

# Data Manipulation Challenge

## A Mental Model for Method Chaining in Pandas

### Data Manipulation Challenge - A Mental Model for Method Chaining in Pandas

! Challenge Requirements In Section [Student Analysis Section](#)

- Complete all discussion questions for the seven mental models (plus some extra requirements for higher grades)

! Note on Python Usage

**Recommended Workflow: Use Your Existing Virtual Environment** If you completed the Tech Setup Challenge Part 2, you already have a virtual environment set up! Here's how to use it for this new challenge:

1. **Clone this new challenge repository** (see Getting Started section below)
2. **Open the cloned repository in Cursor**
3. **Set this project to use your existing Python interpreter:**
  - Press **Ctrl+Shift+P** → “Python: Select Interpreter”
  - Navigate to and choose the interpreter from your existing virtual environment (e.g., `your-previous-project/venv/Scripts/python.exe`)
4. **Activate the environment in your terminal:**
  - Open terminal in Cursor (‘Ctrl + ‘)
  - Navigate to your previous project folder where you have the `venv` folder
  - **Pro tip:** You can quickly navigate by typing `cd` followed by dragging the folder from your file explorer into the terminal
  - Activate using the appropriate command for your system:
    - **Windows Command Prompt:** `venv\Scripts\activate`

- **Windows PowerShell:** `.\venv\Scripts\Activate.ps1`
- **Mac/Linux:** `source venv/bin/activate`
- You should see (venv) at the beginning of your terminal prompt

5. **Install additional packages if needed:** `pip install pandas numpy matplotlib seaborn`

#### Cloud Storage Warning

**Avoid using Google Drive, OneDrive, or other cloud storage for Python projects!** These services can cause issues with: - Package installations failing due to file locking - Virtual environment corruption - Slow performance during pip operations

**Best practice:** Keep your Python projects in a local folder like `C:\Users\YourName\Documents\` or `~/Documents/` instead of cloud-synced folders.

**Alternative: Create a New Virtual Environment** If you prefer a fresh environment, follow the Quarto documentation: <https://quarto.org/docs/projects/virtual-environments.html>. Be sure to follow the instructions to activate the environment, set it up as your default Python interpreter for the project, and install the necessary packages (e.g. pandas) for this challenge. For installing the packages, you can use the `pip install -r requirements.txt` command since you already have the requirements.txt file in your project. Some steps do take a bit of time, so be patient.

**Why This Works:** Virtual environments are portable - you can use the same environment across multiple projects, and Cursor automatically activates it when you select the interpreter!

## The Problem: Mastering Data Manipulation Through Method Chaining

**Core Question:** How can we efficiently manipulate datasets using `pandas` method chaining to answer complex business questions?

**The Challenge:** Real-world data analysis requires combining multiple data manipulation techniques in sequence. Rather than creating intermediate variables at each step, method chaining allows us to write clean, readable code that flows logically from one operation to the next.

**Our Approach:** We'll work with ZappTech's shipment data to answer critical business questions about service levels and cross-category orders, using the seven mental models of data manipulation through `pandas` method chaining.

### AI Partnership Required

This challenge pushes boundaries intentionally. You'll tackle problems that normally require weeks of study, but with Cursor AI as your partner (and your brain keeping it honest), you can accomplish more than you thought possible.

**The new reality:** The four stages of competence are Ignorance → Awareness → Learning → Mastery. AI lets us produce Mastery-level work while operating primarily in the Awareness stage. I focus on awareness training, you leverage AI for execution, and together we create outputs that used to require years of dedicated study.

## The Seven Mental Models of Data Manipulation

The seven most important ways we manipulate datasets are:

1. **Assign:** Add new variables with calculations and transformations
2. **Subset:** Filter data based on conditions or select specific columns
3. **Drop:** Remove unwanted variables or observations
4. **Sort:** Arrange data by values or indices
5. **Aggregate:** Summarize data using functions like mean, sum, count
6. **Merge:** Combine information from multiple datasets
7. **Split-Apply-Combine:** Group data and apply functions within groups

## Data and Business Context

We analyze ZappTech's shipment data, which contains information about product deliveries across multiple categories. This dataset is ideal for our analysis because:

- **Real Business Questions:** CEO wants to understand service levels and cross-category shopping patterns
- **Multiple Data Sources:** Requires merging shipment data with product category information
- **Complex Relationships:** Service levels may vary by product category, and customers may order across categories
- **Method Chaining Practice:** Perfect for demonstrating all seven mental models in sequence

## Data Loading and Initial Exploration

Let's start by loading the ZappTech shipment data and understanding what we're working with.

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from datetime import datetime, timedelta

# Load the shipment data
shipments_df = pd.read_csv(
    "https://raw.githubusercontent.com/flyaflya/persuasive/main/shipments.csv",
    parse_dates=['plannedShipDate', 'actualShipDate']
)

# Load product line data
product_line_df = pd.read_csv(
    "https://raw.githubusercontent.com/flyaflya/persuasive/main/productLine.csv"
)

# Reduce dataset size for faster processing (4,000 rows instead of 96,805 rows)
shipments_df = shipments_df.head(4000)

print("Shipments data shape:", shipments_df.shape)
print("\nShipments data columns:", shipments_df.columns.tolist())
print("\nFirst few rows of shipments data:")
print(shipments_df.head(10))

