df.head()
```

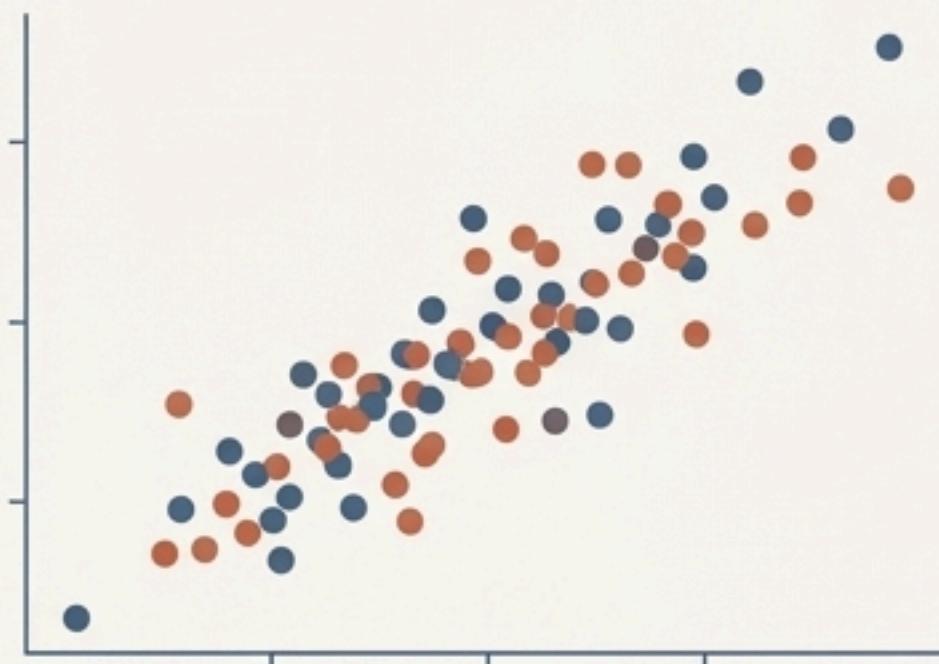
# Matplotlib: Visualizing Data to Uncover Insights

A picture is worth a thousand rows of data. Matplotlib is the workhorse library for creating static, animated, and interactive visualizations in Python. Understanding a few basic plot types is enough to start exploring your data effectively.



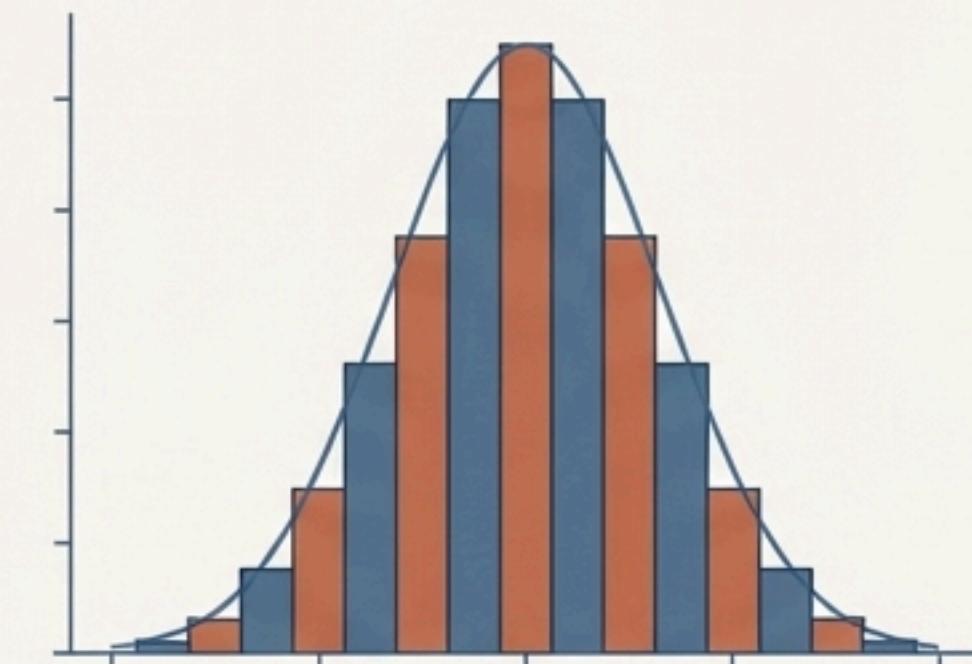
Ideal for showing trends over time.

```
plt.plot(x, y)
```



Perfect for visualizing the relationship between two variables.

```
plt.scatter(x, y)
```



Essential for understanding the distribution of a single variable.

