



Essential Guide to Clean Data

Clean your data, clarify your insights: practical techniques for reliable analysis.

A hands-on guide to cleaning and validating data with Python for reliable analysis.

Ibon Martínez-Arranz | imartinez@labrubiocom

Data Science Manager at Rubió Metabolomics

www.rubiometabolomics.com

Ibon Martínez-Arranz, PhD in Mathematics and Statistics, holds MSc degrees in Applied Statistical Techniques and Mathematical Modeling. He has extensive expertise in statistical modeling and advanced data analysis techniques. Since 2017, he has led the Data Science area at Rubió Metabolomics, driving predictive model development and statistical computing management for metabolomics, data handling, and R&D projects. His doctoral research in Mathematics investigated and adapted genetic algorithms to improve the classification of NAFLD subtypes.

Itziar Mincholé Canals | iminchole@labrubio.com

Data Specialist at Rubió Metabolomics

www.rubiometabolomics.com

Itziar Mincholé Canals has a degree in physical sciences (1999) from the University of Zaragoza (Spain). In 1999 she began developing her professional career as a software analyst and developer, mainly working on web development projects and database programming in different sectors. In 2009 she joined OWL Metabolomics as a bioinformatician, where she has participated in several R&D projects and has developed several software tools for metabolomics. In 2016 she completed the master's degree in applied statistics with R software. After joining the Data Science department, she also supports data analysis.



Essential Guide to Clean Data

Ibon Martínez-Arranz



**Life
Feels
Good**



rubió
Metabolomics

Contents

Essential Guide to Clean Data	1
Introduction — Why Clean Data Matters	3
What Do We Mean by “Clean Data”?	3
Dirty Data in the Real World	4
The Hidden Cost of Not Cleaning	4
Data Cleaning and ALCOA+	5
This Book’s Approach	5
References	5
Loading and Inspecting Raw Data	7
Load it right	7
First visual inspection	8
Common surprises	8
Save a copy	9
Summary	9
Handling Missing Values	11
Why missingness matters	11
Types of missing data	11
Detecting missing values in pandas	12
Strategies to handle missing values	13
Drop them (but be careful)	13
Impute them (simple methods)	13
Impute them (fancy methods)	13

Flag them	14
Real-world example: BMI in a clinical dataset	14
Documentation is key	15
Summary	15
Fixing Data Types and Formats	17
Why it matters	17
Checking types in pandas	18
Fixing numeric columns	18
Converting to datetime	18
Handling categorical data	19
Boolean conversions	20
Dealing with unit conversions	20
Case study: fixing a lab dataset	20
Summary	21
Standardizing and Cleaning Text Data	23
Why text cleaning matters	23
Common problems in text data	24
Basic string cleaning in pandas	24
Replacing variants and typos	24
Dealing with missing or meaningless text	25
Categorical cleanup example	25
Unicode and accent normalization	26
Summary	26
Detecting and Handling Duplicates	29
What is a duplicate?	29
Visual inspection	30
Handling duplicates safely	30
1. Drop exact duplicates	30
2. Drop by key columns	30

3. Aggregate near-duplicates	31
4. Use fuzzy matching (advanced)	31
Logging duplicate removal	31
Case study: clinical visit log	32
Summary	32
Outlier Detection and Correction	33
What is an outlier?	33
Visualizing outliers	33
Boxplots	33
Histograms	34
Scatter plots	34
Z-scores	34
IQR method	34
What to do with outliers	35
1. Investigate	35
2. Correct if clearly wrong	35
3. Cap or clip	35
4. Flag them	36
5. Remove them (last resort)	36
Real-world example: metabolic data	36
Summary	37
Resolving Inconsistencies Across Columns	39
Common types of column inconsistencies	39
Logical checks with dates	39
Recalculating from raw data	40
Contradictory flags	40
Category mismatch	41
Value correlation	41
Case study: verifying timestamps	42
Summary	42

Validating and Enforcing Constraints	43
Why validation matters	43
Types of validations	43
Range checks	43
Allowed values (categorical checks)	44
Null constraints	44
Cross-field dependencies	44
Pattern checks with regex	44
Validating with <code>pandera</code>	45
Using <code>pydantic</code> for row-level validation	45
Referential integrity	46
Constraint enforcement with assertions	46
Case study: validating a clinical dataset	46
Logging and documenting rules	47
Summary	47
Logging, Versioning, and Audit Trails	49
Why traceability matters	49
Keeping raw and cleaned data separate	50
Logging with the <code>logging</code> module	50
Versioning data with DVC	51
Recording code versions	51
Jupyter notebooks and cell history	51
Reproducible scripts	52
Case study: versioned cleaning pipeline	52
Summary	53
Case Studies in Data Cleaning	55
Case 1: Cleaning Clinical EHR Data	55
Problems spotted:	55
Cleaning steps:	56
Results:	56

Case 2: Public Health Indicators (WHO)	57
Problems spotted:	57
Cleaning steps:	57
Results:	58
Case 3: LC-MS Metabolomics Data	58
Issues found:	58
Cleaning actions:	59
Outcome:	59
Summary	59
References and Resources	61
Key References	61
Datasets Used in Case Studies	62
Clinical EHR Dataset (Synthetic Example)	62
WHO Health Indicators Dataset	62
Metabolomics LC-MS Dataset	62
Tools Used Throughout	62

Essential Guide to Clean Data

Introduction — Why Clean Data Matters

Before we even think about training a model, generating a plot, or calculating a correlation, there's one thing that every data scientist, analyst, or researcher needs to do: clean the data.

Cleaning data might not be the most glamorous part of data science, but it is undoubtedly one of the most important. It's the foundation on which everything else is built. No matter how advanced your algorithms are or how elegant your visualizations look, if your data is full of errors, inconsistencies, and missing values, your results will be misleading at best—and completely wrong at worst.

What Do We Mean by “Clean Data”?

“Clean data” is a bit of a vague term, but broadly speaking, it refers to data that:

- Has consistent formatting and structure
- Uses the correct data types (e.g., dates are stored as dates, numbers as numbers)
- Has no missing or malformed values—or at least, missing data has been handled appropriately
- Doesn't contain obvious errors, duplicates, or outliers (unless justified)
- Matches across sources if needed (e.g., foreign keys between tables)

- Is well-documented and reproducible

Clean data isn't perfect—it's realistic, usable, and trustworthy. In short, clean data is **fit for purpose**.

