

Exception Handling, Some Python Packages - Numpy, Pandas, Matplotlib, File Handling, Time-Space Complexity



Martin Wafula
Multimedia University of Kenya

April 23, 2025

Exception Handling

NumPy

Pandas

Matplotlib

File Handling

GUI Development

Machine Learning

Exception Handling



Definition

Exceptions are Python's mechanism for handling runtime errors and exceptional conditions in a controlled way.

Key Characteristics:

- ▶ Disrupt normal program flow
- ▶ Are objects (instances of classes)
- ▶ Follow inheritance hierarchy
- ▶ Can be caught and handled

Common Exception Types:

- ▶ `ValueError`
- ▶ `TypeError`
- ▶ `IndexError`
- ▶ `KeyError`
- ▶ `FileNotFoundException`

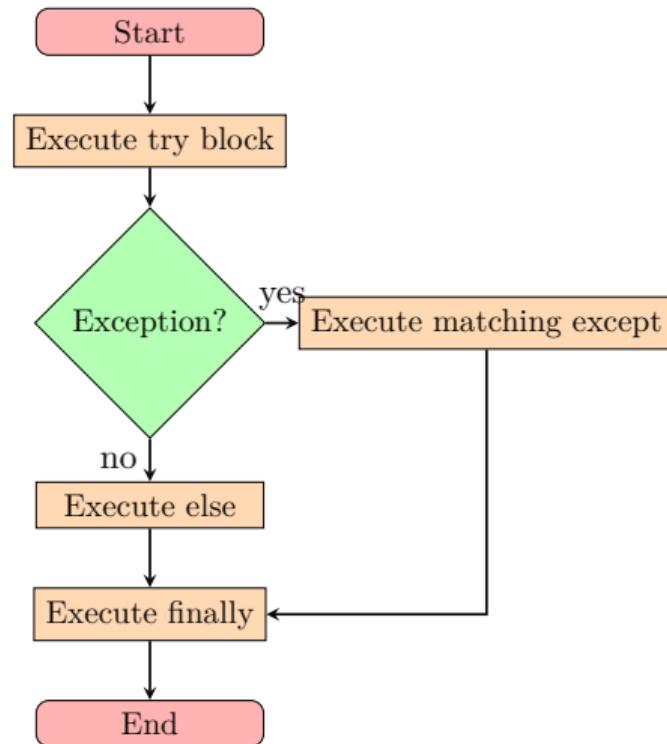
Python Philosophy

EAFP (Easier to Ask for Forgiveness than Permission) style is preferred over LBYL (Look Before You Leap).

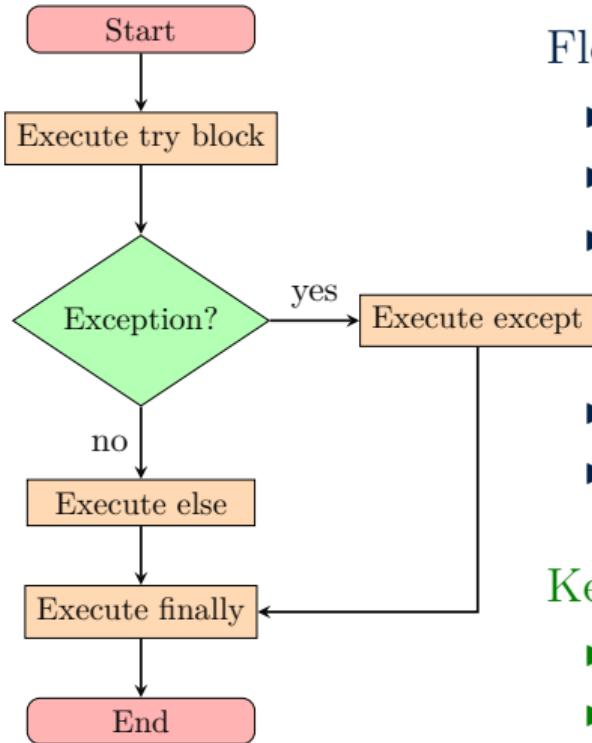
Basic Exception Handling Syntax

```
1 # Basic try-except block
2 try:
3     # Code that might raise exceptions
4     result = 10 / denominator
5     file = open("data.txt")
6     value = int(file.read())
7 except ZeroDivisionError:
8     print("Cannot divide by zero!")
9 except FileNotFoundError:
10    print("File not found!")
11 except ValueError as ve:
12     print(f"Invalid value: {ve}")
13 except Exception as e:
14     print(f"Unexpected error: {type(e).__name__}")
15 else:
16     print("Operation successful!")
17 finally:
18     file.close() # Always executes
19     print("Cleanup complete")
```

Exception Handling Flow



Exception Handling Flow



Flow Explanation

- ▶ **Start:** Begin exception handling block
- ▶ **Try:** Code that might raise exceptions
- ▶ **Exception?:** Check if any exception occurred
 - ▶ **Yes:** Execute matching **except** block
 - ▶ **No:** Execute **else** block (optional)
- ▶ **Finally:** Always executes (cleanup code)
- ▶ **End:** Continue with normal program flow

Key Points

- ▶ **else** runs only if no exceptions occur
- ▶ **finally** runs *always* (even if exceptions propagate)
- ▶ **Order matters:** Specific exceptions first

```
1 class InvalidEmailError(ValueError):
2     """Raised when email format is invalid"""
3     def __init__(self, email):
4         self.email = email
5         super().__init__(f"Invalid email: {email}")
6 class WeakPasswordError(Exception):
7     """Raised when password doesn't meet requirements"""
8     def __init__(self, reason):
9         self.reason = reason
10        super().__init__(f"Weak password: {reason}")
11    def validate_user(email, password):
12        if "@" not in email:
13            raise InvalidEmailError(email)
14        if len(password) < 8:
15            raise WeakPasswordError("too short")
16        if not any(c.isdigit() for c in password):
17            raise WeakPasswordError("no digits")
18        return True
```

Exception Chaining

```
1 def process_file(filename):
2     try:
3         with open(filename) as f:
4             return int(f.read())
5     except FileNotFoundError as fnf_error:
6         raise RuntimeError("Processing failed") from fnf_error
7
8 def main():
9     try:
10         result = process_file("missing.txt")
11     except RuntimeError as err:
12         print(f"Error: {err}")
13         print(f"Caused by: {err.__cause__}")
14         # Access original traceback
15         import traceback
16         traceback.print_tb(err.__traceback__)
17
```