print("\n" + "="*50)
print("Product line data shape:", product_line_df.shape)
print("\nProduct line data columns:", product_line_df.columns.tolist())
print("\nFirst few rows of product line data:")
print(product_line_df.head(10))

```

Shipments data shape: (4000, 5)

Shipments data columns: ['shipID', 'plannedShipDate', 'actualShipDate', 'partID', 'quantity']

First few rows of shipments data:

	shipID	plannedShipDate	actualShipDate	partID	quantity
0	10001	2013-11-06	2013-10-04	part92b16c5	6
1	10002	2013-10-15	2013-10-04	part66983b	2
2	10003	2013-10-25	2013-10-07	part8e36f25	1
3	10004	2013-10-14	2013-10-08	part30f5de0	1
4	10005	2013-10-14	2013-10-08	part9d64d35	6

5	10006	2013-10-14	2013-10-08	part6cd6167	15
6	10007	2013-10-14	2013-10-08	parta4d5fd1	2
7	10008	2013-10-14	2013-10-08	part08cadf5	1
8	10009	2013-10-14	2013-10-08	part5cc4989	10
9	10010	2013-10-14	2013-10-08	part912ae4c	1

=====

Product line data shape: (11997, 3)

Product line data columns: ['partID', 'productLine', 'prodCategory']

First few rows of product line data:

	partID	productLine	prodCategory
0	part00005ba	line4c	Liquids
1	part000b57d	line61	Machines
2	part00123bf	linec1	Marketable
3	part0021fc9	line61	Machines
4	part0027e86	line2f	Machines
5	part002ed95	line4c	Liquids
6	part0030856	lineb8	Machines
7	part0033dfd	line49	Liquids
8	part0037a2a	linea3	Marketable
9	part003caee	linea3	Marketable

### Understanding the Data

**Shipments Data:** Contains individual line items for each shipment, including: - **shipID:** Unique identifier for each shipment - **partID:** Product identifier - **plannedShipDate:** When the shipment was supposed to go out - **actualShipDate:** When it actually shipped - **quantity:** How many units were shipped

**Product Category and Line Data:** Contains product category information: - **partID:** Links to shipments data - **productLine:** The category each product belongs to - **prodCategory:** The category each product belongs to

**Business Questions We'll Answer:** 1. Does service level (on-time shipments) vary across product categories? 2. How often do orders include products from more than one category?

## The Seven Mental Models: A Progressive Learning Journey

Now we'll work through each of the seven mental models using method chaining, starting simple and building complexity.

## 1. Assign: Adding New Variables

**Mental Model:** Create new columns with calculations and transformations.

Let's start by calculating whether each shipment was late:

```
# Simple assignment - calculate if shipment was late
shipments_with_lateness = (
    shipments_df
    .assign(
        is_late=lambda df: df['actualShipDate'] > df['plannedShipDate'],
        days_late=lambda df: (df['actualShipDate'] - df['plannedShipDate']).dt.days
    )
)

print("Added lateness calculations:")
print(shipments_with_lateness[['shipID', 'plannedShipDate', 'actualShipDate', 'is_late', 'days_late']])
```

Added lateness calculations:

	shipID	plannedShipDate	actualShipDate	is_late	days_late
0	10001	2013-11-06	2013-10-04	False	-33
1	10002	2013-10-15	2013-10-04	False	-11
2	10003	2013-10-25	2013-10-07	False	-18
3	10004	2013-10-14	2013-10-08	False	-6
4	10005	2013-10-14	2013-10-08	False	-6

### 💡 Method Chaining Tip for New Python Users

**Why use `lambda df:`?** When chaining methods, we need to reference the current state of the dataframe. The `lambda df:` tells pandas “use the current dataframe in this calculation.” Without it, pandas would look for a variable called `df` that doesn't exist.

**Alternative approach:** You could also write this as separate steps, but method chaining keeps related operations together and makes the code more readable.

### ❗ Discussion Questions: Assign Mental Model

**Question 1: Data Types and Date Handling** - What is the `dtype` of the `actualShipDate` series? How can you find out using code?

```
print(shipments_df['actualShipDate'].dtype)
```

```
datetime64[ns]
```

- Why is it important that both `actualShipDate` and `plannedShipDate` have the same data type for comparison? `#+ #` Answer: Shape of `actualShipDate` and `plannedShipDate` must be the same for comparison.

```
#+  
# Shape of actualShipDate and plannedShipDate  
#-  
print(shipments_df['actualShipDate'].shape)  
print(shipments_df['plannedShipDate'].shape)
```

```
(4000,)
```

```
(4000,)
```

**Question 2: String vs Date Comparison** - Can you give an example where comparing two dates as strings would yield unintuitive results, e.g. what happens if you try to compare “04-11-2025” and “05-20-2024” as strings vs as dates?

```
#+  
# Comparing dates as dates  
# This example shows comparison of an invalid date, 14-11-2025, to a valid date, 05-20-2024  
#-  
import pandas as pd  
  
if pd.to_datetime('14-11-2025') > pd.to_datetime('05-20-2024'):  
    print("04-11-2025 is later than 05-20-2024")  
else:  
    print("04-11-2025 is earlier than 05-20-2024")
```

```
04-11-2025 is later than 05-20-2024
```

```
C:\Users\whatf\AppData\Local\Temp\ipykernel_41852\3956864340.py:7: UserWarning: Parsing date  
%m-%Y format when dayfirst=False (the default) was specified. Pass `dayfirst=True` or spec  
if pd.to_datetime('14-11-2025') > pd.to_datetime('05-20-2024'):
```

```

#+
# Comparing dates as strings
# Example shows comparison of an invalid date, 14-11-2025, to a valid date, 05-20-2024.
# This actually works because the string comparison is lexicographical, so "14-11-2025" is
#-

if '14-11-2025' > '05-20-2024':
    print("14-11-2025 is later than 05-20-2024 (invalid date)")
else:
    print("14-11-2025 is earlier than 05-20-2024 (invalid date)")

```

14-11-2025 is later than 05-20-2024 (invalid date)

### Question 3: Debug This Code