Dirty Data in the Real World

Unfortunately, most real-world datasets are messy. Very messy. Columns might be misnamed, data might be missing without explanation, units might be inconsistent, and errors can sneak in during data entry, export, or transformation. In clinical datasets, it's common to see blood pressure measured in mmHg mixed with values that look suspiciously like they came from a different scale—or simply typos.

This messiness isn't just annoying—it's dangerous. As noted by Karr et al. (2006), poor data quality can invalidate scientific results, introduce bias, and waste enormous amounts of time and resources. Especially in healthcare, misleading conclusions from dirty data can have real-world consequences.

The Hidden Cost of Not Cleaning

Skipping the cleaning process—or doing it half-heartedly—has cascading effects:

- **Garbage in, garbage out:** Bad input leads to bad output. Your model may appear to perform well, but it's learning from flawed assumptions.
- **Reproducibility issues:** Without clear cleaning steps, it becomes hard to reproduce results.
- **Loss of trust:** Colleagues and stakeholders lose confidence in your analysis if inconsistencies pop up.
- **Compliance risks:** In regulated environments like clinical trials, dirty data violates standards such as GCP (Good Clinical Practice) and ALCOA+ principles.

Data Cleaning and ALCOA+

ALCOA+ stands for:

- **Attributable:** It should be clear who recorded the data.
- **Legible:** The data should be readable and permanent.
- **Contemporaneous:** Recorded at the time of the event.
- **Original:** The data must be the source or a certified copy.
- **Accurate:** Free from error.

The “+” includes: Complete, Consistent, Enduring, and Available. Data cleaning isn’t just about technical correctness—it’s about meeting these expectations too.

This Book’s Approach

This book is hands-on. We’re going to walk through real problems and clean them using Python—mostly with pandas, but also with other useful tools. Each chapter tackles a typical problem:

- What to do with missing data?
- How do I handle weird date formats?
- Are outliers always wrong?
- How do I track what changes I made to the data?

We’ll keep things practical, reproducible, and honest. You’ll get code examples, visualizations, and plenty of tips based on real-world experience.

References

- Karr, A. F., Sanil, A. P., & Banks, D. L. (2006). Data quality: A statistical perspective. *Statistical Methodology*, 3(2), 137–173. <https://doi.org/10.1016/j.stamet>

.2005.08.005

- EMA (2010). Reflection paper on expectations for electronic source data and data transcribed to electronic data collection tools in clinical trials. *European Medicines Agency*. <https://www.ema.europa.eu>
- FDA (2018). *Data Integrity and Compliance With Drug CGMP: Questions and Answers*. <https://www.fda.gov/media/119267/download>

Loading and Inspecting Raw Data

Once you get your hands on a dataset, the first instinct is often to run a model or generate a plot. But before doing anything fancy, it's crucial to slow down and look at the data you've got. Not the idea of the data. The actual contents.

Load it right

In Python, we commonly use `pandas.read_csv()` to load tabular data. But things can get tricky fast. Maybe the file uses semicolons instead of commas. Maybe there's an encoding problem. Maybe the headers are actually in row 3. Loading data is often the first sign that something isn't quite right.

```
1 import pandas as pd
2
3 # Default load
4 df = pd.read_csv("data.csv")
5
6 # Custom separator and encoding
7 df = pd.read_csv("data.csv", sep=';', encoding='latin1')
```

Some useful arguments in `read_csv()`:

- `sep, delimiter`: For files using ; or \t
- `encoding`: Try '`utf-8`', '`latin1`', or '`ISO-8859-1`' if you see encoding errors
- `skiprows`: If the header isn't in the first line
- `na_values`: To specify what strings should be treated as missing

First visual inspection

Once loaded, start simple:

```
1 df.head()  
2 df.tail()  
3 df.sample(5)
```

Then go deeper:

```
1 df.info()  
2 df.describe()  
3 df.columns  
4 df.dtypes
```

These commands tell you:

- How many rows and columns there are
- If any columns are missing values
- What data types are being used
- If numerical columns have suspicious max/min values

Common surprises

At this stage, you often spot things like:

- IDs read as float instead of string (e.g., 00123 becomes 123 . 0)
- Columns with mixed types (e.g., numbers and text)
- Dates stored as object instead of datetime
- Missing values coded as “n/a”, “NA”, or “-”

These are all signals that the dataset needs attention before any analysis can begin.

Save a copy

Always save the raw version and keep your cleaning steps separate. This is good scientific practice, but also aligns with GCP and ALCOA+ principles: traceability, reproducibility, and transparency.

```
1 df.to_csv("cleaned/step0_raw_copy.csv", index=False)
```

This way, you can always go back to the original if something goes wrong later.

Summary

- Always inspect the data before doing anything else.
- Be careful with loading parameters—files can be messy.
- Print and scan the first few rows, check data types and summaries.
- Save a copy of the raw data before you begin cleaning.

We haven't modified anything yet—but now we know what we're dealing with. Let's get our hands dirty in the next chapter.

Handling Missing Values

Missing values are everywhere. In clinical data, you might have patients who skipped visits. In survey data, some respondents might refuse to answer certain questions. In IoT devices, a sensor might temporarily go offline. Whatever the source, missing data is a reality—and handling it well is a critical skill.

Why missingness matters

The way we handle missing data can have a big impact on the conclusions we draw. Drop too much, and you may lose signal or bias your sample. Impute the wrong way, and you risk creating false confidence. There's no universal solution, but there are strategies that work well in different situations.

Let's say you're building a model to predict blood glucose levels. If half your rows are missing BMI values, dropping those rows could eliminate important patterns. On the other hand, filling them all with the average might mask real variability. Every decision here influences your final model.

Types of missing data

Statisticians typically distinguish three types:

- **MCAR (Missing Completely At Random):** There's no pattern in the missingness. This is ideal (and rare).

- **MAR (Missing At Random):** The missingness depends on observed variables.
- **MNAR (Missing Not At Random):** The missingness depends on unobserved variables.

For example:

- MCAR: A lab machine randomly failed to record values.
- MAR: Older patients are less likely to answer a digital survey.
- MNAR: Patients with very high blood pressure skipped follow-up visits.

Understanding this classification helps you decide what methods are safe to use.