Key Attributes

- ▶ `__cause__`: Direct cause of exception
- ▶ `__context__`: Context for implicit chaining
- ▶ `__traceback__`: Traceback object

```
1 class DatabaseConnection:
2     def __init__(self, dbname):
3         self.dbname = dbname
4         self.connection = None
5
6     def __enter__(self):
7         print(f"Connecting to {self.dbname}")
8         self.connection = sqlite3.connect(self.dbname)
9         return self.connection
10
11    def __exit__(self, exc_type, exc_val, exc_tb):
12        print("Closing connection")
13        if self.connection:
14            self.connection.close()
15        if exc_type is not None:
16            print(f"Error occurred: {exc_val}")
17        return True # Suppress exception
18
19 # Usage
```

```
1 import logging
2
3 logging.basicConfig(
4     filename='app.log',
5     level=logging.ERROR,
6     format='%(asctime)s - %(levelname)s - %(message)s'
7 )
8
9 def process_data(data):
10     try:
11         result = complex_operation(data)
12     except ValueError as e:
13         logging.error(f"Value error in processing: {e}", exc_info=True)
14         raise
15     except Exception as e:
16         logging.critical(f"Unexpected error: {e}", exc_info=True)
17         raise
18     else:
19         logging.info("Processing completed successfully")
```

Do's

- ▶ Catch specific exceptions first
- ▶ Use custom exceptions for domain errors
- ▶ Clean up resources in `finally` blocks
- ▶ Log exceptions with context
- ▶ Preserve stack traces when re-raising

Don'ts

- ▶ Don't catch `Exception` unless necessary
- ▶ Don't ignore exceptions silently
- ▶ Don't use exceptions for flow control
- ▶ Don't expose sensitive information

Real-World Example: Web API Client

```
1 import requests
2 from requests.exceptions import RequestException
3 class APIClient:
4     def __init__(self, base_url):
5         self.base_url = base_url
6     def get_data(self, endpoint, retries=3):
7         url = f"{self.base_url}/{endpoint}"
8         for attempt in range(retries):
9             try:
10                 response = requests.get(url, timeout=5)
11                 response.raise_for_status()
12                 return response.json()
13             except requests.HTTPError as http_err:
14                 if response.status_code == 404:
15                     raise ValueError("Endpoint not found") from http_err
16                 if attempt == retries - 1:
17                     raise
18             except RequestException as req_err:
19                 if attempt == retries - 1:
```

Database Transaction Example

```
1 import sqlite3
2 from contextlib import contextmanager
3
4 @contextmanager
5 def transaction(db_path):
6     conn = sqlite3.connect(db_path)
7     cursor = conn.cursor()
8     try:
9         yield cursor
10        conn.commit()
11    except Exception as e:
12        conn.rollback()
13        raise e
14    finally:
15        conn.close()
16 # Usage
17 with transaction("app.db") as cursor:
18     cursor.execute("INSERT INTO users VALUES (?, ?)", ("alice", 25))
19     cursor.execute("UPDATE accounts SET balance = balance - 100")
```

Exception Handling in Threads

```
1 import threading
2 import queue
3
4 def worker(task_queue, result_queue):
5     while True:
6         try:
7             task = task_queue.get_nowait()
8             result = process_task(task)
9             result_queue.put(result)
10        except queue.Empty:
11            break
12        except Exception as e:
13            result_queue.put(e)
14            break
15
16 def main():
17     tasks = [1, 2, 3, "invalid", 5]
18     task_queue = queue.Queue()
19     result_queue = queue.Queue()
```

Handling Keyboard Interrupt

```
1 import time
2
3 def long_running_process():
4     try:
5         for i in range(10):
6             print(f"Processing step {i}")
7             time.sleep(1)
8     except KeyboardInterrupt:
9         print("\nReceived interrupt, cleaning up...")
10        # Perform cleanup
11        time.sleep(1)  # Simulate cleanup
12        print("Cleanup complete")
13        raise # Re-raise to exit
14
15 def main():
16     try:
17         long_running_process()
18     except KeyboardInterrupt:
19         print("Program terminated by user")
```

```
1 def validate_user(user):
2     errors = []
3     if len(user.name) < 3:
4         errors.append(ValueError("Name too short"))
5     if not user.email or "@" not in user.email:
6         errors.append(ValueError("Invalid email"))
7     if len(user.password) < 8:
8         errors.append(ValueError("Password too short"))
9     if errors:
10        raise ExceptionGroup("User validation failed", errors)
11
12 try:
13     validate_user(user)
14 except* ValueError as eg:
15     for error in eg.exceptions:
16         print(f"Validation error: {error}")
17
```

Numpy



Core Features

- ▶ Homogeneous N-dimensional arrays
- ▶ Vectorized operations
- ▶ Memory-efficient contiguous storage
- ▶ Broadcasting capabilities

Memory Layout:

- ▶ Contiguous memory block
- ▶ Fixed-type elements
- ▶ Strides for indexing
- ▶ Views vs copies

Performance Benefits:

- ▶ Avoid Python interpreter overhead
- ▶ Optimized C/Fortran code
- ▶ Cache-friendly access
- ▶ Parallelizable operations

```
1 import numpy as np
2 # From Python sequences
3 arr1 = np.array([1, 2, 3])                      # 1D array
4 arr2 = np.array([[1, 2], [3, 4]])                # 2D array
5 # Special arrays
6 zeros = np.zeros((3, 4))                         # 3x4 zeros
7 ones = np.ones((2, 2, 2))                         # 2x2x2 ones
8 identity = np.eye(3)                             # Identity matrix
9 empty = np.empty((2, 3))                          # Uninitialized
10 # Ranges and sequences
11 range_arr = np.arange(0, 10, 0.5)               # With step
12 linspace = np.linspace(0, 1, 5)                 # 5 points in [0,1]
13 # Random arrays
14 random_arr = np.random.rand(2, 3)               # Uniform [0,1)
15 normal_arr = np.random.normal(0, 1, (2, 2))    # Normal dist
16
```

```
1 arr = np.array([[1, 2, 3], [4, 5, 6]])
2
3 # Basic attributes
4 print(arr.shape)      # (2, 3)
5 print(arr.ndim)        # 2
6 print(arr.size)        # 6
7 print(arr.dtype)       # int64
8 print(arr.itemsize)    # 8 (bytes per element)
9
10 # Array methods
11 print(arr.sum())       # 21
12 print(arr.mean(axis=0)) # [2.5 3.5 4.5] (column means)
13 print(arr.max(axis=1)) # [3 6] (row maximums)
14 print(arr.reshape(3, 2))# Reshaped array
15 print(arr.flatten())   # [1 2 3 4 5 6]
16 print(arr.T)           # Transpose
17
```

```
1 arr = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
2 # Basic indexing
3 print(arr[0, 1])      # 2 (row 0, column 1)
4 print(arr[:, 1])       # [2 5 8] (all rows, column 1)
5 print(arr[1:3, :2])    # [[4 5], [7 8]] (rows 1-2, cols 0-1)
6 # Boolean indexing
7 mask = arr > 5
8 print(mask)            # [[False False False], [False False True], [True
                           True True]]
9 print(arr[mask])       # [6 7 8 9]
10 # Fancy indexing
11 print(arr[[0, 2], [1, 0]]) # [2 7] (elements (0,1) and (2,0))
12 # Views vs copies
13 arr_view = arr[:2, :2]    # View (shares memory)
14 arr_copy = arr[:2, :2].copy() # Copy (new memory)
```