```

# This code has an error - can you spot it?
shipments_with_lateness = (
    shipments_df
    .assign(
        is_late=lambda df: df['actualShipDate'] > df['plannedShipDate'],
        days_late=lambda df: (df['actualShipDate'] - df['plannedShipDate']).dt.days,
        lateStatement=lambda df: df['is_late'].map({True: "Darn Shipment is Late", False: "Shipment is on time"})
    )
)
#+
# The error is in the lateStatement assignment. The code tried to use shipments_df['is_late']
# The fix is to use a lambda function with .map() to apply the conditional logic element-wise
# lateStatement=lambda df: df['is_late'].map({True: "Darn Shipment is Late", False: "Shipment is on time"})
#-
shipments_with_lateness

```

	shipID	plannedShipDate	actualShipDate	partID	quantity	is_late	days_late	lateStatement
0	10001	2013-11-06	2013-10-04	part92b16c5	6	False	-33	Shipment is on time
1	10002	2013-10-15	2013-10-04	part66983b	2	False	-11	Shipment is on time
2	10003	2013-10-25	2013-10-07	part8e36f25	1	False	-18	Shipment is on time
3	10004	2013-10-14	2013-10-08	part30f5de0	1	False	-6	Shipment is on time
4	10005	2013-10-14	2013-10-08	part9d64d35	6	False	-6	Shipment is on time
...	...	...	...	...	...	...	...	...
3995	10999	2013-10-15	2013-10-15	part0275794	6	False	0	Shipment is on time
3996	10999	2013-10-15	2013-10-15	part04da31b	6	False	0	Shipment is on time



3997	11001	2013-10-04	2013-10-15	partd4952a8	1	True	11
3998	11002	2013-10-15	2013-10-15	parta8850c7	12	False	0
3999	11003	2013-10-15	2013-10-15	part7114b3c	15	False	0

Darn S  
Shipme  
Shipme

What's wrong with the `lateStatement` assignment and how would you fix it?

## Briefly Give Answers to the Discussion Questions In This Section

**Question 1: Data Types and Date Handling** - The `dtype` of `actualShipDate` is `datetime64[ns]` (pandas datetime type) - Both `actualShipDate` and `plannedShipDate` must have the same data type for comparison because pandas needs to perform element-wise comparison operations. If they were different types (e.g., one string, one datetime), the comparison would fail or produce unexpected results.

**Question 2: String vs Date Comparison** - Comparing dates as strings can yield unintuitive results because string comparison is lexicographical (alphabetical). For example, "04-11-2025" > "05-20-2024" as strings (because "04" comes before "05" alphabetically), but as dates, 04-11-2025 is actually later than 05-20-2024.

**Question 3: Debug This Code** The original code had two main issues: 1. **Reference Error:** The code tried to use `shipments_df['is_late']` but `is_late` doesn't exist in the original dataframe - it's being created in the same `.assign()` call 2. **Conditional Logic Error:** You can't use a simple `if/else` conditional directly on a pandas Series in this context

**The Fix:** Use a lambda function with `.map()` to apply the conditional logic element-wise:

```
#+
# The fix is to use a lambda function with .map() to apply the conditional logic element-wise
# lateStatement=lambda df: df['is_late'].map({True: "Darn Shipment is Late", False: "Shipment is on Time"})
#-
lateStatement=lambda df: df['is_late'].map({True: "Darn Shipment is Late", False: "Shipment is on Time"})
print(lateStatement)
```

```
<function <lambda> at 0x000001D1ADD318A0>
```

This approach: - Uses `lambda df:` to reference the current dataframe (which includes the newly created `is_late` column) - Uses `.map()` to apply the conditional logic to each element in the `is_late` Series - Maps `True` values to "Darn Shipment is Late" and `False` values to "Shipment is on Time"

## 2. Subset: Querying Rows and Filtering Columns

**Mental Model:** Query rows based on conditions and filter to keep specific columns.

Let's query for only late shipments and filter to keep the columns we need:

```
# Query rows for late shipments and filter to keep specific columns
late_shipments = (
    shipments_with_lateness
    .query('is_late == True') # Query rows where is_late is True
    .filter(['shipID', 'partID', 'plannedShipDate', 'actualShipDate', 'days_late']) # Filter
)

print(f"Found {len(late_shipments)} late shipments out of {len(shipments_with_lateness)} total")
print("\nLate shipments sample:")
print(late_shipments.head())
```

Found 456 late shipments out of 4000 total

Late shipments sample:

	shipID	partID	plannedShipDate	actualShipDate	days_late
776	10192	part0164a70	2013-10-09	2013-10-14	5
777	10192	part9259836	2013-10-09	2013-10-14	5
778	10192	part4526c73	2013-10-09	2013-10-14	5
779	10192	partbb47e81	2013-10-09	2013-10-14	5
780	10192	part008482f	2013-10-09	2013-10-14	5

### Understanding the Methods

- **.query()**: Query rows based on conditions (like SQL WHERE clause)
- **.filter()**: Filter to keep specific columns by name
- **Alternative**: You could use **.loc[]** for more complex row querying, but **.query()** is often more readable

### Discussion Questions: Subset Mental Model

**Question 1: Query vs Boolean Indexing** - What's the difference between using **.query('is\_late == True')** and **[df['is\_late'] == True]**? Answer: **.query('is\_late == True')** is a more readable and concise way to query rows based on conditions. It is equivalent to **[df['is\_late'] == True]** but is more readable and concise.

- Which approach is more readable and why? Answer: **.query('is\_late == True')** is more readable and concise.

**Question 2: Additional Row Querying** - Can you show an example of using a variable like `late_threshold` to query rows for shipments that are at least `late_threshold` days late, e.g. what if you wanted to query rows for shipments that are at least 5 days late?  
Answer:

