# Data Cleaning Process in Pandas

**Before Cleaning Data, We Need to Understand the Data and Their Type and Try Identify the Issue/Assessment Manually and Programmatic**

**There Are 2 Type of Assessments**

- Manual - Looking Through the Data Manually in Excel or Google Sheet.
- Programmatic - By Using Pandas Function Such as Info (), Describe () Or Sample ()

**Manual:**

- **Review Your Dataset Carefully and Note Down Problems You Have Identify**
- **Check Number of Row and Columns.**
- **Any Columns Has Any Missing Values.**
- **Misspelling In String Columns**
- **Check Accuracy of Data Like Values in All Column Are Accurate**
- **Valid Data or Not (**New York -> Ny**)**
- **Any Negative Values in Any Columns Like (Weight, Height)**
- **Column Not Contains Multiple Value Which Is Inaccurate Format**
- **Check Missing Data, Etc.**

**Automatic/Programmatic Assessment**

- Head & Tail
- Sample
- Describe
- Info
- Isnull
- Duplicated
- Data type

**There Are 2 Step Involving in Assessment:**

- **Discover** - find problem
- **Document** - write problem to remember it while cleaning data

**Once you discover the problem and create documentation now we need to labialised the problem to process the data cleaning as per the created path to resolve the issue in the data set for analysis. This process will help to clean data fast and safely.**

**There are 4 option which is help to sort your problem to clean your data to improve the data quality and arrange your data in Dirty Data or Messy Data.**

1. **Dirty Data**
   - (Error, Missing Values, Nan Values, Col Formatting, String Formatting, Wrong Word, Misspelled Word)
2. **Messy Data**
   - (Structure Issue, Single Column Contain Multiple Values)

# My manual Assessment and labelling

**Data Quality Dimensions:**

1. Completeness - Is Data Missing?
2. Validity - Is Data Invalid (Negative Height, Duplicate Patients)
3. Accuracy - Data Is Valid But Not Accurate (Weight - 1kg )
4. Consistency - Both Valid & Accurate But Written Differently ( New York -> Ny )

**Order Of Severity (Who Is Dangerous Problem In Among All 4 Data Quality Dimensions)**

| Completene | ← | Validity | ← | Accuracy | ← | Consistency |

**After assigning label to your problems now follow the below process to clean your data in recommended Data cleaning orders as per the label given to your problems.**

## Data cleaning order

1. Quality - Completeness
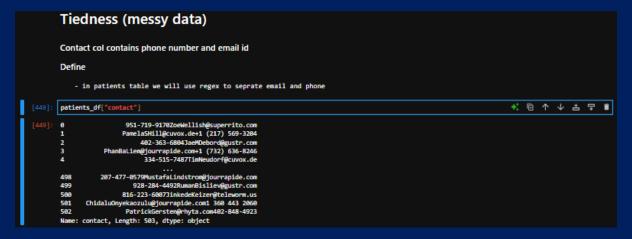2. Tidiness
3. Quality - Validity
4. Quality - Accuracy
5. Quality - consistency

## Step involve in data cleaning

- Define - solution of problem you have found
- Code - write code of the solution
- Test - check code is working to solve the problem.

**Always make sure to create a copy your pandas data frame before starts the cleaning process ***

**Define**: Write how you solve the problem

**Code**: write your code

```python
# code

import re
import numpy as np

def find_contact_details(text: str) -> tuple:
    if not text or text.strip() == "":
        return np.nan, np.nan

    # Phone number pattern: supports international (+xx) and US formats (xxx-xxx-xxxx)
    phone_number_pattern = re.compile(r"(\+[\d]{1,3}\s)?(\(?[\d]{3}\)?\s?-?[\d]{3}\s?-?[\d]{4})")
    phone_number_match = re.search(phone_number_pattern, text)

    if phone_number_match:
        phone_number = phone_number_match.group()
        # Remove phone number from the text before searching for the email
        text = re.sub(phone_number_pattern, '', text).strip()
    else:
        phone_number = np.nan

    # Email pattern to match most valid email addresses
    email_pattern = re.compile(r'[a-zA-Z0-9_.+-]+@[a-zA-Z0-9-]+\.[a-zA-Z0-9-.]+')
    email_address_match = re.search(email_pattern, text)

    if email_address_match:
        email_address = email_address_match.group()
    else:
        email_address = np.nan

    return phone_number, email_address


patients_df["Phone"]=patients_df["contact"].apply(lambda x: find_contact_details(x)).apply(lambda x: x[0])
patients_df["email"]=patients_df["contact"].apply(lambda x: find_contact_details(x)).apply(lambda x: x[1])
```

**Test**: test your output is worked or not as per the define solution

```python
# test
patients_df.sample()
```

| | patient_id | assigned_sex | given_name | surname | address | city | state | zip_code | country | birthdate | weight | height | bmi | Phone | email |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 229 | 230 | male | John | Doe | 123 Main Street | New York | NY | 12345.0 | United States | 1/1/1975 | 180.0 | 72 | 24.4 | 1234567890 | johndoe@email.com |

# Thanks You