15

```
1 a = np.array([1, 2, 3])
2 b = np.array([4, 5, 6])
3 # Element-wise operations
4 print(a + b)          # [5 7 9]
5 print(a * 2)          # [2 4 6]
6 print(np.sqrt(a))     # [1. 1.414 1.732]
7 print(np.exp(a))      # [2.718 7.389 20.085]
8 # Comparison operations
9 print(a == b)         # [False False False]
10 print(a > 1)        # [False True True]
11 # Matrix operations
12 mat = np.array([[1, 2], [3, 4]])
13 print(mat @ mat)     # Matrix multiplication
14 print(np.linalg.inv(mat)) # Matrix inverse
15
```

Broadcasting Definition

NumPy's ability to work with arrays of different shapes during arithmetic operations

```
1 # Broadcasting examples
2 a = np.array([1, 2, 3])
3 b = np.array([[10], [20]])
4
5 print(a + b) # [[11 12 13], [21 22 23]]
6
7 # Rules:
8 # 1. Align array shapes from right
9 # 2. Dimensions must match or be 1
10 # 3. Missing dimensions are treated as 1
11
12 # Example shapes:
13 # (3,) + (2,1) -> (1,3) + (2,1) -> (2,3)
```

14

What are ufuncs?

Fast element-wise operations that support broadcasting, type casting, and other features

```
1 # Built-in ufuncs
2 arr = np.array([0, np.pi/2, np.pi])
3 print(np.sin(arr))      # [0.  1.  0.]
4 print(np.add(arr, 1))   # [1.  2.571  4.142]
5 # Creating custom ufuncs
6 def my_func(x, y):
7     return x**2 + y**2
8 my_ufunc = np.frompyfunc(my_func, 2, 1)
9 print(my_ufunc(arr, 2)) # [4.  6.467 13.87]
```

10

Key Features

- ▶ Optimized C implementations

- ▶ Support all array methods

```
1 data = np.array([(1, 'Alice', 25.5),
2                  (2, 'Bob', 30.2)],
3                  dtype=[('id', 'i4'), ('name', 'U10'), ('age', 'f4')])
4 # Accessing fields
5 print(data['name'])    # ['Alice' 'Bob']
6 print(data[0]['age'])  # 25.5
7 # Record arrays (alternative syntax)
8 recarr = data.view(np.recarray)
9 print(recarr.age)     # [25.5 30.2]
10 # Sorting
11 data_sorted = np.sort(data, order='age')
12
```

Use Cases

- ▶ Tabular data with mixed types
- ▶ Memory-efficient storage
- ▶ Interoperability with other systems

```
1 # Saving and loading arrays
2 arr = np.random.rand(3, 3)
3 # Text files
4 np.savetxt('array.txt', arr, delimiter=',')
5 loaded = np.loadtxt('array.txt', delimiter=',')
6 # Binary files (more efficient)
7 np.save('array.npy', arr)          # Single array
8 np.savez('arrays.npz', a=arr, b=arr*2)  # Multiple arrays
9 loaded = np.load('array.npy')
10 arrays = np.load('arrays.npz')
11 print(arrays['a'])   # Access by name
12 # Memory mapping
13 large_arr = np.memmap('large.dat', dtype='float32',
14                         mode='w+', shape=(1000, 1000))
15
```

```
1 import numpy as np
2 import numexpr as ne
3 # Vectorized vs loop
4 def slow_dot(a, b):
5     result = 0
6     for x, y in zip(a, b):
7         result += x * y
8     return result
9 a = np.random.rand(1000000)
10 b = np.random.rand(1000000)
11 %timeit slow_dot(a, b)  # ~500ms
12 %timeit np.dot(a, b)    # ~1ms
13 # Using numexpr for complex expressions
14 expr = ne.evaluate('sin(a) + log(b) * sqrt(a/b)')
15
```

Optimization Tips

- ▶ Avoid Python loops
- ▶ Use built-in functions
- ▶ Consider memory layout
- ▶ Use specialized libraries

Do's

- ▶ Preallocate arrays when possible
- ▶ Use vectorized operations
- ▶ Specify dtypes explicitly
- ▶ Use views instead of copies
- ▶ Leverage broadcasting

Don'ts

- ▶ Don't treat NumPy arrays like Python lists
- ▶ Don't ignore array shapes
- ▶ Don't mix Python types in arrays
- ▶ Don't forget about memory usage

Real-World Example: Image Processing

```
1 from scipy.misc import face
2 import matplotlib.pyplot as plt
3 # Load sample image (raccoon)
4 img = face()
5 # Convert to grayscale
6 gray = img.mean(axis=2)
7 # Apply filters
8 blurred = np.zeros_like(gray, dtype=float)
9 blurred[1:-1,1:-1] = (gray[:-2,:-2] + gray[:-2,1:-1] + gray[:-2,2:] +
10                      gray[1:-1,:-2] + gray[1:-1,1:-1] + gray[1:-1,2:] +
11                      gray[2:,:-2] + gray[2:,1:-1] + gray[2:,2:]) / 9
12 # Thresholding
13 binary = (gray > gray.mean()).astype(int)
14 # Display
15 plt.imshow(binary, cmap='gray')
16 plt.show()
```

17

Pandas



Series:

- ▶ 1D labeled array
- ▶ Homogeneous data
- ▶ Index alignment
- ▶ Vectorized operations

Series Creation

```
1 s = pd.Series([1, 2, 3],  
2                 index=['a', 'b', 'c'],  
3                 dtype='float32',  
4                 name='values')  
5
```

DataFrame:

- ▶ 2D labeled structure
- ▶ Columns can be different types
- ▶ Row and column indexes
- ▶ SQL-like operations