Detecting missing values in pandas

```
1 df.isna().sum()  
2 df[df['column_name'].isna()]
```

To get a quick overview:

```
1 missing_percent = df.isna().mean().round(2) * 100  
2 print(missing_percent.sort_values(ascending=False))
```

You can also visualize missingness with the `missingno` library:

```
1 import missingno as msno  
2 msno.matrix(df)  
3 msno.heatmap(df)
```

This shows which rows and columns have gaps and whether patterns exist. Heatmaps are great for spotting if two variables tend to be missing together.

Strategies to handle missing values

Drop them (but be careful)

```
1 df_clean = df.dropna()
```

This works if only a few rows are affected and the loss of data won't bias the outcome. But if 30% of your dataset disappears, maybe it's time to pause.

You can also drop only rows missing specific columns:

```
1 df = df.dropna(subset=['age', 'sex'])
```

Or drop columns with too many missing values:

```
1 df = df.dropna(axis=1, thresh=len(df) * 0.5) # keep
columns with >50% non-null
```

Impute them (simple methods)

```
1 df['age'].fillna(df['age'].median(), inplace=True)
2 df['sex'].fillna('Unknown', inplace=True)
```

Be cautious with mean imputation—it can distort distributions, especially with skewed data. Median is often a safer bet.

Impute them (fancy methods)

You can use KNN imputation, regression models, or multiple imputation packages. One example with `sklearn`:

```
1 from sklearn.impute import KNNImputer
2 imputer = KNNImputer(n_neighbors=5)
```

```
3 df[['age', 'bmi']] = imputer.fit_transform(df[['age', 'bmi']])
```

Or with `IterativeImputer` (multivariate regression):

```
1 from sklearn.experimental import enable_iterative_imputer
2 from sklearn.impute import IterativeImputer
3
4 imp = IterativeImputer(random_state=42)
5 df_imputed = imp.fit_transform(df.select_dtypes(include='number'))
```

These techniques try to “guess” missing values using other variables, which is powerful—but make sure to cross-validate your models to avoid overfitting to imputed data.

Flag them

Another trick is to add an indicator column:

```
1 df['age_missing'] = df['age'].isna()
```

This way, your model can learn from the fact that data was missing, which might itself be informative. For example, patients who skipped lab tests might differ systematically from those who didn’t.

Real-world example: BMI in a clinical dataset

Suppose you’re analyzing data from a metabolic clinic and notice 15% of the patients are missing BMI. You:

1. Create a histogram of BMI to check its distribution.
2. Decide to fill missing values with the median.
3. Add a column `bmi_missing = True/False`.

4. Log this step and save an intermediate file.

```
1 bmi_median = df['bmi'].median()  
2 df['bmi_missing'] = df['bmi'].isna()  
3 df['bmi'].fillna(bmi_median, inplace=True)  
4 df.to_csv("cleaned/step1_bmi_imputed.csv", index=False)
```

Now you've preserved the information, filled the blanks, and left a trail others can follow.

Documentation is key

Don't just fill in blanks—write down what you did and why. This step matters for ALCOA+ and for reproducibility. Also, store your imputations in separate scripts or notebook cells so they can be reviewed. If someone else picks up your project, they should be able to trace every change.

Summary

- Missing data is normal, but must be handled thoughtfully.
- Understand the type of missingness before choosing a method.
- Use pandas or missingno to detect and explore missing values.
- Decide whether to drop, impute, or flag—and document your choices.
- Be cautious with automated imputations: always think before you fill.

Handling missing values well is a small step that can make a huge difference later. Your future self—and anyone reviewing your analysis—will thank you for being careful here.

Fixing Data Types and Formats

Getting data into the right format isn't glamorous, but it's essential. If your dates are stored as strings, your categories are treated as free text, and your numbers come in as objects, you're going to hit a wall. Models won't train. Visualizations will fail. Merges will break.

Correcting data types is one of those silent but powerful steps that makes everything else work.

Why it matters

Every column in a dataset has a type: integer, float, string, datetime, boolean, categorical... but when data is loaded from CSVs, Excel files, or databases, pandas has to guess what type each column is. And it often guesses wrong.

For example:

- `object` type might actually be a date, a number, or a categorical variable.
- IDs like 00123 might be turned into 123.0, losing leading zeros.
- `True/False` values might show up as strings: "Yes", "no", "TRUE", "0".

The earlier you detect and fix these issues, the smoother everything else becomes.

Checking types in pandas

```
1 df.dtypes
```

To get more detail:

```
1 print(df.info())
```

You'll see which columns are `object`, `float64`, `int64`, `bool`, `datetime64`, and so on.

Fixing numeric columns

Sometimes numbers are stored as strings:

```
1 df['weight'] = pd.to_numeric(df['weight'], errors='coerce')
```

The `errors='coerce'` option turns non-numeric strings (like “n/a”) into `NaN`, which you can deal with later.

If you're mixing commas and dots:

```
1 df['height'] = df['height'].str.replace(',', '.').astype(float)
```

Converting to datetime

Dates are one of the most common troublemakers. They often come in weird formats.

```
1 df['visit_date'] = pd.to_datetime(df['visit_date'], errors='coerce')
```

You can specify the format if needed:

```
1 df['visit_date'] = pd.to_datetime(df['visit_date'],
                                     format='%d/%m/%Y')
```

Check for `NaT` (missing datetimes) and make sure all values were converted correctly.

To extract parts of a date:

```
1 df['year'] = df['visit_date'].dt.year
2 df['month'] = df['visit_date'].dt.month
3 df['weekday'] = df['visit_date'].dt.day_name()
```

Handling categorical data

Free-text categories are a pain. They lead to duplicates and inconsistencies:

```
1 print(df['sex'].unique())
```

You might see: `['Male', 'male', 'M', 'F', 'Female', 'f', '']`

First, clean and unify:

```
1 df['sex'] = df['sex'].str.strip().str.lower()
2 df['sex'] = df['sex'].replace({'male': 'M', 'female': 'F',
                               'm': 'M', 'f': 'F'})
```

Then convert to category:

```
1 df['sex'] = df['sex'].astype('category')
```

This saves memory and tells pandas that it's not just a string—it's a label with a finite set of values.

