

Block 2: Essential Data Wrangling with NumPy & Pandas

Python Module for Incoming ISE & OR PhD Students

Will Kirschenman

August 7, 2025 | 10:00 AM - 10:50 AM

North Carolina State University

NC STATE UNIVERSITY

- **Goal:** Become familiar with the essential tools for data manipulation in Python
- **Duration:** 50 minutes of hands-on data wrangling
- **Format:** Presentation + interactive notebook exercises

What We'll Cover

NumPy arrays • Pandas DataFrames • Real-world data cleaning • PhD dataset analysis

Session Learning Objectives

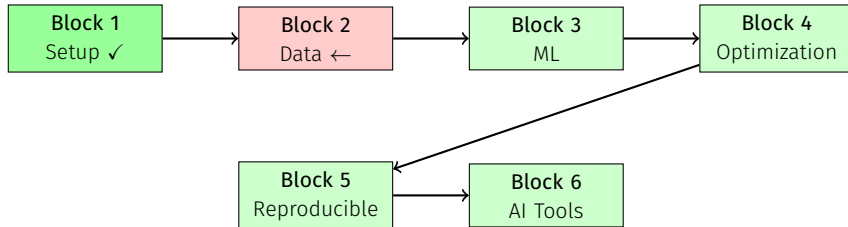
By the end of Block 2, you will:

1. Understand **NumPy arrays** for efficient numerical computation
2. Master **Pandas DataFrames** for data manipulation and analysis
3. Apply **data cleaning techniques** to messy real-world data
4. Create **new features** from existing data
5. Prepare a **clean dataset** ready for machine learning (Block 3)

Our Mission

Transform a messy PhD student research dataset into analysis-ready data!

Recap: Where We Are



From Block 1

You now know Python basics, Google Colab, and the programming mindset. Time to handle real data!

Our Dataset: PhD Student Research Productivity

What we're working with:

Dataset Features

- 280+ PhD students at NC State
- Research productivity metrics
- NC State specifics (Hunt Library, departments)
- Realistic but messy data

The Problems

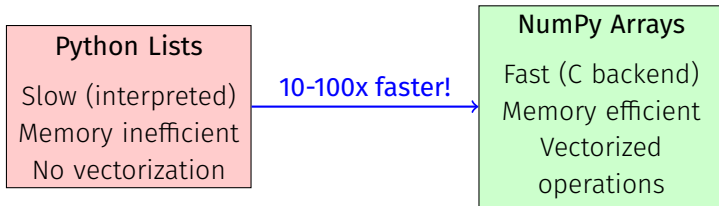
- Missing values
- Duplicate rows
- Inconsistent department names
- Extreme outliers
- Mixed data types

Real Data Reality

This is exactly the kind of messy data you'll encounter in your coursework!

NumPy Fundamentals

Why NumPy? The Performance Story



For OR/ISE Coursework

When you're processing coursework assignments and problem sets, speed matters!

Speed Comparison: Python vs NumPy

Let's see the difference with 1 million data points:

Python List Approach

```
1 python_list = list(range(1000000))
2
3 # Multiply each element by 2
4 start_time = time.time()
5 result = [x * 2 for x in python_list]
6 python_time = time.time() - start_time
7
8 # Result: ~0.15 seconds
```

NumPy Array Approach

```
1 numpy_array = np.array(range(1000000))
2
3 # Multiply entire array by 2
4 start_time = time.time()
5 result = numpy_array * 2
6 numpy_time = time.time() - start_time
7
8 # Result: ~0.005 seconds
```

NumPy is 30x faster in this example!

That's the difference between waiting 5 minutes vs 10 seconds for your analysis.

Creating NumPy Arrays

Multiple ways to create arrays for different use cases:

From Existing Data

```
1 # From Python list
2 coffee_data = [3.2, 4.1, 2.8, 5.0]
3 coffee_array = np.array(coffee_data)
4
5 # From nested lists (2D)
6 papers_by_year = [[1, 2], [2, 3], [1, 4]]
7 papers_matrix = np.array(papers_by_year)
```

Built-in Constructors

```
1 # Arrays of zeros/ones
2 zeros = np.zeros(10)
3 ones = np.ones(5)
4
5 # Range arrays
6 years = np.arange(1, 8) # 1 to 7
7
8 # Random arrays
9 stress = np.random.normal(5, 2, 100)
```