DataFrame Creation

```
1 data = {'A': [1, 2], 'B': ['x', 'y']}
2 df = pd.DataFrame(data,
3                     index=['row1', 'row2'])
4
```

```
1 # Create sample DataFrame
2 df = pd.DataFrame({
3     'Name': ['Alice', 'Bob', 'Charlie'],
4     'Age': [25, 30, 35],
5     'City': ['NY', 'LA', 'Chicago']
6 }, index=['a', 'b', 'c'])
7 # Selection methods
8 df['Name']           # Column selection
9 df.loc['a']           # Row by label
10 df.iloc[0]           # Row by position
11 df[df['Age'] > 28] # Boolean indexing
12 df.query('Age > 28 & City == "LA"') # Query method
13 # Multi-indexing
14 df.set_index(['City', 'Name']).sort_index()
15 df.xs('NY', level='City') # Cross-section
16 # Adding/removing columns
17 df['Senior'] = df['Age'] > 30
18 df.drop(columns=['Senior'], inplace=True)
```

```
1 # Create DataFrame with missing values
2 df = pd.DataFrame({
3     'A': [1, 2, np.nan],
4     'B': [5, np.nan, np.nan],
5     'C': ['a', 'b', None]
6 })
7 # Detection
8 df.isna()
9 df.notna()
10 # Handling
11 df.dropna()          # Drop rows with any NA
12 df.fillna(0)         # Fill NA with 0
13 df.fillna(df.mean()) # Fill with mean
14 df.interpolate()    # Interpolate values
15 # Advanced
16 df['A'].fillna(df['A'].mean(), inplace=True)
17 df['B'].fillna(method='ffill', inplace=True)
18
```

```
1 # Example messy data
2 data = {
3     'date': ['2022-01-01', '2022/02/01', '01-03-2022'],
4     'price': ['$10.50', '15.99', '20'],
5     'quantity': [1, 'two', 3]
6 }
7 df = pd.DataFrame(data)
8 # Cleaning steps
9 df['date'] = pd.to_datetime(df['date'], errors='coerce')
10 df['price'] = df['price'].str.replace('$', '').astype(float)
11 df['quantity'] = pd.to_numeric(df['quantity'], errors='coerce')
12 # Outlier handling
13 df = df[(df['price'] > 0) & (df['price'] < 100)]
14 df = df.dropna()
15 # String operations
16 df['price_category'] = np.where(df['price'] < 15, 'low', 'high')
17
```

```
1 # Sample sales data
2 sales = pd.DataFrame({
3     'region': ['East', 'West', 'East', 'West', 'East'],
4     'product': ['A', 'B', 'A', 'B', 'A'],
5     'revenue': [100, 150, 200, 50, 300],
6     'cost': [80, 120, 150, 40, 250]
7 })
8
9 # Basic grouping
10 grouped = sales.groupby('region')
11 grouped.sum()
12 grouped.mean()
13
14 # Multiple aggregations
15 sales.groupby('region').agg({
16     'revenue': ['sum', 'mean'],
17     'cost': ['sum', 'mean']
18 })
```

```
1 # Sample DataFrames
2 customers = pd.DataFrame({
3     'id': [1, 2, 3],
4     'name': ['Alice', 'Bob', 'Charlie']
5 })
6
7 orders = pd.DataFrame({
8     'order_id': [101, 102, 103],
9     'customer_id': [1, 3, 4],
10    'amount': [50, 100, 75]
11 })
12
13 # Merge options
14 pd.merge(customers, orders, left_on='id', right_on='customer_id')
15 pd.merge(customers, orders, left_on='id', right_on='customer_id', how='left')
16 pd.merge(customers, orders, left_on='id', right_on='customer_id', how='outer')
```

```
1 # Create time series
2 date_rng = pd.date_range(start='1/1/2022', end='1/10/2022', freq='D')
3 ts = pd.Series(np.random.randn(len(date_rng)), index=date_rng)
4
5 # Time-based operations
6 ts.resample('3D').mean() # Resample to 3-day periods
7 ts.rolling(window=3).mean() # Moving average
8
9 # Timezone handling
10 ts = ts.tz_localize('UTC')
11 ts = ts.tz_convert('US/Eastern')
12
13 # Shifting
14 ts.shift(1) # Shift forward
15 ts.shift(-1) # Shift backward
16
17 # Time deltas
18 ts.index + pd.Timedelta(days=2)
19
```

```
1 # Reading from files
2 df_csv = pd.read_csv('data.csv')
3 df_excel = pd.read_excel('data.xlsx', sheet_name='Sheet1')
4 df_json = pd.read_json('data.json')
5 df_sql = pd.read_sql('SELECT * FROM table', con)
6
7 # Writing to files
8 df.to_csv('output.csv', index=False)
9 df.to_excel('output.xlsx', sheet_name='Results')
10 df.to_json('output.json', orient='records')
11 df.to_sql('table', con, if_exists='replace')
12
13 # Performance tips
14 pd.read_csv('large.csv', chunksize=10000) # Process in chunks
15 pd.read_csv('data.csv', usecols=['col1', 'col2']) # Read only needed
   columns
16
```

```
1 # Pivot tables
2 sales.pivot_table(index='region', columns='product',
3                     values='revenue', aggfunc='sum')
4
5 # Melt
6 pd.melt(df, id_vars=['Name'], value_vars=['Math', 'Science'])
7
8 # Apply functions
9 df['Age'].apply(lambda x: x + 1)
10 df.apply(np.mean, axis=0) # Column-wise
11 df.apply(np.mean, axis=1) # Row-wise
12
13 # Vectorized operations
14 df['revenue'] = df['price'] * df['quantity']
15 df['profit'] = df['revenue'] - df['cost']
16
```

```
1 # Creating categories
2 df['grade'] = pd.Categorical(
3     ['A', 'B', 'C', 'A', 'B'],
4     categories=['A', 'B', 'C', 'D'],
5     ordered=True
6 )
7
8 # Operations
9 df['grade'] > 'B' # Comparison works with ordered
10 df['grade'].cat.codes # Numeric codes
11 df['grade'].cat.categories # View categories
12
13 # Memory efficiency
14 df['large_category'] = df['large_category'].astype('category')
15
16 # Grouping benefits
17 df.groupby('grade').size()
18
```

```
1 # Data types matter
2 df['id'] = df['id'].astype('int32')    # Save memory
3 # Avoid iterrows()
4 # Instead of:
5 for index, row in df.iterrows():
6     process(row)
7
8 # Use:
9 df.apply(process, axis=1)
10 # Or vectorized operations:
11 df['new_col'] = df['col1'] + df['col2']
12 # Query optimization
13 df.query('age > 30 & salary < 50000')  # Often faster than boolean
14     indexing
15 # Use eval() for complex expressions
16 df.eval('revenue = price * quantity - cost', inplace=True)
17
```

```
1 # Load and clean data
2 sales = pd.read_csv('sales.csv', parse_dates=['date'])
3 sales['revenue'] = sales['price'] * sales['quantity']
4 # Monthly analysis
5 monthly = sales.resample('M', on='date').agg({
6     'revenue': 'sum',
7     'quantity': 'sum',
8     'product': 'nunique'
9 })
10 # Top products
11 top_products = sales.groupby('product')['revenue'].sum().nlargest(5)
12 # Customer segmentation
13 sales['customer_type'] = pd.cut(
14     sales['revenue'],
15     bins=[0, 100, 500, float('inf')],
16     labels=['small', 'medium', 'large']
17 )
18 # Visualization
19 monthly['revenue'].plot(kind='bar')
```