Boolean conversions

```
1 df['is_smoker'] = df['is_smoker'].map({'yes': True, 'no': False})
```

Or more safely:

```
1 df['is_smoker'] = df['is_smoker'].str.lower().map(lambda x: x in ['yes', 'true', '1'])
```

Dealing with unit conversions

Another hidden source of trouble is mismatched units. Imagine weight in kilograms and pounds mixed in the same column.

```
1 # Assuming weights over 250 are probably in pounds
2 mask = df['weight'] > 250
3 df.loc[mask, 'weight'] = df.loc[mask, 'weight'] * 0.453592
```

Add a flag to track which rows were changed:

```
1 df['weight_converted'] = mask
```

Always document when and why you did this. Unit mismatches can break even well-trained models.

Case study: fixing a lab dataset

You receive a dataset with the following issues:

- Dates stored as text in European format
- Blood pressure as strings: “120/80”

- BMI as a mix of floats and “n/a”

Step-by-step:

```
1 # Convert dates
2 df['visit_date'] = pd.to_datetime(df['visit_date'],
3                                   dayfirst=True, errors='coerce')
4
5 # Split blood pressure
6 df[['sbp', 'dbp']] = df['blood_pressure'].str.split('/', expand=True)
7 df['sbp'] = pd.to_numeric(df['sbp'], errors='coerce')
8 df['dbp'] = pd.to_numeric(df['dbp'], errors='coerce')
9
10 df['bmi'] = pd.to_numeric(df['bmi'], errors='coerce')
```

In less than 10 lines, you've fixed three major format problems and prepared the dataset for modeling or exploration.

Summary

- Always check your data types as soon as you load the data.
- Convert text to numeric, datetime, category, or boolean as needed.
- Watch for leading zeros, inconsistent formats, and strange values.
- Use `.astype()`, `pd.to_numeric()`, and `pd.to_datetime()` to take control.
- Don't forget unit mismatches—they can be subtle but serious.

When your data types are correct, everything else becomes easier. Clean types = clean analysis.

Standardizing and Cleaning Text Data

Messy text is one of the most underestimated data problems. Unlike numbers, text is inherently flexible and ambiguous. A single column might contain values like “Male”, “male”, “MALE”, “M”, or even just “m”—and they all mean the same thing. Unless we clean and standardize that text, our models and summaries won’t understand it.

Text data cleaning is about consistency. We’re not doing NLP here—we’re just getting labels and free-form strings into shape so we can group, filter, or model them effectively.

Why text cleaning matters

Let’s say you’re analyzing clinical trial data and want to group patients by diagnosis. If the column contains:

- 1 Diabetes
- 2 Type 2 diabetes
- 3 T2DM
- 4 type ii diabetes

...you’ll get four different groups unless you normalize them.

Common problems in text data

- Inconsistent capitalization
- Trailing or leading whitespace
- Typos and abbreviations
- Accent marks and special characters
- Mixing formats (e.g., “kg” vs “kilograms”)
- Empty strings or placeholder text (“n/a”, “–”)

Basic string cleaning in pandas

```
1 df['diagnosis'] = df['diagnosis'].str.strip()  
2 df['diagnosis'] = df['diagnosis'].str.lower()  
3 df['diagnosis'] = df['diagnosis'].str.replace('-', ' ')
```

Use `.str` methods to handle most common issues:

- `str.strip()` – remove whitespace
- `str.lower()` – standardize case
- `str.replace()` – fix formatting and symbols
- `str.title()` – useful for names
- `str.contains()` – find patterns

Replacing variants and typos

If you know the most common variants:

```
1 mapping = {  
2     't2dm': 'type 2 diabetes',  
3     'type ii diabetes': 'type 2 diabetes',  
4     'diabetes type 2': 'type 2 diabetes'}
```

```
5 }
6 df['diagnosis'] = df['diagnosis'].replace(mapping)
```

You can also use regex for pattern-based replacements:

```
1 df['code'] = df['code'].str.replace(r'^[A-Z0-9]', '',
                                     regex=True)
```

This removes non-alphanumeric characters, often useful in cleaning IDs.

Dealing with missing or meaningless text

Sometimes cells contain strings like “unknown”, “n/a”, or just “-”. You can treat them as missing:

```
1 df['notes'] = df['notes'].replace(['n/a', 'na', '-', 'unknown'],
                                    pd.NA)
```

Once standardized, you can handle them as proper missing values.

Categorical cleanup example

You receive a column of patient-reported smoking status:

```
1 Yes
2 yes
3 No
4 no
5 unknown
6 1
7 0
8 Y
9 N
```

Let's clean and standardize it:

```

1 def clean_smoking(value):
2     val = str(value).strip().lower()
3     if val in ['yes', 'y', '1', 'true']: return 'Yes'
4     elif val in ['no', 'n', '0', 'false']: return 'No'
5     else: return 'Unknown'
6
7 df['smoker'] = df['smoker'].apply(clean_smoking)
8 df['smoker'] = df['smoker'].astype('category')

```

Now you can group or analyze smoking status confidently.

Unicode and accent normalization

Text may include invisible characters or different encodings (especially in multilingual data). Normalize accents using `unicodedata`:

```

1 import unicodedata
2
3 def remove_accents(text):
4     if isinstance(text, str):
5         return ''.join(
6             c for c in unicodedata.normalize('NFKD', text)
7                 )
8         if not unicodedata.combining(c)
9     )
10    return text
11 df['city'] = df['city'].apply(remove_accents)

```

Now “Málaga” and “Malaga” will match.

Summary

- Text fields are full of inconsistencies—don’t trust them blindly.

- Use `.str` methods and mappings to clean and unify values.
- Replace placeholder text like “n/a” with proper nulls.
- Normalize accents and symbols if needed.
- Always convert cleaned strings to categorical when possible.

This kind of cleaning isn’t flashy—but it’s the difference between chaos and clarity.

Detecting and Handling Duplicates

Duplicate records are like cockroaches—if you see one, there's probably more. They sneak in through system exports, user errors, or merging datasets. Sometimes duplicates are exact copies. Other times they differ by a timestamp, a typo, or an extra space.

Not all duplicates are bad. You might expect repeated lab results, or multiple visits per patient. But when duplicates are unintended, they can skew statistics, inflate sample sizes, or confuse downstream processing.

What is a duplicate?

In pandas, a duplicate row is one where all column values match another row. But you can define duplicates based on a subset of columns if needed.

