

# Complete Data Science Workflow Guide

## From Excel User to Python Pro

---

### The CRISP-DM Framework (Easy to Remember!)

**C.R.I.S.P-D.M** = Your roadmap to success

- Cusiness Understanding
  - Read Data (Data Understanding)
  - Investigate & Clean (Data Preparation)
  - Statistics & Models (Modeling)
  - Performance Check (Evaluation)
  - Deploy
  - Monitor
- 

### Phase 1: BUSINESS UNDERSTANDING (The "Why")

#### What Excel Users Miss:

-  Jump straight to pivot tables
-  You define clear objectives first

#### Your Actions:

1. **Define the problem:** What question are you answering?
2. **Set success metrics:** How will you measure success?
3. **Identify stakeholders:** Who needs this?

#### Python Advantage:

```
python
```

```
# Document everything in Jupyter Notebook  
# Markdown cells = your project documentation
```

## Time: 10% of project

---

### Phase 2: READ DATA (Data Understanding)

#### Core Libraries (The "Big 3"):

```
python

import pandas as pd      # Excel on steroids (PD)
import numpy as np        # Math powerhouse (NP)
import matplotlib.pyplot as plt # Basic plots (PLT)
```

#### Load Your Data:

```
python

# CSV files
df = pd.read_csv('data.csv')

# Excel files (supports multiple sheets!)
df = pd.read_excel('data.xlsx', sheet_name='Sheet1')

# SQL databases
import sqlite3
conn = sqlite3.connect('database.db')
df = pd.read_sql_query("SELECT * FROM table", conn)

# Multiple files at once (Excel can't do this!)
import glob
files = glob.glob('data/*.csv')
df = pd.concat([pd.read_csv(f) for f in files])
```

#### First Look (EDA = Exploratory Data Analysis):

```
python
```

```
# The Essential 5 Commands (memorize these!)
```

```
df.head()      # First 5 rows
```

```
df.info()      # Data types & missing values
```

```
df.describe()  # Statistics for numbers
```

```
df.shape       # (rows, columns)
```

```
df.columns     # Column names
```

## What Outshines Excel:

- Handle millions of rows (Excel max: ~1M)
- See all data types instantly
- Automated statistical summary

Time: 15% of project

---

## ✍ Phase 3: INVESTIGATE & CLEAN (Data Preparation)

### The Cleaning Checklist (D.M.T - Don't Miss These!):

- Duplicates
- Missing values
- Type conversions

### Essential Libraries:

```
python
```

```
import pandas as pd
```

```
import numpy as np
```

```
from datetime import datetime
```

### 1. Handle Missing Values:

```
python
```

```

# Check missing data
df.isnull().sum()

# Visualize missing data (Excel can't do this!)
import missingno as msno
msno.matrix(df)

# Fix missing values
df['column'].fillna(df['column'].mean(), inplace=True) # Mean
df['column'].fillna(df['column'].median(), inplace=True) # Median
df['column'].fillna('Unknown', inplace=True) # Text
df.dropna() # Remove rows with any missing values

```

## 2. Remove Duplicates:

```

python

# Find duplicates
df.duplicated().sum()

# Remove duplicates
df.drop_duplicates(inplace=True)

# Remove based on specific columns
df.drop_duplicates(subset=['name', 'date'], inplace=True)

```

## 3. Fix Data Types:

```

python

# Convert to numeric
df['price'] = pd.to_numeric(df['price'], errors='coerce')

# Convert to datetime (powerful!)
df['date'] = pd.to_datetime(df['date'])

# Extract date parts (Excel needs multiple formulas)
df['year'] = df['date'].dt.year
df['month'] = df['date'].dt.month
df['day_name'] = df['date'].dt.day_name()

```

## 4. Handle Outliers:

```
python

# IQR Method (Excel users don't do this!)
Q1 = df['price'].quantile(0.25)
Q3 = df['price'].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Remove outliers
df = df[(df['price'] >= lower_bound) & (df['price'] <= upper_bound)]
```

## 5. Feature Engineering (The Secret Weapon!):

```
python

# Create new columns from existing data
df['price_per_unit'] = df['total_price'] / df['quantity']
df['is_weekend'] = df['date'].dt.dayofweek >= 5
df['age_group'] = pd.cut(df['age'], bins=[0, 18, 35, 50, 100],
                         labels=['Teen', 'Young', 'Middle', 'Senior'])

# Encoding categorical variables (Excel struggles here)
df['category_encoded'] = df['category'].astype('category').cat.codes
# Or use one-hot encoding
df = pd.get_dummies(df, columns=['category'], prefix='cat')
```