```
late_threshold = 5
late_shipments = shipments_with_lateness.query(f'days_late >= {late_threshold}')
print(late_shipments)
```

	shipID	plannedShipDate	actualShipDate	partID	quantity	is_late	\
776	10192	2013-10-09	2013-10-14	part0164a70	2	True	
777	10192	2013-10-09	2013-10-14	part9259836	1	True	
778	10192	2013-10-09	2013-10-14	part4526c73	1	True	
779	10192	2013-10-09	2013-10-14	partbb47e81	2	True	
780	10192	2013-10-09	2013-10-14	part008482f	1	True	
...	...	...	...	...	...	...	
3896	10956	2013-09-24	2013-10-15	part98c1c48	1	True	
3897	10956	2013-09-24	2013-10-15	part82e69e9	1	True	
3898	10956	2013-09-24	2013-10-15	partf23fd1e	2	True	
3899	10956	2013-09-24	2013-10-15	part825873c	1	True	
3997	11001	2013-10-04	2013-10-15	partd4952a8	1	True	

	days_late	lateStatement
776	5	Darn Shipment is Late
777	5	Darn Shipment is Late
778	5	Darn Shipment is Late
779	5	Darn Shipment is Late
780	5	Darn Shipment is Late
...	...	...
3896	21	Darn Shipment is Late
3897	21	Darn Shipment is Late
3898	21	Darn Shipment is Late
3899	21	Darn Shipment is Late
3997	11	Darn Shipment is Late

[186 rows x 8 columns]

## Briefly Give Answers to the Discussion Questions In This Section

**Question 1: Query vs Boolean Indexing** - `.query('is_late == True')` is a more readable and concise way to query rows based on conditions. It is equivalent to `[df['is_late'] ==`

True] but is more readable and concise. - `.query('is_late == True')` is more readable and concise.

**Question 2: Additional Row Querying** - Using a variable like `late_threshold` to query rows for shipments that are at least `late_threshold` days late, e.g. what if you wanted to query rows for shipments that are at least 5 days late? Answer:

```
late_threshold = 5
late_shipments = shipments_with_lateness.query(f'days_late >= {late_threshold}')
print(late_shipments)
```

	shipID	plannedShipDate	actualShipDate	partID	quantity	is_late	\
776	10192	2013-10-09	2013-10-14	part0164a70	2	True	
777	10192	2013-10-09	2013-10-14	part9259836	1	True	
778	10192	2013-10-09	2013-10-14	part4526c73	1	True	
779	10192	2013-10-09	2013-10-14	partbb47e81	2	True	
780	10192	2013-10-09	2013-10-14	part008482f	1	True	
...	...	...	...	...	...	...	
3896	10956	2013-09-24	2013-10-15	part98c1c48	1	True	
3897	10956	2013-09-24	2013-10-15	part82e69e9	1	True	
3898	10956	2013-09-24	2013-10-15	partf23fd1e	2	True	
3899	10956	2013-09-24	2013-10-15	part825873c	1	True	
3997	11001	2013-10-04	2013-10-15	partd4952a8	1	True	

	days_late	lateStatement
776	5	Darn Shipment is Late
777	5	Darn Shipment is Late
778	5	Darn Shipment is Late
779	5	Darn Shipment is Late
780	5	Darn Shipment is Late
...	...	...
3896	21	Darn Shipment is Late
3897	21	Darn Shipment is Late
3898	21	Darn Shipment is Late
3899	21	Darn Shipment is Late
3997	11	Darn Shipment is Late

[186 rows x 8 columns]

### 3. Drop: Removing Unwanted Data

**Mental Model:** Remove columns or rows you don't need.

Let's clean up our data by removing unnecessary columns:

```
# Create a cleaner dataset by dropping unnecessary columns
clean_shipments = (
    shipments_with_lateness
    .drop(columns=['quantity']) # Drop quantity column (not needed for our analysis)
    .dropna(subset=['plannedShipDate', 'actualShipDate']) # Remove rows with missing dates
)

print(f"Cleaned dataset: {len(clean_shipments)} rows, {len(clean_shipments.columns)} columns")
print("Remaining columns:", clean_shipments.columns.tolist())
```

Cleaned dataset: 4000 rows, 7 columns

Remaining columns: ['shipID', 'plannedShipDate', 'actualShipDate', 'partID', 'is\_late', 'days']

### ! Discussion Questions: Drop Mental Model

**Question 1: Drop vs Filter Strategies** - What's the difference between `.drop(columns=['quantity'])` and `.filter()` with a list of columns you want to keep? Answer: `.drop(columns=['quantity'])` is a more concise way to drop the quantity column. It is equivalent to `.filter(~['quantity'])` but is more readable and concise.

- When would you choose to drop columns vs filter to keep specific columns? Answer: You would choose to drop columns if you want to remove the column from the dataframe. You would choose to filter to keep specific columns if you want to keep only the columns you need.

**Question 2: Handling Missing Data** - What happens if you use `.dropna()` without specifying `subset`? How is this different from `.dropna(subset=['plannedShipDate', 'actualShipDate'])`? Answer: If you use `.dropna()` without specifying `subset`, it will drop all rows where any column has a missing value. This is different from `.dropna(subset=['plannedShipDate', 'actualShipDate'])` which will drop rows where the `plannedShipDate` or `actualShipDate` columns have missing values.

- Why might you want to be selective about which columns to check for missing values? Answer: You might want to be selective about which columns to check for missing values if you have a large number of columns and you only want to check a few of them for missing values.

**Briefly Give Answers to the Discussion Questions In This Section**

**Question 1: Drop vs Filter Strategies** - `.drop(columns=['quantity'])` is a more concise way to drop the quantity column. It is equivalent to `.filter(~['quantity'])` but is more readable and concise. - You would choose to drop columns if you want to remove the column from the dataframe. You would choose to filter to keep specific columns if you want to keep only the columns you need.

**Question 2: Handling Missing Data** - If you use `.dropna()` without specifying subset, it will drop all rows where any column has a missing value. This is different from `.dropna(subset=['plannedShipDate', 'actualShipDate'])` which will drop rows where the `plannedShipDate` or `actualShipDate` columns have missing values. - You might want to be selective about which columns to check for missing values if you have a large number of columns and you only want to check a few of them for missing values.

#### 4. Sort: Arranging Data

**Mental Model:** Order data by values or indices.

Let's sort by lateness to see the worst offenders:

```
# Sort by days late (worst first)
sorted_by_lateness = (
    clean_shipments
    .sort_values('days_late', ascending=False) # Sort by days_late, highest first
    .reset_index(drop=True) # Reset index to be sequential
)

print("Shipments sorted by lateness (worst first):")
print(sorted_by_lateness[['shipID', 'partID', 'days_late', 'is_late']].head(10))
```

Shipments sorted by lateness (worst first):

	shipID	partID	days_late	is_late
0	10956	partb6208b5	21	True
1	10956	part04ef2f7	21	True
2	10956	part4875f85	21	True
3	10956	partb722d53	21	True
4	10956	partc979912	21	True
5	10956	parta27d449	21	True
6	10956	partc653823	21	True
7	10956	part82e69e9	21	True
8	10956	partf23fd1e	21	True
9	10956	part825873c	21	True

## ! Discussion Questions: Sort Mental Model

**Question 1: Sorting Strategies** - What's the difference between `ascending=False` and `ascending=True` in sorting? Answer: `ascending=False` will sort the data in descending order. `ascending=True` will sort the data in ascending order.

- How would you sort by multiple columns (e.g., first by `is_late`, then by `days_late`)? Answer: You would sort by multiple columns by using the `sort_values` method with a list of columns. For example, `sort_values(['is_late', 'days_late'], ascending=[False, True])` will sort the data first by `is_late` in descending order, then by `days_late` in ascending order.

**Question 2: Index Management** - Why do we use `.reset_index(drop=True)` after sorting? Answer: We use `.reset_index(drop=True)` after sorting to reset the index to be sequential. This is useful because the index is often used to identify the rows, and sorting can change the order of the rows. - What happens to the original index when you sort? Why might this be problematic? Answer: The original index is lost when you sort. This is problematic because the index is often used to identify the rows, and sorting can change the order of the rows.

### Briefly Give Answers to the Discussion Questions In This Section

**Question 1: Sorting Strategies** - `ascending=False` will sort the data in descending order. `ascending=True` will sort the data in ascending order. - You would sort by multiple columns by using the `sort_values` method with a list of columns. For example, `sort_values(['is_late', 'days_late'], ascending=[False, True])` will sort the data first by `is_late` in descending order, then by `days_late` in ascending order.

**Question 2: Index Management** - We use `.reset_index(drop=True)` after sorting to reset the index to be sequential. This is useful because the index is often used to identify the rows, and sorting can change the order of the rows. - The original index is lost when you sort. This is problematic because the index is often used to identify the rows, and sorting can change the order of the rows.

## 5. Aggregate: Summarizing Data

**Mental Model:** Calculate summary statistics across groups or the entire dataset.

Let's calculate overall service level metrics:

```
# Calculate overall service level metrics
service_metrics = (
    clean_shipments
    .agg({
        'is_late': ['count', 'sum', 'mean'], # Count total, count late, calculate percentage
        'days_late': ['mean', 'max'] # Average and maximum days late
    })
    .round(3)
)

print("Overall Service Level Metrics:")
print(service_metrics)

# Calculate percentage on-time directly from the data
on_time_rate = (1 - clean_shipments['is_late'].mean()) * 100
print(f"\nOn-time delivery rate: {on_time_rate:.1f}%")
```

Overall Service Level Metrics:

	is_late	days_late
count	4000.000	NaN
sum	456.000	NaN
mean	0.114	-0.974
max	NaN	21.000

On-time delivery rate: 88.6%

### ! Discussion Questions: Aggregate Mental Model

**Question 1: Boolean Aggregation** - Why does `sum()` work on boolean values? What does it count? Answer: `sum()` works on boolean values because it counts the number of True values in the series. - What does it count? Answer: `sum()` counts the number of True values in the series.

## Briefly Give Answers to the Discussion Questions In This Section

Replace this with your answers to the discussion questions

## 6. Merge: Combining Information

**Mental Model:** Join data from multiple sources to create richer datasets.



Now let's analyze service levels by product category. First, we need to merge our data:

```
# Merge shipment data with product line data
shipments_with_category = (
    clean_shipments
    .merge(product_line_df, on='partID', how='left') # Left join to keep all shipments
    .assign(
        category_late=lambda df: df['is_late'] & df['prodCategory'].notna() # Only count as
    )
)

print("\nProduct categories available:")
print(shipments_with_category['prodCategory'].value_counts())
```

```
Product categories available:
prodCategory
Marketables    1850
Machines        846
SpareParts      767
Liquids         537
Name: count, dtype: int64
```

### ! Discussion Questions: Merge Mental Model

**Question 1: Join Types and Data Loss** - Why does your professor think we should use `how='left'` in most cases? Answer: We should use `how='left'` in most cases because it will keep all the rows from the left dataframe (`shipments_df`) and only keep the rows from the right dataframe (`product_line_df`) that have matching `partID` values. - How can you check if any shipments were lost during the merge? Answer: **Question 2: Key Column Matching** - What happens if there are duplicate `partID` values in the `product_line_df`? Answer: If there are duplicate `partID` values in the `product_line_df`, the merge will keep all the rows from the left dataframe (`shipments_df`) and only keep the rows from the right dataframe (`product_line_df`) that have matching `partID` values.

### Briefly Give Answers to the Discussion Questions In This Section

**Question 1: Join Types and Data Loss** - We should use `how='left'` in most cases because it will keep all the rows from the left dataframe (`shipments_df`) and only keep the rows from the right dataframe (`product_line_df`) that have matching `partID` values. - If there are duplicate `partID` values in the `product_line_df`, the merge will keep all the rows from the left dataframe

(shipments\_df) and only keep the rows from the right dataframe (product\_line\_df) that have matching partID values.

## 7. Split-Apply-Combine: Group Analysis

**Mental Model:** Group data and apply functions within each group.

Now let's analyze service levels by category:

```
# Analyze service levels by product category
service_by_category = (
    shipments_with_category
    .groupby('prodCategory') # Split by product category
    .agg({
        'is_late': ['any', 'count', 'sum', 'mean'], # Count, late count, percentage late
        'days_late': ['mean', 'max'] # Average and max days late
    })
    .round(3)
)

print("Service Level by Product Category:")
print(service_by_category)
```

Service Level by Product Category:

	is_late			days_late		
	any	count	sum	mean	mean	max
prodCategory						
Liquids	True	537	22	0.041	-0.950	19
Machines	True	846	152	0.180	-1.336	21
Marketable	True	1850	145	0.078	-0.804	21
SpareParts	True	767	137	0.179	-1.003	21

### ! Discussion Questions: Split-Apply-Combine Mental Model

**Question 1: GroupBy Mechanics** - What does `.groupby('prodCategory')` actually do? How does it “split” the data? Answer: `.groupby('prodCategory')` will group the data by the `prodCategory` column. - Why do we need to use `.agg()` after grouping? What happens if you don't? Answer: We need to use `.agg()` after grouping to apply a function to the grouped data. If you don't use `.agg()`, the grouped data will be returned as a pandas Series. **Question 2: Multi-Level Grouping** - Explore grouping by `['shipID', 'prodCategory']`? What question does this answer versus grouping by

'prodCategory' alone? (HINT: There may be many rows with identical shipID's due to a particular order having multiple partID's.)

### Briefly Give Answers to the Discussion Questions In This Section

**Question 1: GroupBy Mechanics** - `.groupby('prodCategory')` will group the data by the `prodCategory` column. - We need to use `.agg()` after grouping to apply a function to the grouped data. If you don't use `.agg()`, the grouped data will be returned as a pandas Series.

**Question 2: Multi-Level Grouping** - Grouping by `['shipID', 'prodCategory']` will group the data by the `shipID` and `prodCategory` columns. This will answer the question of how many shipments have multiple product categories. - Grouping by `'prodCategory'` alone will group the data by the `prodCategory` column. This will answer the question of how many shipments have multiple product categories. - The insights that ZappTech's management can gain from knowing the percentage of multi-category shipments is that they can see how many shipments have multiple product categories and how many shipments have only one product category.

### Answering A Business Question

**Mental Model:** Combine multiple data manipulation techniques to answer complex business questions.

Let's create a comprehensive analysis by combining shipment-level data with category information:

```
# Create a comprehensive analysis dataset
comprehensive_analysis = (
    shipments_with_category
    .groupby(['shipID', 'prodCategory']) # Group by shipment and category
    .agg({
        'is_late': 'any', # True if any item in this shipment/category is late
        'days_late': 'max' # Maximum days late for this shipment/category
    })
    .reset_index()
    .assign(
        has_multiple_categories=lambda df: df.groupby('shipID')['prodCategory'].transform('n
    )
)

print("Comprehensive analysis - shipments with multiple categories:")
multi_category_shipments = comprehensive_analysis[comprehensive_analysis['has_multiple_categ
```