Do's

- ▶ Use vectorized operations
- ▶ Specify data types
- ▶ Use categoricals for strings
- ▶ Leverage built-in methods
- ▶ Use query() for complex filtering

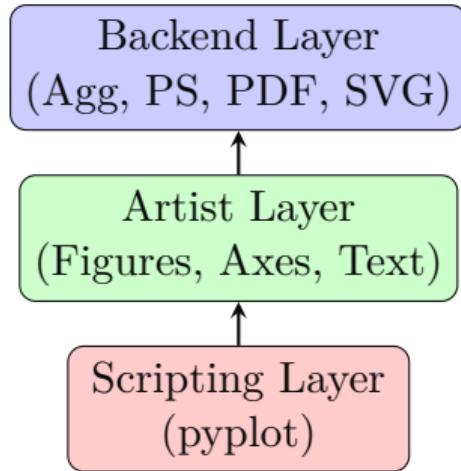
Don'ts

- ▶ Don't use iterrows()/itertuples() unless necessary
- ▶ Don't ignore chained assignment warnings
- ▶ Don't modify data during iteration
- ▶ Don't forget to set indexes when appropriate

```
1 # Chunk processing
2 chunk_size = 10000
3 results = []
4 for chunk in pd.read_csv('huge.csv', chunks=chunk_size):
5     results.append(process_chunk(chunk))
6 final = pd.concat(results)
7
8 # Dask alternative
9 import dask.dataframe as dd
10 ddf = dd.read_csv('huge_*.csv')
11 result = ddf.groupby('category').size().compute()
12
13 # Memory optimization
14 df = pd.read_csv('data.csv', dtype={'id': 'int32', 'price': 'float32'})
15 df = df.select_dtypes(include=['number']) # Keep only numeric columns
16
17 # Sparse data
18 df = df.to_sparse() # For data with many NaN/zeros
19
```

Matplotlib





Three Layer Architecture:

1. Scripting: Quick plotting (MATLAB-style)
2. Artist: Object-oriented interface
3. Backend: Actual rendering

Basic Plotting Example

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 # Create data
5 x = np.linspace(0, 2*np.pi, 100)
6 y1 = np.sin(x)
7 y2 = np.cos(x)
8
9 # Create figure and axes
10 fig, ax = plt.subplots(figsize=(10, 5))
11
12 # Plot data
13 ax.plot(x, y1, label='sin(x)', color='blue', linestyle='--')
14 ax.plot(x, y2, label='cos(x)', color='red', linestyle='--')
15
16 # Customize
17 ax.set_title("Trigonometric Functions")
18 ax.set_xlabel("x (radians)")
19 ax.set_ylabel("Function Value")
```

```
1 # Create figure with multiple subplots
2 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4))
3
4 # First subplot
5 ax1.plot(x, np.sin(x))
6 ax1.set_title('Sine Wave')
7 ax1.set_facecolor('#f0f0f0') # Light gray background
8 ax1.spines['top'].set_visible(False) # Remove top spine
9
10 # Second subplot
11 ax2.plot(x, np.cos(x), 'r--')
12 ax2.set_title('Cosine Wave')
13 ax2.set_ylim(-1.5, 1.5)
14 ax2.grid(True, linestyle=':')
15
16 # Figure-level settings
17 fig.suptitle('Trigonometric Functions', fontsize=16)
18 fig.tight_layout() # Adjust spacing
19
```

```
1 # Line plot
2 plt.plot(x, y)
3 # Scatter plot
4 plt.scatter(x, y, c=colors, s=sizes, alpha=0.5)
5 # Bar chart
6 plt.bar(categories, values)
7 # Histogram
8 plt.hist(data, bins=20)
9 # Box plot
10 plt.boxplot([data1, data2])
11 # Pie chart
12 plt.pie(sizes, labels=labels, autopct='%1.1f%%')
13 # Heatmap
14 plt.imshow(matrix, cmap='viridis')
15 plt.colorbar()
16 # Contour plot
17 plt.contour(X, Y, Z, levels=20)
18
```

```
1 # Create figure
2 fig = plt.figure(figsize=(10, 6))
3 ax = fig.add_subplot(111, projection='polar') # Polar plot
4 # Plot data
5 theta = np.linspace(0, 2*np.pi, 100)
6 r = np.abs(np.sin(5*theta))
7 ax.plot(theta, r, linewidth=3)
8 # Customize
9 ax.set_theta_offset(np.pi/2)
10 ax.set_theta_direction(-1)
11 ax.set_rgrids([0.2, 0.4, 0.6, 0.8, 1.0])
12 # Annotations
13 ax.annotate('Peak', xy=(np.pi/2, 1), xytext=(0.5, 1.1),
14             arrowprops=dict(facecolor='black', shrink=0.05))
15
16 # Text and arrows
17 ax.text(0, 0, 'Center', ha='center')
18
```

```
1 # Available styles
2 print(plt.style.available)
3 # ['ggplot', 'seaborn', 'dark_background', 'bmh', ...]
4 # Using styles
5 plt.style.use('ggplot')
6 # Custom style dictionary
7 mystyle = {
8     'axes.facecolor': 'white',
9     'axes.grid': True,
10    'grid.color': '.8',
11    'xtick.color': '.15',
12    'ytick.color': '.15'
13 }
14 plt.style.use(mystyle)
15 # Context manager for temporary style
16 with plt.style.context('dark_background'):
17     plt.plot(x, np.sin(x))
18     plt.show()
```