```
plt.hist(data)
```

# The Workshop: Preparing the Titanic Dataset

## The Challenge:

We have a raw dataset containing information about passengers on the Titanic. Our goal is to prepare this data for a machine learning model that can predict whether a passenger survived.

## The Process:

Raw data is almost never model-ready. It's often messy, with missing values, text, and irrelevant information. We will now use our "toolbox" (NumPy, Pandas, Matplotlib) to systematically clean and transform the data, following the same process experts use.

```
>>> df.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	11	0	1	Braund, Mr. Owen Harris	male	NaN	0	0	1921926	13.00	NaN	E
1	12	1	2	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	NaN	1	0	2888753	80.38	NaN	Y
2	15	1	3	Heikkinen, Miss. Laina	female	NaN	0	0	1322	7.40	NaN	F
3	16	0	2	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	NaN	1	1	TICKET4	254.00	NaN	S
4	17	0	2	Allen, Mr. William Henry	male	NaN	0	0	20718	1.00	NaN	A

# Step 1: Exploratory Data Analysis (EDA) - The First Diagnosis

Before we change anything, we must understand our data. The `.info()` method in Pandas provides a concise summary of the DataFrame, revealing its structure and potential issues.

```
# Get a summary of the dataset  
df_titanic.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 891 entries, 0 to 890  
Data columns (total 12 columns):  
 #   Column      Non-Null Count Dtype  
 ---  -----  
 0   PassengerId 891 non-null   int64  
 1   Survived     891 non-null   int64  
 2   Pclass       891 non-null   int64  
 3   Name         891 non-null   object  
 4   Sex          891 non-null   object  
 5   Age          714 non-null   float64  
 6   SibSp        891 non-null   int64  
 7   Parch        891 non-null   int64  
 8   Ticket       891 non-null   object  
 9   Fare          891 non-null   float64  
 10  Cabin         204 non-null   object  
 11  Embarked     889 non-null   object
```

**\*Missing Values\***: We have 891 rows, but not 891 values here.

**\*Incorrect Data Types\***: Models can't process 'object' (text) data. This needs to be converted to numbers.

# Step 2: Imputation - Handling Missing Values

## The Strategy:

We can't feed a model data with gaps ('NaNs'). Our approach depends on the column:

- **Age (Numerical):** We will fill missing ages with the *median* age of all passengers.
- **Embarked (Categorical):** We'll fill the two missing values with the *mode* (the most common port of embarkation).
- **Cabin (Too Many Missing):** With over 77% of its values missing, this column is unsalvageable. The best course of action is to drop it entirely.

### Age Imputation (Median)

```
// Fill 'Age' with the median  
median_age = df_titanic['Age'].median()  
df_titanic['Age'].fillna(median_age, inplace=True)
```

### Embarked Imputation (Mode)

```
// Fill 'Embarked' with the mode  
mode_embarked = df_titanic['Embarked'].mode()  
df_titanic['Embarked'].fillna(mode_embarked, inplace=True)
```