Example:

```
1 # Exact duplicates
2 df.duplicated().sum()
3
4 # Drop them
5 df = df.drop_duplicates()
```

To check by specific columns:

```
1 df.duplicated(subset=['patient_id', 'visit_date']).sum()
```

Visual inspection

Start by sorting:

```
1 df.sort_values(by=['patient_id', 'visit_date']).head(10)
```

Sometimes duplicates are not identical but *nearly* so:

- Same patient ID and date, but different lab values
- Same name and birthdate, but different phone number

This is where your domain knowledge comes in.

Handling duplicates safely

1. Drop exact duplicates

```
1 df = df.drop_duplicates()
```

You can also keep the *last* occurrence instead of the first:

```
1 df = df.drop_duplicates(keep='last')
```

2. Drop by key columns

```
1 df = df.drop_duplicates(subset=['patient_id', 'visit_date'])
```

Be careful—this may drop rows that are valid but repeated.

3. Aggregate near-duplicates

Sometimes it's better to collapse rows:

```
1 df = df.groupby(['patient_id', 'visit_date']).agg({  
2     'glucose': 'mean',  
3     'bmi': 'first'  
4 }).reset_index()
```

You choose which fields to average, count, or take the first value from.

4. Use fuzzy matching (advanced)

When names or text fields are close but not identical, you can use [fuzzywuzzy](#) or [recordlinkage](#). Example:

```
1 from fuzzywuzzy import fuzz  
2 fuzz.ratio("John Smith", "Jon Smith") # 96
```

This helps detect records that look like duplicates even if they're not exact matches.

Logging duplicate removal

Always log how many rows were removed and why:

```
1 before = len(df)  
2 df = df.drop_duplicates()  
3 after = len(df)  
4 print(f"Removed {before - after} duplicate rows")
```

If working under ALCOA+, document criteria used and keep an original backup.

Case study: clinical visit log

You receive an export of visits from a hospital system. Some patients have 2 or 3 identical entries for the same date. After talking to the team, you learn it's a system glitch that repeated the same record.

```
1 duplicates = df[df.duplicated(subset=['patient_id', 'visit_date'])]
2 print("Duplicates found:", len(duplicates))
3
4 # Drop safely
5 df = df.drop_duplicates(subset=['patient_id', 'visit_date'])
```

You document the issue, clean the dataset, and save a clean copy:

```
1 df.to_csv("cleaned/step2_duplicates_removed.csv", index=False)
```

Summary

- Duplicates can distort analysis if left unchecked.
- Use `duplicated()` and `drop_duplicates()` to detect and clean.
- Decide if duplicates are true errors or expected repetitions.
- Consider grouping and aggregating near-duplicates.
- Always log what you removed and why.

Handling duplicates is about being precise, careful, and curious. They're not just technical artifacts—they're clues to how your data was created.

Outlier Detection and Correction

Outliers are values that look suspiciously far from the rest. Maybe they're errors. Maybe they're just rare events. Either way, they can have a big impact—especially on statistics like mean and standard deviation, or models sensitive to scale.

But here's the tricky part: not all outliers are wrong. In medicine, a BMI of 60 is unusual, but it could be correct. In finance, a sudden spike in sales might reflect a real event. So before deleting or correcting anything, we need to **understand** why an outlier exists.

What is an outlier?

There's no single definition, but here are a few signs:

- A value far outside the normal range
- A point that falls outside the interquartile range (IQR)
- A value more than 3 standard deviations from the mean
- A point that breaks a business or clinical rule (e.g., heart rate > 300)

Visualizing outliers

Boxplots

```
1 import seaborn as sns  
2 sns.boxplot(x=df['bmi'])
```

Boxplots are great for seeing the spread of the data and identifying points outside the whiskers.

Histograms

```
1 df['glucose'].hist(bins=50)
```

Helpful for spotting spikes, long tails, or gaps.

Scatter plots

Use when outliers depend on context:

```
1 sns.scatterplot(x='age', y='cholesterol', data=df)
```

Z-scores

For numerical columns:

```
1 from scipy.stats import zscore  
2 z = zscore(df['bmi'].dropna())  
3 outliers = df.loc[(z < -3) | (z > 3)]
```

IQR method

```
1 Q1 = df['bmi'].quantile(0.25)  
2 Q3 = df['bmi'].quantile(0.75)  
3 IQR = Q3 - Q1
```

```
4  
5 outliers = df[(df['bmi'] < Q1 - 1.5 * IQR) | (df['bmi'] >  
Q3 + 1.5 * IQR)]
```

What to do with outliers

1. Investigate

- Is it a data entry error?
- Does it break known constraints?
- Could it be a valid extreme case?

2. Correct if clearly wrong

```
1 # Example: age of 250 is likely a typo  
2 df.loc[df['age'] > 120, 'age'] = pd.NA
```

Or set to null to review later:

```
1 mask = df['glucose'] > 1000  
2 df.loc[mask, 'glucose'] = pd.NA
```

3. Cap or clip

Limit values to a reasonable range:

```
1 df['bmi'] = df['bmi'].clip(lower=10, upper=60)
```

This preserves the row but limits the influence of the extreme value.

4. Flag them

Mark outliers for later review:

```
1 Q1 = df['bmi'].quantile(0.25)
2 Q3 = df['bmi'].quantile(0.75)
3 IQR = Q3 - Q1
4
5 df['bmi_outlier'] = (df['bmi'] < Q1 - 1.5 * IQR) | (df['bmi'] > Q3 + 1.5 * IQR)
```

5. Remove them (last resort)

Only when justified:

```
1 df = df[~df['bmi_outlier']]
```

Don't delete rows just because they make your plot prettier.

Real-world example: metabolic data

You're analyzing fasting glucose levels. Most values are between 70 and 120 mg/dL. But a few values are over 1000—more than ten times the typical range.

Step-by-step:

1. Plot a histogram.
2. Use IQR to define outliers.
3. Set outliers to NaN.
4. Add an outlier flag.
5. Document the change.

```
1 Q1 = df['glucose'].quantile(0.25)
2 Q3 = df['glucose'].quantile(0.75)
```

```
3 IQR = Q3 - Q1
4
5 mask = (df['glucose'] < Q1 - 1.5 * IQR) | (df['glucose'] > Q3 + 1.5 * IQR)
6 df['glucose_outlier'] = mask
7 df.loc[mask, 'glucose'] = pd.NA
```

Summary

- Outliers can be real or errors—investigate before acting.
- Use boxplots, histograms, z-scores, or IQR to detect them.
- Consider capping, flagging, or removing—but document all steps.
- Never assume that a strange value is wrong without context.

