

We'll do this project in a **Jupyter Notebook**, which is perfect for this kind of step-by-step analysis and documentation. If you have Anaconda installed, you can easily launch a Jupyter Notebook from there.

Let's begin with **Step 1: Data Loading and Initial Exploration**.

The goal here is to load the dataset and get a first impression of its structure, content, and overall health. I'll write the Python code to do this. You can follow along and run this code in your Jupyter Notebook.

Here is the plan for this phase:

1. **Import pandas:** The essential library for data manipulation in Python.
2. **Load the CSV:** Read your `employee_wellness_dataset.csv` file into a pandas DataFrame.
3. **Initial Inspection:** We will look at the first few rows, check the data types of each column, and see a statistical summary of the numerical columns.

Let's execute the code for this.

```
Python
import pandas as pd

# Load the dataset
df = pd.read_csv('employee_wellness_dataset.csv')

# Display the first 5 rows of the dataframe
print("First 5 rows of the dataset:")
print(df.head())

# Display the summary of the dataframe
print("\nDataFrame Info:")
df.info()

# Display the descriptive statistics of the dataframe
print("\nDescriptive Statistics:")
print(df.describe())

# Check for missing values
print("\nMissing Values Count:")
print(df.isnull().sum())
```

Code output

First 5 rows of the dataset:

	S.No	Timestamp	Age	Gender	Country	state	self_employed	family_history	treatment	work_interfere	no_employees	remote_work	tech_company	benefits	care_options	wellness_program	seek_help	anonymity	leave	mental_health_consequence	phys_health_consequence	coworkers	supervisor	mental_health_interview	phys_health_interview	mental_vs_physical	obs_consequence	comments	
0	1	8/27/2014 11:29	37	Female	United States	IL	NaN	No	Yes	Often	25-Jun	No	Yes	Not sure	No	Yes	Yes	Somewhat easy	No	Maybe	Yes	No	NaN	No	Some of them	Yes	No	Maybe	Yes
1	2	8/27/2014 11:29	44	M	United States	IN	NaN	No	No	Rarely	More than 1000	No	No	No	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know
2	3	8/27/2014 11:29	32	Male	Canada	NaN	NaN	No	No	Rarely	25-Jun	No	Yes	No	No	Don't know	Somewhat difficult	Yes	No	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes
3	4	8/27/2014 11:29	31	Male	United Kingdom	NaN	NaN	Yes	Yes	Often	26-100	No	Yes	Yes	No	No	Somewhat difficult	Yes	No	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes
4	5	8/27/2014 11:30	31	Male	United States	TX	NaN	No	No	Never	100-500	Yes	Yes	No	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know	Don't know

DataFrame Info:

```
<class 'pandas.core.frame.DataFrame'>
```

RangeIndex: 1048 entries, 0 to 1047

Data columns (total 28 columns):

#	Column	Non-Null Count	Dtype
0	S.No	1048 non-null	int64
1	Timestamp	1048 non-null	object
2	Age	1048 non-null	int64
3	Gender	1048 non-null	object
4	Country	1048 non-null	object
5	state	636 non-null	object
6	self_employed	1030 non-null	object
7	family_history	1048 non-null	object
8	treatment	1048 non-null	object

9	work_interfere	812	non-null	object
10	no_employees	1048	non-null	object
11	remote_work	1048	non-null	object
12	tech_company	1048	non-null	object
13	benefits	1035	non-null	object
14	care_options	1048	non-null	object
15	wellness_program	1044	non-null	object
16	seek_help	1048	non-null	object
17	anonymity	1048	non-null	object
18	leave	1044	non-null	object
19	mental_health_consequence	1048	non-null	object
20	phys_health_consequence	1048	non-null	object
21	coworkers	1048	non-null	object
22	supervisor	1048	non-null	object
23	mental_health_interview	1048	non-null	object
24	phys_health_interview	1048	non-null	object
25	mental_vs_physical	1048	non-null	object
26	obs_consequence	1048	non-null	object
27	comments	132	non-null	object

dtypes: int64(2), object(26)

memory usage: 229.4+ KB

Descriptive Statistics:

	S.No	Age
count	1048.000000	1.048000e+03
mean	524.500000	9.541988e+07
std	302.675844	3.089010e+09
min	1.000000	-1.726000e+03
25%	262.750000	2.700000e+01
50%	524.500000	3.100000e+01
75%	786.250000	3.600000e+01
max	1048.000000	1.000000e+11

Missing Values Count:

S.No	0
Timestamp	0
Age	0
Gender	0
Country	0
state	412
self_employed	18
family_history	0
treatment	0
work_interfere	236
no_employees	0
remote_work	0
tech_company	0
benefits	13

```

care_options          0
wellness_program      4
seek_help             0
anonymity             0
leave                 4
mental_health_consequence 0
phys_health_consequence 0
coworkers             0
supervisor           0
mental_health_interview 0
phys_health_interview 0
mental_vs_physical    0
obs_consequence       0
comments              916
dtype: int64

```

Initial Findings and Observations

Excellent! We've successfully loaded the data. Here's what we can see so far:

- **Structure:** The dataset has **1048 rows** and **28 columns**. This is a good amount of data to work with.
- **Data Types:** Most columns are of type `object`, which usually means they contain text or categorical data. `Age` and `S.No` are integers, as expected. `Timestamp` is currently an `object`, and we should convert it to a `datetime` format later.
- **Missing Values:** We have several columns with missing values:
 - `comments` has a very high number of missing values (916), which is expected for a comments field. We'll likely drop this column.
 - `state` is missing 412 values. This makes sense as the question was only for US employees.
 - `work_interfere` has 236 missing values.
 - `self_employed`, `benefits`, `wellness_program`, and `leave` have a small number of missing values.
- **Age Column Issues:** The `Descriptive Statistics` for the `Age` column look strange. The `min` age is -1726 and the `max` is a very large number. This indicates there are invalid entries in the 'Age' column that we need to clean.
- **Gender Column:** From the `head()` output, we see 'Female', 'M', and 'Male'. We'll need to standardize this column to have consistent values (e.g., 'Female' and 'Male').

Next Steps: Phase 2 - Data Cleaning

Based on our initial findings, our next immediate task is to clean the data. This is a critical step to ensure our analysis is accurate. Here's the plan:

1. **Clean the Age column:** Remove or correct the invalid age values.
2. **Standardize the Gender column:** Consolidate the different representations of gender into a uniform format.
3. **Handle Missing Values:** We'll decide on a strategy for columns like `self_employed`, `work_interfere`, and `state`.
4. **Drop Unnecessary Columns:** We will remove columns that are not useful for our analysis, like `S.No` and `comments`.
5. **Convert Timestamp:** Change the `Timestamp` column to a datetime object, which might be useful for time-based analysis.