```
1 # Basic saving
2 plt.savefig('plot.png') # Default format from extension
3
4 # Quality and size control
5 plt.savefig('high_res.png', dpi=300, bbox_inches='tight')
6
7 # Different formats
8 plt.savefig('plot.pdf') # Vector format
9 plt.savefig('plot.svg') # Scalable Vector Graphics
10
11 # Transparent background
12 plt.savefig('transparent.png', transparent=True)
13 # Multiple plots to PDF
14 from matplotlib.backends.backend_pdf import PdfPages
15 with PdfPages('multipage.pdf') as pdf:
16     plt.plot(x, y1)
17     pdf.savefig()
18     plt.plot(x, y2)
19     pdf.savefig()
```

```
1 # Interactive mode
2 plt.ion() # Turn on interactive mode
3
4 # Create plot
5 fig, ax = plt.subplots()
6 line, = ax.plot(x, np.sin(x))
7
8 # Update plot in loop
9 for phase in np.linspace(0, 2*np.pi, 20):
10     line.set_ydata(np.sin(x + phase))
11     fig.canvas.draw()
12     fig.canvas.flush_events()
13     plt.pause(0.1)
14
15 # Widgets
16 from matplotlib.widgets import Slider
17 ax_slider = plt.axes([0.2, 0.1, 0.6, 0.03])
18 slider = Slider(ax_slider, 'Phase', 0, 2*np.pi)
```

```
1 from mpl_toolkits.mplot3d import Axes3D  
2  
3 # Create 3D axes  
4 fig = plt.figure(figsize=(8, 6))  
5 ax = fig.add_subplot(111, projection='3d')  
6  
7 # Generate data  
8 theta = np.linspace(0, 2*np.pi, 100)  
9 z = np.linspace(0, 10, 100)  
10 r = z  
11 x = r * np.sin(theta)  
12 y = r * np.cos(theta)  
13  
14
```

```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 # Plot
4 ax.plot(x, y, z, label='3D spiral')
5 ax.set_xlabel('X')
6 ax.set_ylabel('Y')
7 ax.set_zlabel('Z')
8 ax.legend()
9
10 # Surface plot
11 X, Y = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
12 Z = np.sin(np.sqrt(X**2 + Y**2))
13 ax.plot_surface(X, Y, Z, cmap='viridis')
14
```

Real-World Example: Stock Analysis

```
1 import pandas as pd
2 import matplotlib.dates as mdates
3
4 # Load stock data
5 df = pd.read_csv('stock.csv', parse_dates=['Date'])
6 df = df.set_index('Date')
7
8 # Create figure
9 fig, (ax1, ax2) = plt.subplots(2, 1, figsize=(12, 8))
10
11 # Price plot
12 ax1.plot(df.index, df['Close'], label='Closing Price')
13 ax1.set_ylabel('Price ($)')
14 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m'))
15 ax1.grid(True)
16
17
```

```
1 # Volume plot
2 ax2.bar(df.index, df['Volume'], width=5, label='Volume')
3 ax2.set_ylabel('Volume')
4 ax2.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m'))
5
6
7 # Formatting
8 fig.autofmt_xdate()
9 fig.suptitle('Stock Performance')
10 plt.tight_layout()
11
```

Do's

- ▶ Use the OO interface for complex plots
- ▶ Set figure size early
- ▶ Use `tight_layout()` to prevent overlaps
- ▶ Label your axes and add titles
- ▶ Choose appropriate color maps

Don'ts

- ▶ Don't use the pyplot interface for complex figures
- ▶ Don't forget to close figures in scripts
- ▶ Don't use jet colormap
- ▶ Don't overload your plots

```
1 import seaborn as sns
2
3 # Set style
4 sns.set_style("whitegrid")
5 # Create complex plots easily
6 tips = sns.load_dataset("tips")
7 sns.relplot(x="total_bill", y="tip", hue="day",
8             col="time", data=tips)
9 # Statistical plots
10 sns.boxplot(x="day", y="total_bill", data=tips)
11 sns.violinplot(x="day", y="total_bill", data=tips)
12 # Regression plots
13 sns.lmplot(x="total_bill", y="tip", data=tips)
14 # Matrix plots
15 flights = sns.load_dataset("flights").pivot("month", "year", "passengers")
16 sns.heatmap(flights, annot=True, fmt="d")
17
```

File Handling



```
1 # Writing to a file
2 with open('example.txt', 'w', encoding='utf-8') as f:
3     f.write('Hello World!\n')
4     f.write('Second line\n')
5 # Reading from a file
6 with open('example.txt', 'r') as f:
7     content = f.read()      # Entire content
8     lines = f.readlines()  # List of lines
9 # Appending to a file
10 with open('example.txt', 'a') as f:
11     f.write('Appended line\n')
12
```

Best Practice

Always use `with` statement for automatic file closing

```
1 # Common file modes
2 with open('file.txt', 'r') as f:
3     # Read (default)
4     content = f.read()
5 with open('output.txt', 'w') as f:
6     # Write (truncate)
7     f.write("Hello\nWorld")
8 with open('data.bin', 'wb') as f:
9     # Binary write
10    f.write(bytes([65, 66, 67]))
11 with open('log.txt', 'a') as f:
12     # Append
13     f.write("New entry\n")
14 with open('data.txt', 'r+') as f:
15     # Read/write
16     content = f.read()
17     f.seek(0)
18     f.write("Updated")
```

Common Encodings

- ▶ utf-8 (recommended)
- ▶ ascii (basic English)
- ▶ latin-1 (Western Europe)
- ▶ utf-16 (alternative Unicode)

Encoding Issues

Always specify encoding when opening text files to avoid platform-dependent behavior

```
1 import struct
2 # Packing binary data
3
4 data = struct.pack('3i2f', 1, 2, 3, 4.5, 5.5)
5
6 # 3 ints, 2 floats
7 with open('data.bin', 'wb') as f:
8     f.write(data)
9
10 # Unpacking binary data
11 with open('data.bin', 'rb') as f:
12     data = f.read()
13     a, b, c, d, e = struct.unpack('3i2f', data)
14     print(f"Values: {a}, {b}, {c}, {d}, {e}")
15
16 # Working with bytes
17 bytes_data = b'Hello\x00World'
18 with open('binary.txt', 'wb') as f:
19     f.write(bytes_data)
```

Struct Format Characters

- i** integer
- f** float
- d** double
- s** char[]

```
1 import csv
2 # Reading CSV
3 with open('data.csv', 'r') as f:
4     reader = csv.reader(f)
5     header = next(reader) # Skip header
6     for row in reader:
7         print(row[0], float(row[1])) # Access columns
8 # Writing CSV
9 with open('output.csv', 'w', newline='') as f:
10    writer = csv.writer(f)
11    writer.writerow(['Name', 'Age', 'City']) # Header
12    writer.writerow(['Alice', 25, 'NY'])
13    writer.writerow(['Bob', 30, 'LA'])
14 # DictReader/DictWriter
15 with open('data.csv', 'r') as f:
16     reader = csv.DictReader(f)
17     for row in reader:
18         print(row['Name'], row['Age'])
```

Working with JSON