Outliers are warning lights. They tell you where to look more closely—and sometimes, where the real story is hiding.

Resolving Inconsistencies Across Columns

Some problems aren't obvious by looking at individual columns—they live in the relationships between columns. Maybe a patient has a discharge date that comes before their admission. Or a BMI that doesn't match their height and weight. Or someone marked as deceased but with follow-up lab results.

These kinds of cross-field inconsistencies are common in real datasets. Fixing them is a bit like detective work—you need to compare, calculate, and ask “does this make sense?”

Common types of column inconsistencies

- **Date logic errors:** end date before start date, birth date in the future
- **Redundant fields mismatch:** BMI doesn't match height and weight
- **Flag contradictions:** `pregnant = True` and `sex = male`
- **Linked variables drift:** weight increased by 20 kg but height stayed the same

Logical checks with dates

Start with simple comparisons:

```

1 df['admit_before_discharge'] = df['admit_date'] <= df['
    discharge_date']
2 invalid_dates = df[~df['admit_before_discharge']]
```

Flag invalid rows and review:

```
1 print(invalid_dates[['patient_id', 'admit_date', ''
    'discharge_date']])
```

Set them to `NaT` if needed:

```
1 df.loc[~df['admit_before_discharge'], 'discharge_date'] =
    pd.NaT
```

Recalculating from raw data

If you have height and weight, you can recompute BMI:

```

1 df['height_m'] = df['height_cm'] / 100
2 df['bmi_calc'] = df['weight_kg'] / df['height_m']**2
```

Compare with recorded BMI:

```

1 bmi_diff = (df['bmi'] - df['bmi_calc']).abs()
2 df['bmi_mismatch'] = bmi_diff > 1.5
```

Now you can decide whether to trust the calculated or recorded value—or keep both.

Contradictory flags

Example: someone marked `pregnant = True` and `sex = male`.

```
1 mask = (df['pregnant'] == True) & (df['sex'] == 'M')
2 df.loc[mask, ['pregnant', 'sex']]
```

If it's a data entry mistake, you may want to:

- Set `pregnant = False`
- Set it to `NaN`
- Flag the row for review

Always log how many contradictions were found:

```
1 print("Contradictions found:", mask.sum())
```

Category mismatch

In some cases, values don't match expected category pairs:

- A "follow-up" record without an initial visit
- A "discharged" flag without an admission
- A "completed" task with a null timestamp

You can detect these with conditional checks:

```
1 mask = (df['visit_type'] == 'follow-up') & df['
    initial_visit_id'].isna()
2 df['visit_issue'] = mask
```

Value correlation

If two variables should be tightly correlated, check them:

```
1 sns.scatterplot(x='total_cholesterol', y='ldl')
```

If LDL is higher than total cholesterol for multiple cases, something's wrong.

You can also calculate correlation:

```
1 corr = df[['total_cholesterol', 'ldl']].corr()
2 print(corr)
```

Case study: verifying timestamps

You're working with ICU data. Each row has `start_time`, `end_time`, and `duration_minutes`.

```
1 # Recalculate duration
2 df['duration_calc'] = (df['end_time'] - df['start_time'])
   .dt.total_seconds() / 60
3
4 # Compare with recorded
5 df['duration_diff'] = (df['duration_calc'] - df['
   duration_minutes']).abs()
6 df['duration_mismatch'] = df['duration_diff'] > 5
```

Now you've spotted where duration was miscalculated, possibly from rounding errors or manual edits.

Summary

- Some errors only appear when comparing multiple columns.
- Use logic and simple math to detect contradictions.
- Recalculate values from raw inputs when possible.
- Flag inconsistencies, correct when obvious, and always document.

These aren't just bugs—they're clues. Data that doesn't make sense across columns tells you a story. Sometimes that story is about biology. Sometimes it's about a typo.

Validating and Enforcing Constraints

Data validation is about making sure the values in your dataset make sense—not just in isolation, but according to defined rules. These rules might come from your domain knowledge, regulatory standards, or logical expectations. Think of it as teaching your dataset some basic manners.

If your analysis is a house, validation is the plumbing check before you move in.

Why validation matters

Anyone can make a mistake. A user might enter a temperature in Fahrenheit instead of Celsius. A timestamp might get truncated. A value might be typed as “999” instead of “99”. These small errors can cascade into bigger problems if left unchecked.

By enforcing constraints early, you catch issues before they cause confusion downstream.

Types of validations

Range checks

Ensure values fall within acceptable limits:

```
1 mask = (df['age'] >= 0) & (df['age'] <= 120)
2 df['age_invalid'] = ~mask
```

Allowed values (categorical checks)

```
1 valid_values = ['M', 'F']
2 df['sex_valid'] = df['sex'].isin(valid_values)
```

Null constraints

Make sure required fields aren't missing:

```
1 df['admit_date_valid'] = df['admit_date'].notna()
```

Cross-field dependencies

```
1 # Pregnancy flag only allowed if sex is female
2 mask = (df['sex'] == 'F') | (df['pregnant'] != True)
3 df['pregnancy_constraint_ok'] = mask
```

Pattern checks with regex

```
1 # Email format
2 pattern = r'^[\w\.-]+@[\\w\.-]+\.\w{2,4}$'
3 df['email_valid'] = df['email'].str.match(pattern)
```

Validating with pandera

pandera is a Python library to define schemas and validate pandas DataFrames.

Example:

```
1 import pandera as pa
2 from pandera import Column, DataFrameSchema
3
4 schema = DataFrameSchema({
5     "age": Column(pa.Int, checks=pa.Check.in_range(0,
6                 120)),
7     "sex": Column(pa.String, checks=pa.Check.isin(["M", "F"])),
8     "bmi": Column(pa.Float, nullable=True),
9 })
10 validated_df = schema.validate(df)
```

If something's wrong, it raises a clear error.

Using pydantic for row-level validation

If your data comes row by row (e.g., from an API), `pydantic` is a great option.

```
1 from pydantic import BaseModel, Field, validator
2
3 class Patient(BaseModel):
4     age: int = Field(..., ge=0, le=120)
5     sex: str
6
7     @validator("sex")
8     def check_sex(cls, v):
9         if v not in ["M", "F"]:
10             raise ValueError("Invalid sex")
11         return v
```

Referential integrity

If you're working with multiple tables (e.g., patients and visits), you need to ensure foreign keys match.

```
1 valid_ids = patients['patient_id']
2 df['id_in_patients'] = df['patient_id'].isin(valid_ids)
```

If you find rows that reference a non-existent patient, that's a red flag.