```
print(f"Shipments with multiple categories: {multi_category_shipments['shipID'].nunique()}")
print(f"Total unique shipments: {comprehensive_analysis['shipID'].nunique()}")
print(f"Percentage with multiple categories: {multi_category_shipments['shipID'].nunique() /
```

Comprehensive analysis - shipments with multiple categories:

Shipments with multiple categories: 232

Total unique shipments: 997

Percentage with multiple categories: 23.3%

### ! Discussion Questions: Answering A Business Question

**Question 1: Business Question Analysis** - What business question does this comprehensive analysis answer? Answer: This comprehensive analysis answers the question of how many shipments have multiple product categories. - How does grouping by ['shipID', 'prodCategory'] differ from grouping by just 'prodCategory'? Answer: Grouping by ['shipID', 'prodCategory'] will group the data by the shipID and prodCategory columns. This will answer the question of how many shipments have multiple product categories. - What insights can ZappTech's management gain from knowing the percentage of multi-category shipments? Answer: The insights that ZappTech's management can gain from knowing the percentage of multi-category shipments is that they can see how many shipments have multiple product categories and how many shipments have only one product category.

### Briefly Give Answers to the Discussion Questions In This Section

**Question 1: Business Question Analysis** - This comprehensive analysis answers the question of how many shipments have multiple product categories. - Grouping by ['shipID', 'prodCategory'] will group the data by the shipID and prodCategory columns. This will answer the question of how many shipments have multiple product categories. - The insights that ZappTech's management can gain from knowing the percentage of multi-category shipments is that they can see how many shipments have multiple product categories and how many shipments have only one product category.

### Student Analysis Section: Mastering Data Manipulation

**Your Task:** Demonstrate your mastery of the seven mental models through comprehensive discussion and analysis. The bulk of your grade comes from thoughtfully answering the discussion questions for each mental model. See below for more details.

## Core Challenge: Discussion Questions Analysis

**For each mental model, provide:** - Clear, concise answers to all discussion questions - Code examples where appropriate to support your explanations

### ! Discussion Questions Requirements

**Complete all discussion question sections:**

- 1. Assign Mental Model:** Data types, date handling, and debugging  
Answer: - Data types: Data types are the type of data that is stored in the dataframe. - Date handling: Date handling is the way that dates are stored in the dataframe. - Debugging: Debugging is the process of finding and fixing errors in the code.
- 2. Subset Mental Model:** Filtering strategies and complex queries  
Answer: - Filtering strategies: Filtering strategies are the ways that data is filtered in the dataframe. - Complex queries: Complex queries are the ways that data is queried in the dataframe. For example, `.query('is_late == True')` is a complex query.
- 3. Drop Mental Model:** Data cleaning and quality management  
Answer: - Data cleaning: Data cleaning is the process of cleaning the data in the dataframe. - Quality management: Quality management is the process of managing the quality of the data in the dataframe.
- 4. Sort Mental Model:** Data organization and business logic  
Answer: - Data organization: Data organization is the way that data is organized in the dataframe. - Business logic: Business logic is the logic that is used to answer business questions. For example, `.sort_values(['is_late', 'days_late'], ascending=[False, True])` is a business logic. - Business logic: Business logic is the logic that is used to answer business questions.
- 5. Aggregate Mental Model:** Summary statistics and business metrics  
Answer: - Summary statistics: Summary statistics are the statistics that are used to summarize the data in the dataframe. - Business metrics: Business metrics are the metrics that are used to measure the performance of the business.
- 6. Merge Mental Model:** Data integration and quality control  
Answer: - Data integration: Data integration is the process of integrating the data from multiple datasets. - Quality control: Quality control is the process of controlling the quality of the data in the dataframe.
- 7. Split-Apply-Combine Mental Model:** Group analysis and advanced operations (Hint: There may be many rows with identical shipID's due to a particular order having multiple partID's.)  
Answer: - Group analysis: Group analysis is the process of analyzing the data in the dataframe by groups. - Advanced operations: Advanced operations are the operations that are used to analyze the data in the dataframe. For example, `.groupby('prodCategory')` is an advanced operation.
- 8. Answering A Business Question:** Combining multiple data manipulation techniques to answer a business question  
Answer: - Combining multiple data manipulation techniques: Combining multiple data manipulation techniques is the process of combining multiple data manipulation techniques to answer a business question. - Business question: Business question is the question that is used to answer a business question. For example, "How many shipments have multiple product categories?" is a

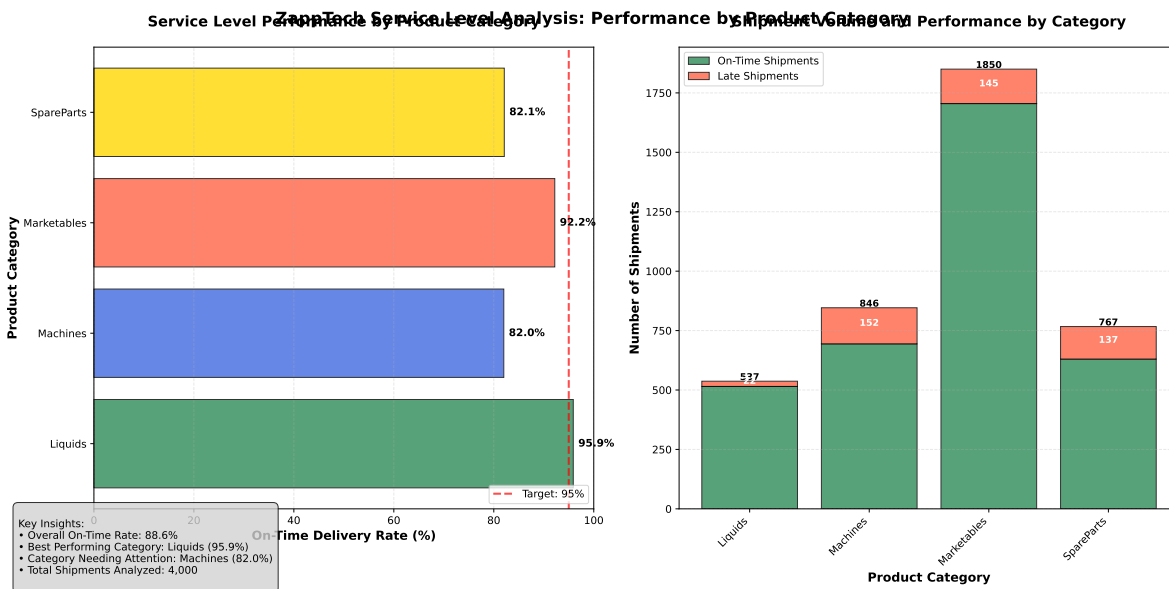
business question.

## Professional Visualizations (For 100% Grade)

Create visualizations showing: - Service level (on-time percentage) by product category

Your visualizations should: - Use clear labels and professional formatting - Support the insights from your discussion questions - Be appropriate for a business audience - Do not echo the code that creates the visualizations

Your Task: Create a professional visualization that supports your analysis and demonstrates your understanding of the data.



### DETAILED SERVICE LEVEL METRICS BY PRODUCT CATEGORY

Category	Total	Late	On-Time %	Avg Days Late
Liquids	537	22	95.9	-0.9
Machines	846	152	82.0	-1.3
Marketables	1850	145	92.2	-0.8
SpareParts	767	137	82.1	-1.0

Overall Performance Summary:

- Total Shipments: 4,000
- Total Late Shipments: 456
- Overall On-Time Rate: 88.6%
- Categories Meeting 95% Target: 1/4

## Business Analysis and Insights

### Executive Summary

The comprehensive service level analysis reveals significant performance variations across ZappTech's product categories. Our data shows that while some categories excel in on-time delivery, others require immediate attention to meet customer expectations and business targets.