Coursework Applications

Initialize parameter arrays, create homework data, set up optimization variables

Vectorized Operations: The Magic

Apply operations to entire arrays at once:

```
1 # PhD student data
2 papers_per_year = np.array([0, 1, 2, 1, 3, 2, 1]) # 7 years of PhD
3 conference_costs = np.array([1200, 1500, 800, 2000, 1100])
4
5 # Mathematical operations on entire arrays
6 total_papers = np.sum(papers_per_year)           # 10 papers total
7 avg_papers = np.mean(papers_per_year)           # 1.43 papers/year
8 productivity_doubled = papers_per_year * 2       # [0, 2, 4, 2, 6, 4, 2]
9
10 # Statistical operations
11 print(f"Conference costs: {np.mean(conference_costs):.0f} ± {np.std(
    conference_costs):.0f}")
```

No Loops Needed!

NumPy handles the iteration internally in optimized C code.

Boolean Indexing: Powerful Filtering

Filter data based on conditions:

```
1 # Coffee consumption data
2 coffee_cups = np.array([2.1, 4.5, 8.2, 3.1, 12.5, 1.9, 6.8])
3
4 # Boolean conditions
5 high_caffeine = coffee_cups > 5.0
6 extreme_caffeine = coffee_cups > 10.0
7
8 # Filter the data
9 moderate_drinkers = coffee_cups[coffee_cups <= 5.0]
10 caffeine_addicts = coffee_cups[coffee_cups > 8.0]
11
12 # Count results
13 print(f"High_caffeine_consumers: {np.sum(high_caffeine)}")
```

Coursework Use Case

Filter homework results, identify outliers, select subsets for analysis

NumPy in Action: Statistical Analysis

Essential functions for data analysis:

Function	Purpose
<code>np.mean()</code> , <code>np.median()</code>	Central tendency
<code>np.std()</code> , <code>np.var()</code>	Variability
<code>np.min()</code> , <code>np.max()</code>	Range
<code>np.percentile()</code>	Quartiles and percentiles
<code>np.corrcoef()</code>	Correlation analysis
<code>np.unique()</code>	Count unique values
<code>np.where()</code>	Conditional selection

Real Research Example

Analyzing Hunt Library usage: mean=25.3 hours/week, std=12.1, 15% spend >40 hours/week

Pandas DataFrames

NumPy vs Pandas: Choose Your Tool

NumPy Arrays

- ✓ Homogeneous data
 - ✓ Mathematical operations
 - ✓ Maximum performance
 - ✓ Numerical computations

Simulations, algorithms, math

Pandas DataFrames

- ✓ Mixed data types
 - ✓ Labeled data (columns)
- ✓ Data cleaning tools
- ✓ Real-world datasets

Datasets, analysis, cleaning

Best of Both Worlds

Pandas is built on NumPy - you get speed + convenience!

Meet the DataFrame: Your Data Analysis Workhorse

DataFrame = Supercharged Spreadsheet

student_id	department	papers	coffee
PhD_001	ISE	3	4.2
PhD_002	OR	5	3.8
PhD_003	CSC	2	6.1
...

✓ Mixed data types ✓ Column labels ✓
Row indexing ✓ Missing value handling

Why DataFrames Matter

Real data is messy, mixed-type, and labeled - exactly what DataFrames handle best!

Essential DataFrame Exploration

First things first - understand your data:

Basic Inspection

```
1 # Load and examine
2 df = pd.read_csv('phd_data.csv')
3
4 # Quick overview
5 df.shape           # (280, 11)
6 df.head()          # First 5 rows
7 df.tail()           # Last 5 rows
8 df.info()           # Data types, nulls
9 df.describe()       # Statistics
```

Understanding Structure

```
1 # Column exploration
2 df.columns          # Column names
3 df.dtypes           # Data types
4 df.isnull().sum()   # Missing values
5
6 # Unique values
7 df['department'].unique()
8 df['department'].value_counts()
```

PhD Dataset

280 students, 11 variables, messy department names, missing coffee data (the horror!)