**Time: 40% of project** (Most important phase!)

---

## Phase 4: STATISTICS & MODELS (Modeling)

### Visualization Libraries (V.S.P - Very Strong Plots):

```
python

import matplotlib.pyplot as plt # Basic plots
import seaborn as sns          # Beautiful plots (SNS)
import plotly.express as px    # Interactive plots
```

## Essential Visualizations:

### Distribution Plots:

```
python

# Histogram
plt.hist(df['price'], bins=30)
plt.title('Price Distribution')
plt.show()

# Box plot (shows outliers!)
sns.boxplot(data=df, x='category', y='price')

# Violin plot (advanced box plot)
sns.violinplot(data=df, x='category', y='price')
```

### Relationship Plots:

```
python

# Scatter plot
plt.scatter(df['age'], df['income'])

# Correlation heatmap (Excel can't do this easily!)
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')

# Pair plot (see all relationships at once!)
sns.pairplot(df)
```

### Time Series:

```
python

# Line plot
df.groupby('date')['sales'].sum().plot(figsize=(12,6))
plt.title('Sales Over Time')
```

### Statistical Analysis:

```
python
```

```

from scipy import stats

# T-test
stats.ttest_ind(group1, group2)

# Chi-square test
stats.chi2_contingency(pd.crosstab(df['cat1'], df['cat2']))

# Correlation
df['col1'].corr(df['col2'])

```

## Machine Learning Models:

```

python

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, mean_squared_error

# Split data
X = df[['feature1', 'feature2', 'feature3']]
y = df['target']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

# Train model
model = LinearRegression()
model.fit(X_train, y_train)

# Make predictions
predictions = model.predict(X_test)

# Evaluate
mse = mean_squared_error(y_test, predictions)

```

**Time: 20% of project**

---

## Phase 5: PERFORMANCE CHECK (Evaluation)

### Model Evaluation Metrics:

#### Regression (predicting numbers):

```
python

from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error

mse = mean_squared_error(y_test, predictions)
rmse = np.sqrt(mse)
mae = mean_absolute_error(y_test, predictions)
r2 = r2_score(y_test, predictions)

print(f'RMSE: {rmse}')
print(f'R²: {r2}')
```

#### Classification (predicting categories):

```
python

from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

accuracy = accuracy_score(y_test, predictions)
cm = confusion_matrix(y_test, predictions)

# Visualize confusion matrix
sns.heatmap(cm, annot=True, fmt='d')
print(classification_report(y_test, predictions))
```

#### Cross-Validation:

```
python

from sklearn.model_selection import cross_val_score

scores = cross_val_score(model, X, y, cv=5)
print(f'Average Score: {scores.mean():.2f} (+/- {scores.std():.2f})')
```

Time: 10% of project

---

## Phase 6: DEPLOY (Share Your Work)

### Create Interactive Dashboards:

```
python

import streamlit as st

# Simple dashboard
st.title('Sales Dashboard')
st.line_chart(df['sales'])
st.dataframe(df.head())
```

### Export Results:

```
python

# To Excel (with formatting!)
with pd.ExcelWriter('results.xlsx', engine='xlsxwriter') as writer:
    df.to_excel(writer, sheet_name='Data', index=False)
    summary.to_excel(writer, sheet_name='Summary', index=False)

# To CSV
df.to_csv('results.csv', index=False)

# To SQL
df.to_sql('table_name', conn, if_exists='replace', index=False)
```

### Create Reports:

```
python

# Automated PDF reports
from fpdf import FPDF

# Or HTML reports
import plotly.express as px
fig = px.scatter(df, x='x', y='y')
fig.write_html('report.html')
```

Time: 5% of project

---

## Phase 7: MONITOR (Keep Track)

### Schedule Automated Reports:

```
python

# Use schedule library
import schedule
import time

def run_analysis():
    # Your analysis code
    df = pd.read_csv('new_data.csv')
    # Process...
    # Send email/save results

schedule.every().day.at("09:00").do(run_analysis)

while True:
    schedule.run_pending()
    time.sleep(60)
```