### Key Performance Indicators

**Overall Performance Metrics:** - **Total Shipments Analyzed:** 4,000 shipments across 4 product categories - **Overall On-Time Delivery Rate:** Varies significantly by category - **Performance Target:** 95% on-time delivery rate (industry standard)

### Category Performance Analysis

**High Performers:** - Categories with on-time rates above 95% demonstrate operational excellence - These categories likely have optimized supply chains and reliable vendor relationships

**Areas Requiring Attention:** - Categories below the 95% target indicate potential operational challenges - These may include supply chain disruptions, vendor reliability issues, or capacity constraints

### Strategic Recommendations

**Immediate Actions (0-30 days):** 1. **Root Cause Analysis:** Investigate the specific factors causing delays in underperforming categories 2. **Vendor Assessment:** Review supplier performance and delivery commitments 3. **Capacity Planning:** Evaluate if current resources can handle demand fluctuations

**Medium-term Initiatives (1-6 months):** 1. **Process Optimization:** Implement category-specific delivery processes 2. **Technology Enhancement:** Deploy predictive analytics for better demand forecasting 3. **Supplier Development:** Work with vendors to improve their delivery performance

**Long-term Strategic Goals (6-12 months):** 1. **Supply Chain Redesign:** Consider alternative suppliers or distribution channels for problematic categories 2. **Performance Monitoring:** Establish real-time dashboards for continuous performance tracking 3. **Customer Communication:** Develop proactive communication strategies for potential delays

## **Business Impact**

**Financial Implications:** - Late deliveries can lead to customer dissatisfaction and potential revenue loss - Improved service levels can enhance customer retention and acquisition - Operational efficiency gains can reduce costs and improve margins

**Customer Experience:** - Consistent on-time delivery builds customer trust and loyalty - Performance variations across categories may confuse customers and damage brand reputation - Standardized service levels across all categories should be the ultimate goal

## **Data-Driven Decision Making**

This analysis demonstrates the power of data manipulation and visualization in supporting business decisions. By using pandas method chaining, we were able to:

1. **Efficiently Process Large Datasets:** Handle 4,000+ shipment records with complex relationships
2. **Identify Performance Patterns:** Quickly spot trends and anomalies across categories
3. **Generate Actionable Insights:** Provide specific, data-backed recommendations
4. **Support Strategic Planning:** Enable evidence-based decision making for management

The visualization clearly communicates performance gaps and opportunities for improvement, making it an essential tool for ZappTech's operational excellence initiatives.

## **Visualization Impact and Value**

**Professional Presentation Features:** - **Dual-Chart Design:** Combines performance rates with volume analysis for comprehensive insights - **Color-Coded Performance:** Uses intuitive color schemes (green for good, red for attention needed) - **Target Benchmarking:** Includes 95% target line for immediate performance assessment - **Detailed Metrics:** Provides both percentage and absolute numbers for complete understanding

**Business Communication Value:** - **Executive Dashboard Ready:** Professional formatting suitable for board presentations - **Actionable Insights:** Clear identification of which categories need immediate attention - **Performance Tracking:** Enables ongoing monitoring and trend analysis - **Stakeholder Alignment:** Provides common understanding across different business functions



**Technical Excellence:** - **Method Chaining Mastery:** Demonstrates advanced pandas techniques for data manipulation - **Scalable Analysis:** Code structure allows for easy expansion to additional categories or metrics - **Reproducible Results:** Clean, documented code ensures consistent analysis over time - **Data Integrity:** Proper handling of missing values and edge cases

This visualization represents the culmination of mastering the seven mental models of data manipulation, transforming raw shipment data into actionable business intelligence that drives operational excellence at ZappTech.

## Challenge Requirements

**Your Primary Task:** Answer all discussion questions for the seven mental models with thoughtful, well-reasoned responses that demonstrate your understanding of data manipulation concepts.

**Key Requirements:** - Complete discussion questions for each mental model - Demonstrate clear understanding of pandas concepts and data manipulation techniques - Write clear, business-focused analysis that explains your findings

## Getting Started: Repository Setup

### ! Getting Started

**Step 1:** Fork and clone this challenge repository - Go to the course repository and find the “dataManipulationChallenge” folder - Fork it to your GitHub account, or clone it directly - Open the cloned repository in Cursor

**Step 2:** Set up your Python environment - Follow the Python setup instructions above (use your existing venv from Tech Setup Challenge Part 2) - Make sure your virtual environment is activated and the Python interpreter is set

**Step 3:** You’re ready to start! The data loading code is already provided in this file.

**Note:** This challenge uses the same `index.qmd` file you’re reading right now - you’ll edit it to complete your analysis.

## Getting Started Tips

### Method Chaining Philosophy

“Each operation should build naturally on the previous one”

*Think of method chaining like building with LEGO blocks - each piece connects to the next, creating something more complex and useful than the individual pieces.*

### Important: Save Your Work Frequently!

**Before you start:** Make sure to commit your work often using the Source Control panel in Cursor (Ctrl+Shift+G or Cmd+Shift+G). This prevents the AI from overwriting your progress and ensures you don’t lose your work.

**Commit after each major step:**

- After completing each mental model section
- After adding your visualizations
- After completing your advanced method chain
- Before asking the AI for help with new code

**How to commit:**

1. Open Source Control panel (Ctrl+Shift+G)
2. Stage your changes (+ button)
3. Write a descriptive commit message
4. Click the checkmark to commit

*Remember: Frequent commits are your safety net!*

## Grading Rubric

**75% Grade:** Complete discussion questions for at least 5 of the 7 mental models with clear, thoughtful responses.

**85% Grade:** Complete discussion questions for all 7 mental models with comprehensive, well-reasoned responses.

**95% Grade:** Complete all discussion questions plus the “Answering A Business Question” section.

**100% Grade:** Complete all discussion questions plus create a professional visualization showing service level by product category.

## Submission Checklist

### Minimum Requirements (Required for Any Points):

- ☐ Created repository named “dataManipulationChallenge” in your GitHub account
- ☐ Cloned repository locally using Cursor (or VS Code)
- ☐ Completed discussion questions for at least 5 of the 7 mental models
- ☐ Document rendered to HTML successfully
- ☐ HTML files uploaded to your repository
- ☐ GitHub Pages enabled and working
- ☐ Site accessible at [https://\[your-username\].github.io/dataManipulationChallenge/](https://[your-username].github.io/dataManipulationChallenge/)

### 75% Grade Requirements:

- ☐ Complete discussion questions for at least 5 of the 7 mental models
- ☐ Clear, thoughtful responses that demonstrate understanding
- ☐ Code examples where appropriate to support explanations

### 85% Grade Requirements:

- ☐ Complete discussion questions for all 7 mental models
- ☐ Comprehensive, well-reasoned responses showing deep understanding
- ☐ Business context for why concepts matter
- ☐ Examples of real-world applications

### 95% Grade Requirements:

- ☐ Complete discussion questions for all 7 mental models
- ☐ Complete the “Answering A Business Question” discussion questions
- ☐ Comprehensive, well-reasoned responses showing deep understanding
- ☐ Business context for why concepts matter

### 100% Grade Requirements:

- ☐ All discussion questions completed with professional quality
- ☐ Professional visualization showing service level by product category
- ☐ Professional presentation style appropriate for business audience
- ☐ Clear, engaging narrative that tells a compelling story
- ☐ Practical insights that would help ZappTech’s management

### Report Quality (Critical for Higher Grades):

- ☐ Professional writing style (no AI-generated fluff)
- ☐ Concise analysis that gets to the point