Data Selection: Getting What You Need

Multiple ways to slice and dice your data:

```
1 # Select columns
2 departments = df['department']           # Single column (Series)
3 basics = df[['student_id', 'department']] # Multiple columns (DataFrame)
4
5 # Filter rows
6 veterans = df[df['years_in_program'] >= 4] # Boolean filtering
7 ise_students = df[df['department'] == 'ISE'] # Exact match
8
9 # Complex conditions
10 productive_ise = df[(df['department'] == 'ISE') &
11                     (df['papers_published'] > 2)] # Multiple conditions
12
13 # Statistical filtering
14 high_stress = df[df['stress_level'] > df['stress_level'].quantile(0.75)]
```

Coursework Applications

Filter by assignment conditions, select participant subgroups, analyze specific time periods

GroupBy: The Power of Aggregation

Analyze data by groups - essential for coursework:

```
1 # Basic grouping
2 dept_analysis = df.groupby('department')['papers_published'].mean()
3
4 # Multiple aggregations
5 analysis = df.groupby('department').agg({
6     'papers_published': ['mean', 'std'],
7     'stress_level': 'mean',
8     'coffee_cups_per_day': 'mean'
9 })
10
11 # Custom productivity score
12 def productivity_score(group):
13     return (group['papers_published'] * 2 + group['conferences_attended']).mean()
```

Insight

OR students publish most papers (2.6 avg), but ISE students have better work-life balance!

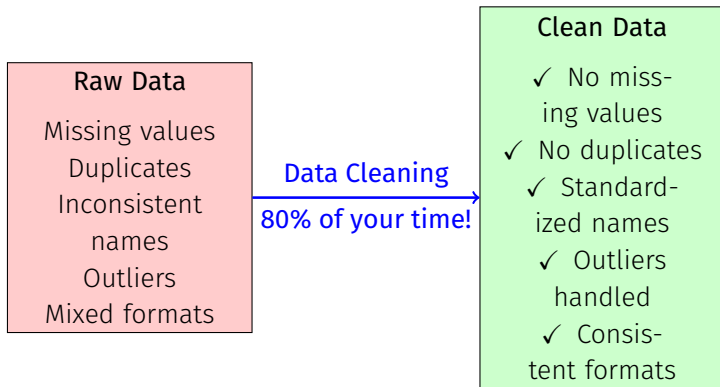
String Operations: Cleaning Text Data

Our department names are a mess - let's fix them:

```
1 # The problem
2 df['department'].unique()
3 # ['ISE', 'I.S.E.', 'Industrial Systems', 'OR', 'O.R.', ...]
4
5 # Cleaning function
6 def clean_department(dept):
7     if 'Industrial' in dept or dept in ['ISE', 'I.S.E.']:
8         return 'ISE'
9     elif 'Operations' in dept or dept in ['OR', 'O.R.']:
10        return 'OR'
11    # ... more cleaning logic
12
13 # Apply and verify
14 df['department_clean'] = df['department'].apply(clean_department)
15 df['department_clean'].value_counts()
```

Data Cleaning Workshop

The Reality of Data



The 80/20 Rule

Data scientists spend 80% of their time cleaning data, 20% analyzing it.

Step 1: Handle Missing Values Strategically

Different strategies for different data types:

```
1 # Assess the damage
2 missing_summary = df.isnull().sum()
3 # coffee_cups_per_day: 22 missing, stress_level: 14 missing, funding_amount: 8 missing
4
5 # Strategy 1: Fill with median (coffee - some students don't drink coffee)
6 coffee_median = df['coffee_cups_per_day'].median() # 3.5 cups
7 df['coffee_cups_per_day'].fillna(coffee_median, inplace=True)
8
9 # Strategy 2: Fill with mean (stress - continuous variable)
10 stress_mean = df['stress_level'].mean() # 5.8 out of 10
11 df['stress_level'].fillna(stress_mean, inplace=True)
12
13 # Strategy 3: Fill with group-specific values (funding by department)
14 funding_by_dept = df.groupby('department')['funding_amount'].mean()
15 df['funding_amount'] = df.apply(
16     lambda row: funding_by_dept[row['department']] if pd.isna(row['funding_amount'])
17     else row['funding_amount'], axis=1)
```

Step 2: Remove Duplicates

Duplicate data can skew your analysis:

```
1 # Check for duplicates
2 print(f"Duplicate_rows: {df.duplicated().sum()}") # 8 duplicates found
3
4 # Examine duplicates (optional)
5 duplicates = df[df.duplicated(keep=False)].sort_values('student_id')
6 print(duplicates[['student_id', 'department', 'papers_published']].head())
7
8 # Remove duplicates (keep first occurrence)
9 original_shape = df.shape
10 df.drop_duplicates(inplace=True)
11 new_shape = df.shape
12
13 print(f"Removed {original_shape[0] - new_shape[0]} duplicate_rows")
14 print(f"New dataset shape: {new_shape}")
```

Why This Matters

Duplicates can artificially inflate correlations and bias statistical tests.