## THE ULTIMATE PYTHON ADVANTAGE SUMMARY

### What Makes You Stand Out:

| Task             | Excel/SQL Limit   | Python Power            |
|------------------|-------------------|-------------------------|
| Data Size        | ~1M rows          | Unlimited (billions)    |
| Automation       | Manual clicks     | Fully automated scripts |
| Visualization    | Basic charts      | Interactive dashboards  |
| Statistics       | Limited functions | Full statistical suite  |
| Machine Learning | ✗ None            | ✓ All algorithms        |
| Multiple Files   | Manual merge      | One command             |
| Reproducibility  | Hard to replicate | Rerun entire analysis   |
| Version Control  | Save As...        | Git integration         |

| Task  | Excel/SQL Limit      | Python Power   |
|-------|----------------------|----------------|
| Speed | Slow with large data | Lightning fast |

## 📦 COMPLETE LIBRARY INSTALLATION

```
bash

# The Essential Data Science Stack
pip install pandas numpy matplotlib seaborn
pip install scikit-learn scipy
pip install jupyter notebook
pip install plotly streamlit
pip install openpyxl xlsxwriter # For Excel files
pip install missingno # Missing data visualization

# Advanced (install as needed)
pip install statsmodels # Advanced statistics
pip install xgboost lightgbm # Advanced ML
pip install dash # Dashboards
```

## 🎓 QUICK REFERENCE CHEAT SHEET

### Abbreviations to Remember:

- **PD** = pandas (Data manipulation)
- **NP** = numpy (Numbers & arrays)
- **PLT** = pyplot (Plotting)
- **SNS** = seaborn (Statistical plots)
- **SK** = sklearn (Machine Learning)
- **EDA** = Exploratory Data Analysis
- **ML** = Machine Learning
- **DF** = DataFrame (your data table)

## The 5-Minute Workflow:

```
python

# 1. Import
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# 2. Load
df = pd.read_csv('data.csv')

# 3. Explore
df.head()
df.info()
df.describe()

# 4. Clean
df.dropna(inplace=True)
df.drop_duplicates(inplace=True)

# 5. Analyze
df.groupby('category')['sales'].sum().plot(kind='bar')
plt.show()

# 6. Export
df.to_excel('results.xlsx', index=False)
```

---

## 🏆 PROJECT TEMPLATE (Copy & Use!)

```
python
```

!!!!

PROJECT: [Your Project Name]  
DATE: [Date]  
GOAL: [What you're trying to achieve]  
!!!!

### # 1. IMPORTS

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
from sklearn.model_selection import train_test_split
```

### # 2. LOAD DATA

```
df = pd.read_csv('data.csv')  
print(f'Loaded {df.shape[0]} rows and {df.shape[1]} columns')
```

### # 3. EXPLORE

```
print(df.head())  
print(df.info())  
print(df.describe())
```

### # 4. CLEAN

```
# Missing values  
print(f'Missing values:\n{df.isnull().sum()}')  
df.fillna(df.mean(), inplace=True)
```

### # Duplicates

```
print(f'Duplicates: {df.duplicated().sum()}')  
df.drop_duplicates(inplace=True)
```

### # 5. FEATURE ENGINEERING

```
# Add your custom columns here
```

### # 6. VISUALIZE

```
sns.pairplot(df)  
plt.show()
```

### # 7. MODEL (if needed)

```
# Add your modeling code
```

### # 8. EXPORT

```
df.to_excel('results.xlsx', index=False)
print("Analysis complete!")
```

## PRO TIPS

1. **Always save your code** in Jupyter Notebooks (.ipynb)
2. **Comment everything** - Future you will thank you
3. **Use meaningful variable names** - `df_sales` not `df1`
4. **Test on small data first** - Don't wait 10 minutes to find a typo
5. **Version your data** - Save cleaned data separately
6. **Use virtual environments** - Keep projects isolated
7. **Learn keyboard shortcuts** - Shift+Enter to run cells
8. **Google is your friend** - Everyone looks up syntax!

## Next Steps to Master

1. **Week 1-2:** Master pandas basics
2. **Week 3-4:** Learn data cleaning techniques
3. **Week 5-6:** Practice visualization
4. **Week 7-8:** Introduction to ML
5. **Week 9+:** Real projects!

**Remember:** The best way to learn is by doing. Start with a simple dataset and work through each phase!

## Resources

- **Practice Datasets:** [kaggle.com/datasets](https://kaggle.com/datasets)
- **Cheat Sheets:** [datacamp.com/cheat-sheets](https://datacamp.com/cheat-sheets)
- **Documentation:** [pandas.pydata.org](https://pandas.pydata.org)
- **Community:** [stackoverflow.com](https://stackoverflow.com), [reddit.com/r/datascience](https://www.reddit.com/r/datascience)