### Embarked Imputation (Mode)

```
// Fill the 'Cabin' column  
df_titanic.drop('Cabin', axis=1, inplace=True)
```

### Cabin Removal (Drop)

```
// Drop the 'Cabin' column  
df_titanic.drop('Cabin', axis=1, inplace=True)
```

### Expert Insight

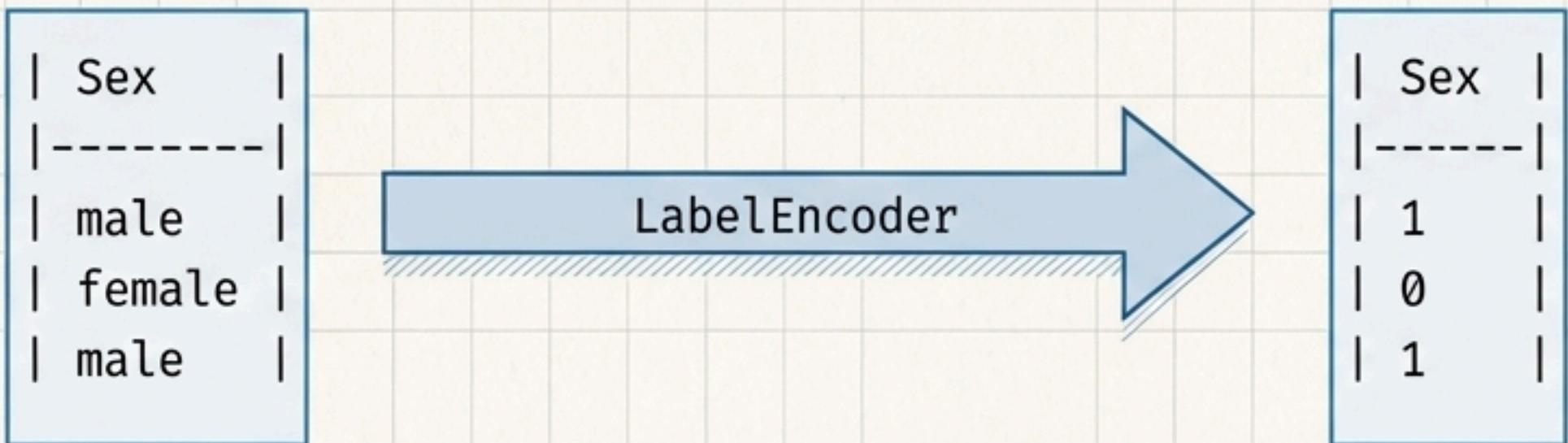


#### Why Median, Not Mean?

The 'Age' column might have outliers (a few very old passengers). The mean is sensitive to outliers, while the median is more robust, making it a safer choice for imputation here.

# Step 3: Encoding - Converting Text to Numbers

Machine learning models are mathematical; they operate on numbers, not text like 'male' or 'female'. We must *\*encode\** these categorical features into a numerical format.



```
from sklearn.preprocessing import LabelEncoder  
  
le = LabelEncoder()  
df_titanic['Sex'] = le.fit_transform(df_titanic['Sex'])  
df_titanic['Embarked'] = le.fit_transform(df_titanic['Embarked'])
```

## A Word of Caution on Label Encoding

LabelEncoder assigns 0, 1, 2... This can imply an ordinal relationship (2 > 1 > 0) where none exists (e.g., ports of embarkation). For features with no natural order, 'OneHotEncoder' is often a better choice, but 'LabelEncoder' is a common starting point.

# Step 4: Feature Engineering & Selection - The Human Element

**The Idea:** Not all data is useful. Some columns provide no predictive value, while others can be engineered to create more powerful signals. This step requires critical thinking about what each feature represents.

## Our Decisions:

- Drop `PassengerId`: A unique ID number is just a random label; it has no relationship with survival.
- Drop `Name`: While unique to each person, the name itself is too complex for a simple model.
- Drop `Ticket`: The ticket numbers are mostly unique and don't contain a consistent, useful pattern.

```
# Drop columns that are not useful for prediction  
df_titanic.drop(['PassengerId', 'Name', 'Ticket'], axis=1, inplace=True)
```

## Expert Insight:



### The Hidden Signal in 'Name'

We dropped 'Name' for simplicity, but an expert might extract titles ('Mr.', 'Mrs.', 'Dr.') from it. A passenger's title could be a powerful proxy for age, gender, and social status. This is a classic example of feature engineering—creating new features from existing ones.

# Step 5: Feature Scaling - Leveling the Playing Field

## The Problem

Our numerical features, like "Age" (0-80) and "Fare" (0-512), exist on vastly different scales. Many machine learning algorithms are sensitive to this and will incorrectly give more weight to features with larger values.

## The Solution

We use a scaler, like `StandardScaler` from scikit-learn, to transform our data. This method rescales features to have a mean of 0 and a standard deviation of 1, without changing the shape of their distribution.



# The Transformation: From Raw Data to Model-Ready Matrix

BEFORE

df_titanic.head()												
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	18	0	1	Hrr. Rerzin Nanxan Delhiness, ...	females	NaN	0	0	12312883	235.3468	NaN	NaN
1	17	1	1	Rlchano Wanta Pasmenoy	females	22.8	0	1	1234536	3.8328	NaN	L
3	18	0	1	Nama Slgak Sertewan Name	females	26.8	1	1	2049593	226.6299		N
4	19	0	1	Bate Jerey Person	females	33.8	1	1	3905651	158.6388		Q
5	18	1	1	Mr. Berice Venadsn Clddy	females	24.6	1	1	2359852	72.6730	NaN	E
6	28	1	1	Make Shave Agoency	females	26.8	0	0	22511952	78.3899	NaN	S
7	28	2	1	Nellirky Alexanden Boaste	female	NaN	0	0	2899853	229.8862	NaN	NaN