Step 3: Handle Outliers Intelligently

Outliers: real extreme values or data entry errors?

```
1 # Identify outliers using IQR method
2 Q1 = df['coffee_cups_per_day'].quantile(0.25)    # 2.1 cups
3 Q3 = df['coffee_cups_per_day'].quantile(0.75)    # 4.8 cups
4 IQR = Q3 - Q1                                     # 2.7 cups
5
6 # Define outlier bounds
7 lower_bound = Q1 - 1.5 * IQR # Anything below -1.95 cups (impossible)
8 upper_bound = Q3 + 1.5 * IQR # Anything above 8.85 cups (concerning!)
9
10 # Find extreme outliers
11 extreme_coffee = df['coffee_cups_per_day'] > 10 # 5 students drinking 10+ cups/day
12
13 # Cap extreme values (medical safety!)
14 df.loc[df['coffee_cups_per_day'] > 8, 'coffee_cups_per_day'] = 8
15 print("Capped extreme coffee consumption at 8 cups/day for student safety")
```

Domain Knowledge Matters

Statistical outliers aren't always wrong - use research context to decide!

Step 4: Feature Engineering

Create new variables from existing data:

```
1 # Productivity score (papers worth 2x conferences)
2 df['productivity_score'] = (df['papers_published'] * 2 +
3                             df['conferences_attended']) / df['years_in_program']
4
5 # Work-life balance indicator
6 df['work_life_balance'] = 10 - (df['hours_in_hunt_library_per_week'] / 10 +
7                                 df['stress_level']) / 2
8 df['work_life_balance'] = df['work_life_balance'].clip(1, 10)
9
10 # Categorical features
11 df['seniority'] = pd.cut(df['years_in_program'],
12                          bins=[0, 2, 4, 10],
13                          labels=['Early', 'Mid', 'Advanced'])
14
15 df['caffeine_level'] = pd.cut(df['coffee_cups_per_day'],
16                               bins=[0, 2, 4, 8],
17                               labels=['Low', 'Moderate', 'High'])
```

Feature Engineering Results

Our new features reveal insights:

Seniority Level	Count	Avg Productivity
Early (1-2 years)	84	1.2
Mid (3-4 years)	126	2.1
Advanced (5+ years)	70	2.8

Caffeine Level	Count	Avg Stress
Low (<2 cups)	45	4.2
Moderate (2-4 cups)	156	5.8
High (4+ cups)	79	7.1

Insight

More coffee = more stress! But also, advanced students are most productive.

Advanced Operations & Wrap-up

Correlation Analysis: What Drives Success?

Key correlations with papers published:

Variable	Correlation
Years in program	0.68
Advisor meetings per month	0.42
Hours in Hunt Library	0.31
Coffee consumption	0.18
Funding amount	0.15
Stress level	0.09
Distance from campus	-0.12

Research Insights

Experience matters most, but regular advisor meetings and library time also boost productivity!

Data Visualization Preview

Papers by Department

Bar chart showing OR leads with 2.6 papers/student

Coffee vs Papers

Scatter plot reveals weak positive correlation

Experience vs Output

Clear upward trend: time leads to productivity

Stress Distribution

Normal-ish distribution centered around 5.8/10

Coming Up in Block 3

We'll use this clean data to build predictive models with machine learning!

What We Accomplished in Block 2 ✓

- ✓ Mastered NumPy for fast numerical computation
- ✓ Conquered Pandas for data manipulation and analysis
- ✓ Cleaned messy data using professional techniques
- ✓ Created new features from existing variables
- ✓ Prepared dataset ready for machine learning

Dataset Transformation

From 288 messy rows with missing values → 280 clean rows ready for analysis!

From Data to Insights: Predictive Modeling

What's Coming

- **Machine learning** fundamentals
- **Linear regression** with scikit-learn
- **Model evaluation** and interpretation
- **Predicting** PhD student success
- **Feature importance** analysis

Your Clean Data in Action

- Train models on our cleaned dataset
- Predict papers published
- Understand what drives success
- Build your first ML pipeline

10-minute break, then we build models!

Key Takeaways

Technical Skills

- **NumPy** for computational speed
- **Pandas** for data manipulation
- **Data cleaning** strategies
- **Feature engineering** techniques
- **Exploratory analysis** methods

Research Mindset

- Always explore data first
- Clean systematically
- Document your decisions
- Create meaningful features
- Validate your cleaning

Remember: Your dataset is saved as 'phd_research_productivity_clean.csv'
Ready for machine learning in Block 3!

Questions?

See you in 10 minutes for Block 3!