```
1 import json
2 # Writing JSON
3 data = {
4     'name': 'Alice',
5     'age': 25,
6     'courses': ['Math', 'Physics']
7 }
8 with open('data.json', 'w') as f:
9     json.dump(data, f, indent=2)
10 # Reading JSON
11 with open('data.json', 'r') as f:
12     loaded = json.load(f)
13     print(loaded['name'])
14 # Custom serialization
15 class Person:
16     def __init__(self, name, age):
17         self.name = name
18         self.age = age
19 def person_encoder(obj):
```

```
1 import xml.etree.ElementTree as ET
2 # Parsing XML
3 tree = ET.parse('data.xml')
4 root = tree.getroot()
5 # Accessing elements
6 for child in root:
7     print(child.tag, child.attrib)
8     for subchild in child:
9         print(subchild.text)
10 # Creating XML
11 root = ET.Element('catalog')
12 book = ET.SubElement(root, 'book', id='1')
13 title = ET.SubElement(book, 'title')
14 title.text = 'Python Programming'
15 ET.dump(root)
16 # Writing to file
17 tree = ET.ElementTree(root)
18 tree.write('output.xml', encoding='utf-8', xml_declaration=True)
```

Working with Excel Files

```
1 import openpyxl
2 # Reading Excel
3 wb = openpyxl.load_workbook('data.xlsx')
4 sheet = wb['Sheet1']
5
6 for row in sheet.iter_rows(values_only=True):
7     print(row[0], row[1])
8 # Writing Excel
9 wb = openpyxl.Workbook()
10 sheet = wb.active
11 sheet['A1'] = 'Name'
12 sheet['B1'] = 'Age'
13 sheet.append(['Alice', 25])
14 wb.save('output.xlsx')
15 # Using pandas
16 df = pd.read_excel('data.xlsx', sheet_name='Sheet1')
17 df.to_excel('output.xlsx', sheet_name='Results', index=False)
18
```

```
1 import os
2 import shutil
3
4 # Directory operations
5 os.mkdir('new_dir')
6 os.rmdir('empty_dir')
7 shutil.rmtree('dir_with_content') # Recursive delete
8
9 # File operations
10 os.rename('old.txt', 'new.txt')
11 os.remove('file.txt')
12 shutil.copy('source.txt', 'dest.txt')
13 shutil.move('source.txt', 'new_location/')
14
15
16
```

File System Operations [cont...]

```
1 import os
2 import shutil
3
4
5 # Path operations
6 path = os.path.join('dir', 'subdir', 'file.txt')
7 if os.path.exists(path):
8     print(f"Size: {os.path.getsize(path)} bytes")
9     print(f"Modified: {os.path.getmtime(path)}")
10
11 # Walking directory tree
12 for root, dirs, files in os.walk('.'):
13     for file in files:
14         print(os.path.join(root, file))
15
```

```
1 import zipfile  
2  
3 # Creating ZIP archive  
4 with zipfile.ZipFile('archive.zip', 'w') as zipf:  
5     zipf.write('file1.txt')  
6     zipf.write('file2.txt', arcname='renamed.txt')  
7 # Extracting ZIP  
8 with zipfile.ZipFile('archive.zip', 'r') as zipf:  
9     zipf.extractall('extracted')  
10 # Reading ZIP contents  
11 with zipfile.ZipFile('archive.zip', 'r') as zipf:  
12     print(zipf.namelist())  
13     with zipf.open('file1.txt') as f:  
14         content = f.read().decode('utf-8')  
15 # Adding compression  
16 with zipfile.ZipFile('compressed.zip', 'w', zipfile.ZIP_DEFLATED) as zipf:  
17     :  
18     zipf.write('large file.txt')
```

```
1 import configparser
2 # Writing config
3 config = configparser.ConfigParser()
4 config['DEFAULT'] = {'debug': 'True', 'log_level': 'INFO'}
5 config['DATABASE'] = {'host': 'localhost', 'port': '5432'}
6
7 with open('config.ini', 'w') as f:
8     config.write(f)
9 # Reading config
10 config = configparser.ConfigParser()
11 config.read('config.ini')
12
13 debug = config['DEFAULT'].getboolean('debug')
14 db_host = config['DATABASE']['host']
15 port = config['DATABASE'].getint('port')
16 # Using YAML (requires PyYAML)
17 import yaml
18 with open('config.yaml') as f:
19     config = yaml.safe_load(f)
```

Do's

- ▶ Use `with` statements for file handling
- ▶ Specify encoding for text files
- ▶ Handle exceptions (`FileNotFoundException`, `PermissionError`)
- ▶ Use `os.path` for path operations
- ▶ Close files explicitly when not using `with`

Don'ts

- ▶ Don't assume files exist - always check
- ▶ Don't use string concatenation for paths
- ▶ Don't ignore file closing
- ▶ Don't load huge files into memory all at once

Real-World Example: Log Parser

```
1 import re
2 from collections import defaultdict
3
4 log_pattern = re.compile(r'(\d{4}-\d{2}-\d{2}) (\d{2}:\d{2}:\d{2}) (\w+)
5   (.*)')
6
7 def parse_log_file(filename):
8     stats = defaultdict(int)
9     errors = []
10
11    with open(filename, 'r') as f:
12        for line in f:
13            match = log_pattern.match(line)
14            if match:
15                date, time, level, message = match.groups()
16                stats[level] += 1
17                if level == 'ERROR':
18                    errors.append((date, time, message))
```

Graphical User Interface (GUI) Development



```
1 import tkinter as tk
2 from tkinter import ttk
3
4 # Create main window
5 root = tk.Tk()
6 root.title("Simple GUI")
7 root.geometry("300x200")
8
9 # Add widgets
10 label = ttk.Label(root, text="Enter your name:")
11 entry = ttk.Entry(root)
12 button = ttk.Button(root, text="Greet",
13                      command=lambda: print(f"Hello, {entry.get()}!"))
14
15 # Layout widgets
16 label.pack(pady=10)
17 entry.pack(pady=5)
18 button.pack(pady=10)
```