Constraint enforcement with assertions

You can add assertions at the top of notebooks or scripts:

```
1 assert df['age'].between(0, 120).all(), "Age out of range
  !"
```

This helps ensure things don't break silently.

Case study: validating a clinical dataset

You receive a dataset with demographics, vital signs, and lab results. You want to ensure:

- Age is 0–120
- Sex is M or F
- All rows have a patient ID
- Glucose values are realistic

Step-by-step:

```
1 assert df['age'].between(0, 120).all()
2 assert df['sex'].isin(['M', 'F']).all()
```

```
3 assert df['patient_id'].notna().all()
4 assert df['glucose'].between(30, 800).all()
```

Alternatively, you could wrap these into a `pandera` schema.

Logging and documenting rules

Whenever you apply validations, document them. You can:

- Keep a `constraints.md` file describing each rule
- Store validation functions in a module (e.g. `validate.py`)
- Add logging to your pipeline

```
1 import logging
2 logging.info("Validated sex field: %d invalid entries",
   (~df['sex'].isin(['M', 'F'])).sum())
```

Summary

- Validation is the practice of enforcing logic and sanity checks on your data.
- Use Python logic, regex, or tools like `pandera` and `pydantic`.
- Catch issues early—don’t wait for your model to fail.
- Document all rules applied.

Clean data isn’t just about missing values or duplicates—it’s also about values that make sense. Validation is your last line of defense before analysis begins.

Logging, Versioning, and Audit Trails

You've cleaned your data. It looks great. But how will someone else know what you changed—and why? And more importantly: how will **you** remember what you did three weeks from now?

That's where logging, versioning, and audit trails come in. These aren't just "nice to have" in regulated environments—they're essential for reproducibility, trust, and scientific integrity.

Why traceability matters

In clinical research, data must comply with ALCOA+ principles:

- **Attributable:** Who made the change?
- **Legible:** Can we read and understand it?
- **Contemporaneous:** Was it recorded at the time?
- **Original:** Do we have the raw version?
- **Accurate:** Is it correct?

The “+” adds: Complete, Consistent, Enduring, and Available.

These principles don't just apply to the data—they apply to your **code, decisions, and process**.

Keeping raw and cleaned data separate

Never overwrite the original file. Always save cleaned versions with clear names:

```
1 df.to_csv("cleaned/step3_validated.csv", index=False)
```

Include a README or changelog describing what changed:

```
1 2025-05-08 | Removed outliers and fixed BMI mismatches.  
           | Stored in step3_validated.csv
```

You can even log the number of rows affected:

```
1 removed = initial_len - len(df)  
2 print(f"Removed {removed} rows with invalid age")
```

Logging with the logging module

Instead of just printing, log messages properly:

```
1 import logging  
2 logging.basicConfig(filename='logs/cleaning.log', level=  
                      logging.INFO)  
3 logging.info("Step 2: Replaced missing BMI values with  
               median")
```

Log unexpected values too:

```
1 bad_sex = df[~df['sex'].isin(['M', 'F'])]  
2 logging.warning(f"Found {len(bad_sex)} invalid sex  
                  entries")
```

Versioning data with DVC

DVC is like Git for data. You can version .csv or .parquet files and track them with your code.

Initialize a DVC project:

```
1 dvc init  
2 dvc add cleaned/step3_validated.csv  
3 git add cleaned/step3_validated.csv.dvc .gitignore  
4 git commit -m "Track validated dataset with DVC"
```

This allows rollback, comparison, and reproducibility.

Recording code versions

Use Git to track your notebooks and scripts. Commit often:

```
1 git commit -am "Cleaned text fields and recoded categories"
```

Tag important stages:

```
1 git tag -a v1.0 -m "Final cleaned dataset before modeling"  
"
```

Jupyter notebooks and cell history

Jupyter keeps execution counts, but it's easy to rerun cells in a different order and lose traceability. Use these tips:

- Run notebooks top to bottom
- Export to HTML or PDF after each milestone
- Use `nbconvert` or tools like `papermill` or `jupytext`

Reproducible scripts

Move cleaning steps into scripts:

```
1 python clean_data.py
```

This ensures anyone can re-run the pipeline, not just click through a notebook.

Include command-line logging:

```
1 print("Step 1 complete: dropped duplicates")
```

Or structured logging:

```
1 import logging  
2 logging.info("Step 3 complete: imputed missing glucose")
```

Case study: versioned cleaning pipeline

You create a script `cleaning_pipeline.py` that:

- Loads raw data from `data/raw.csv`
- Cleans types, missing values, and outliers
- Saves to `cleaned/step_final.csv`
- Logs all actions
- Is tracked by Git and DVC

Now, anyone on your team (or a reviewer) can:

- Re-run the entire process
- See exactly what was changed
- Compare results before and after cleaning

Summary

- Always keep raw and cleaned data separate
- Use logging to document every change
- Version data files with DVC or commit intermediate steps manually
- Track your code with Git and tag important milestones
- Package your cleaning steps into scripts for full reproducibility

Clean data is great. But clean **processes** are even better. That's what makes your work solid, trusted, and ready for the real world.

Case Studies in Data Cleaning

Let's wrap up with some real-world examples. These case studies combine many of the techniques we've explored so far, and show how they're applied to messy, realistic data. Each one starts from a raw dataset and walks through the steps to make it usable.