AFTER

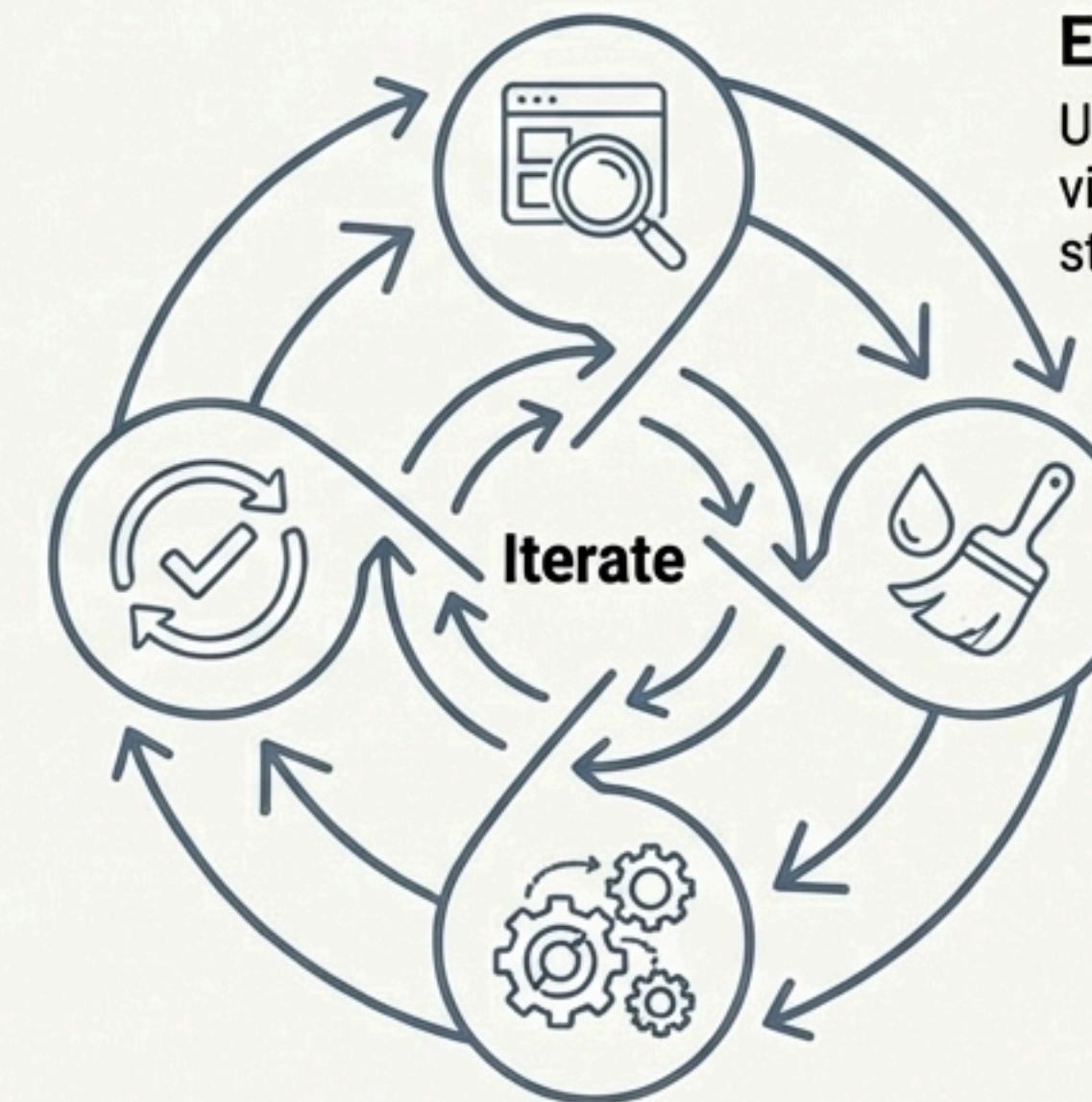
Survived	Pclass	Age_scaled	SibSp_scaled	Parch_scaled	Fare_scaled	Sex_male	Embs_frced_Q	Embs_frced_S
0	3	-0.54	-0.30	-0.46	-0.54	0	1	1
1	2	1.23	0.36	0.31	1.23	1	0	0
0	3	-0.58	0.43	-0.70	-1.03	0	1	1
0	3	-1.10	-1.26	-0.24	-1.20	1	0	0
0	2	0.83	-0.33	0.31	2.29	0	1	1
1	3	1.23	-0.82	0.35	1.23	1	0	1
0	4	-0.58	-0.36	-1.30	-0.56	1	1	0
0	5	-0.54	-0.61	0.01	-0.63	0	1	1

- No Missing Values
- Fully Numerical
- Scaled Features
- Relevant Features Only

The result of our workshop is a clean, structured dataset ready to be fed into a machine learning model. Every step was a deliberate choice to improve the quality of our data.

# The Data Scientist's Method: A Repeatable Workflow

The process we followed wasn't just about the Titanic; it's a fundamental workflow for nearly any data science project. It's an iterative cycle of questioning and refining your data.



## Validate & Iterate

(Looks ahead) This prepared data is now ready for modeling. The model's performance will tell you if you need to revisit earlier steps.

## Explore & Diagnose

Use tools like `.info()`, `.describe()`, and visualizations to understand the data's strengths and weaknesses.

## Clean & Impute

Address structural problems like missing values and duplicates.

## Transform & Engineer

Convert data into a usable format (encoding, scaling) and create new, more powerful features.

**Mastering the tools is the start. Mastering the *process* is what turns data into discovery.**