Case 1: Cleaning Clinical EHR Data

You're given an extract of electronic health records (EHR) with the following columns:

- `patient_id`
- `sex`
- `birth_date`
- `visit_date`
- `height_cm, weight_kg, bmi`
- `glucose, blood_pressure, notes`

Problems spotted:

- `sex` values: “M”, “F”, “male”, “FEMALE”, “unknown”
- `birth_date` stored as text, some in DD/MM/YYYY, some in MM-DD-YYYY
- `bmi` doesn't match height and weight in many cases
- Notes contain typos and irregular whitespace

Cleaning steps:

```

1 # Standardize sex
2 sex_map = {'male': 'M', 'female': 'F', 'f': 'F', 'm': 'M'}
3 df['sex'] = df['sex'].str.strip().str.lower().replace(
    sex_map)
4 df['sex'] = df['sex'].where(df['sex'].isin(['M', 'F']),
    pd.NA)
5
6 # Parse birth_date
7 from dateutil import parser
8
9 def try_parse(date_str):
10     try:
11         return parser.parse(date_str, dayfirst=True)
12     except:
13         return pd.NaT
14
15 df['birth_date'] = df['birth_date'].apply(try_parse)
16
17 # Recalculate BMI
18 height_m = df['height_cm'] / 100
19 df['bmi_calc'] = df['weight_kg'] / (height_m ** 2)
20 df['bmi_diff'] = (df['bmi'] - df['bmi_calc']).abs()
21 df['bmi_flag'] = df['bmi_diff'] > 1.5
22
23 # Clean notes
24 import re
25 df['notes'] = df['notes'].str.strip().str.lower().str.
    replace(r'\s+', ' ', regex=True)

```

Results:

- Cleaned categorical fields
- Consistent date format
- Mismatched BMI flagged for review

- Notes simplified for NLP or grouping
-

Case 2: Public Health Indicators (WHO)

You're working with a dataset of country-level health metrics. It includes:

- `country, year`
- `life_expectancy, population, gdp_per_capita`
- `urban_pct, mortality_rate, vaccination_coverage`

Problems spotted:

- Missing `gdp_per_capita` in some countries
- Duplicate rows for some (`country, year`) pairs
- `urban_pct` has values over 100
- `vaccination_coverage` includes “n/a”

Cleaning steps:

```
1 # Drop duplicates
2 df = df.drop_duplicates(subset=['country', 'year'])
3
4 # Convert to numeric
5 df['gdp_per_capita'] = pd.to_numeric(df['gdp_per_capita'],
6                                     errors='coerce')
6 df['vaccination_coverage'] = pd.to_numeric(df['
7     vaccination_coverage'], errors='coerce')
7
8 # Cap invalid values
9 df['urban_pct'] = df['urban_pct'].clip(upper=100)
```

```
10  
11 # Impute GDP using country median  
12 df['gdp_per_capita'] = df.groupby('country')[  
    'gdp_per_capita'].transform(  
        lambda x: x.fillna(x.median()))  
13  
14 )
```

Results:

- Cleaned dataset ready for time series analysis
- No impossible values (e.g., 120% urban population)
- Reproducible script with clear assumptions

Case 3: LC-MS Metabolomics Data

You're analyzing untargeted LC-MS metabolomics data with compound intensities across multiple samples:

- `sample_id`, `compound_name`, `retention_time`, `area`, `flag`

Issues found:

- Some compounds have duplicate peaks
- `retention_time` varies slightly between replicates
- Some flags are “bb”, “MM”, or missing
- Intensities vary by 10x between replicates

Cleaning actions:

```
1 # Filter based on flag
2 valid_flags = [''] # empty string means no issues
3 df = df[df['flag'].isin(valid_flags)]
4
5 # Group by sample and compound
6 summary = df.groupby(['sample_id', 'compound_name']).agg(
7     {
8         'area': 'mean',
9         'retention_time': 'median'
10    }).reset_index()
11
12 # Log number of removed flagged rows
13 removed = len(df_original) - len(df)
14 print(f"Removed {removed} rows based on flag status")
```

Outcome:

- Aggregated dataset with clean peaks per sample
 - Median retention time per compound
 - Flags documented and discarded safely
-

Summary

These case studies show that real data is messy in many ways—text, dates, numbers, duplicates, logic, and structure. Cleaning it well takes patience, creativity, and rigor. But once you do, you unlock insights you can trust—and models that work.

You're not just fixing a spreadsheet. You're preparing evidence.

References and Resources

A well-cleaned dataset is a silent achievement—it often goes unnoticed, but it makes everything else possible. Below are key references, tools, and datasets mentioned throughout the book, especially in the case studies.

Key References

- Karr, A. F., Sanil, A. P., & Banks, D. L. (2006). *Data quality: A statistical perspective*. Statistical Methodology, 3(2), 137–173. <https://doi.org/10.1016/j.stamet.2005.08.005>
- EMA (2010). *Reflection paper on expectations for electronic source data and data transcribed to electronic data collection tools in clinical trials*. European Medicines Agency. https://www.ema.europa.eu/en/documents/scientific-guideline/reflection-paper-expectations-electronic-source-data-data-transcribed-electronic-data-collection_en.pdf
- FDA (2018). *Data Integrity and Compliance With Drug CGMP: Questions and Answers*. <https://www.fda.gov/media/119267/download>
- Pandera documentation: <https://pandera.readthedocs.io>
- Pydantic documentation: <https://docs.pydantic.dev/>
- DVC (Data Version Control): <https://dvc.org>
- Missingno: <https://github.com/ResidentMario/missingno>

Datasets Used in Case Studies

Clinical EHR Dataset (Synthetic Example)

To replicate the clinical cleaning steps, you can use a synthetic but realistic dataset like:

- [Synthea open EHR data](#) — Fully synthetic patient records, downloadable in CSV format.
- Sample dataset used in SeerPy (cleaned EHR structure): <https://github.com/SEERData/SeerPy> (exploration purposes)

WHO Health Indicators Dataset

- [WHO Global Health Observatory Data Repository](#)
- Specific indicators like life expectancy, mortality, and immunization: <https://www.who.int/data/gho/indicator-metadata-registry>

Metabolomics LC-MS Dataset

- Example LC-MS data (MTBLS files):
 - [MetaboLights - MTBLS135](#) — Plasma metabolome data.
 - [MTBLS404](#) - Human serum NMR and MS data

These provide compound tables with retention time, area, and sample annotations.

Tools Used Throughout

Essential Guide to Clean Data

Tool	Purpose	Link
pandas	Data manipulation	https://pandas.pydata.org
seaborn/matplotlib	Visualization	https://seaborn.pydata.org
scikit-learn	Imputation, scaling, model prep	https://scikit-learn.org
pandera	DataFrame schema validation	https://pandera.readthedocs.io
pydantic	Row-level validation and models	https://docs.pydantic.dev
DVC	Data versioning and reproducibility	https://dvc.org
missingno	Visualizing missing data	https://github.com/ResidentMario/missingno

If you want to go further, consider browsing the datasets on:

- Kaggle Datasets
- UCI Machine Learning Repository
- OpenML

And remember: behind every trusted insight is a dataset that someone took the time to clean properly. Now, that someone is you.