```
1 # Label
2 label = tk.Label(root, text="Information", font=('Arial', 12))
3
4 # Button
5 button = tk.Button(root, text="Click", command=callback)
6
7 # Entry
8 entry = tk.Entry(root, show="*") # For passwords
9
10 # Text (multi-line)
11 text = tk.Text(root, height=5, width=30)
12
13 # Checkbutton
14 check = tk.Checkbutton(root, text="Agree", variable=tk.BooleanVar())
15
16 # Radiobutton
17 radio1 = tk.Radiobutton(root, text="Option 1", variable=var, value=1)
18 radio2 = tk.Radiobutton(root, text="Option 2", variable=var, value=2)
```

pack():

```
1 # Simple vertical layout
2 label.pack(side="top", fill="x",
            padx=5, pady=5)
3 button.pack(side="bottom")
4
```

grid():

```
1 # Table-like layout
2 label.grid(row=0, column=0, sticky="ew")
3 entry.grid(row=0, column=1)
4 button.grid(row=1, columnspan=2)
5
```

place():

```
1 # Absolute positioning
2 label.place(x=10, y=10)
3 button.place(relx=0.5, rely=0.5,
              anchor="center")
4
```

Best Practices

- ▶ Use grid for forms
- ▶ Use pack for simple layouts
- ▶ Avoid mixing managers in same container

```
1 # Button command
2 button = tk.Button(root, text="Click", command=callback)
3
4 # Binding events
5 def on_key(event):
6     print(f"Key pressed: {event.char}")
7
8 root.bind('<Key>', on_key)
9
10 # Event types:
11 # <Button-1> - Left mouse click
12 # <Return> - Enter key
13 # <Motion> - Mouse movement
14 # <Configure> - Window resize
15
16 # Event object attributes:
17 # widget - The widget that received the event
18 # x, y - Mouse coordinates
19 # char - Character for keyboard events
```

```
1 from tkinter import messagebox, filedialog, simpledialog
2
3 # Message boxes
4 messagebox.showinfo("Title", "Information")
5 messagebox.showwarning("Warning", "Be careful!")
6 messagebox.showerror("Error", "Something went wrong")
7 answer = messagebox.askyesno("Question", "Continue?")
8
9 # File dialogs
10 filename = filedialog.askopenfilename()
11 directory = filedialog.askdirectory()
12 save_as = filedialog.asksaveasfilename()
13
14 # Input dialogs
15 name = simpledialog.askstring("Input", "Enter your name")
16 age = simpledialog.askinteger("Input", "Enter your age")
17
```

```
1 class ValidatedEntry(ttk.Entry):
2     def __init__(self, parent, *args, **kwargs):
3         super().__init__(parent, *args, **kwargs)
4         self.configure(validate='key')
5         self.configure(validatecommand=
6             self.register(self.validate_input), '%P'))
7
8     def validate_input(self, proposed):
9         # Only allow digits
10        return proposed.isdigit() or proposed == ""
11
12 # Usage
13 entry = ValidatedEntry(root)
14 entry.pack()
15
16 # Custom composite widget
17 class LabelledEntry(ttk.Frame):
18     def __init__(self, parent, text):
19         super().__init__(parent)
```

```
1 from tkinter import ttk
2
3 # Using themed widgets
4 style = ttk.Style()
5 style.configure('TButton', font=('Arial', 12))
6 style.map('TButton',
7           foreground=[('pressed', 'red'), ('active', 'blue')],
8           background=[('pressed', '!disabled', 'black'), ('active', 'white',
9             ')])
10
11 # Available themes
12 print(style.theme_names()) # ('clam', 'alt', 'default', 'classic')
13 style.theme_use('clam')
14
15 # Custom style
16 style.configure('My.TButton', padding=10, relief="flat")
17 button = ttk.Button(root, text="Styled", style='My.TButton')
```

```
1 # Create menu bar
2 menubar = tk.Menu(root)
3 root.config(menu=menubar)
4
5 # File menu
6 file_menu = tk.Menu(menubar, tearoff=0)
7 file_menu.add_command(label="Open", command=open_file)
8 file_menu.add_command(label="Save", command=save_file)
9 file_menu.add_separator()
10 file_menu.add_command(label="Exit", command=root.quit)
11 menubar.add_cascade(label="File", menu=file_menu)
12
13 # Edit menu
14 edit_menu = tk.Menu(menubar, tearoff=0)
15 edit_menu.add_command(label="Cut")
16 edit_menu.add_command(label="Copy")
17 menubar.add_cascade(label="Edit", menu=edit_menu)
18
19 # Context menu
```

```
1 canvas = tk.Canvas(root, width=400, height=300)
2 canvas.pack()
3
4 # Drawing shapes
5 canvas.create_line(10, 10, 100, 100, fill="red", width=2)
6 canvas.create_rectangle(50, 50, 150, 150, fill="blue")
7 canvas.create_oval(100, 100, 200, 200, fill="green")
8 canvas.create_text(200, 50, text="Hello", font=('Arial', 14))
9
10 # Interactive drawing
11 def start_draw(event):
12     global last_x, last_y
13     last_x, last_y = event.x, event.y
14
15 def draw(event):
16     global last_x, last_y
17     canvas.create_line(last_x, last_y, event.x, event.y, width=2)
18     last_x, last_y = event.x, event.y
```

Real-World Example: Text Editor

```
1 class TextEditor(tk.Tk):
2     def __init__(self):
3         super().__init__()
4         self.title("Simple Text Editor")
5         self.geometry("600x400")
6
7     # Create widgets
8     self.text = tk.Text(self, wrap="word")
9     scrollbar = ttk.Scrollbar(self, command=self.text.yview)
10    self.text.config(yscrollcommand=scrollbar.set)
11
12    # Menu
13    menubar = tk.Menu(self)
14    file_menu = tk.Menu(menubar, tearoff=0)
15    file_menu.add_command(label="Open", command=self.open_file)
16    file_menu.add_command(label="Save", command=self.save_file)
17    menubar.add_cascade(label="File", menu=file_menu)
18    self.config(menu=menubar)
```

Do's

- ▶ Use `ttk` widgets for native look
- ▶ Organize code with classes
- ▶ Use geometry managers consistently
- ▶ Separate UI and business logic
- ▶ Use resource files for strings

Don'ts

- ▶ Don't use `place()` unless necessary
- ▶ Don't mix `tk` and `ttk` widgets arbitrarily
- ▶ Don't block the main thread
- ▶ Don't create widgets in a loop

```
1 import threading
2
3 class LongRunningTask:
4     def __init__(self, root):
5         self.root = root
6         self.progress = ttk.Progressbar(root, mode='indeterminate')
7         self.progress.pack()
8
9     def start(self):
10        self.progress.start()
11        thread = threading.Thread(target=self.run_task)
12        thread.start()
13        self.check_thread(thread)
14
15    def run_task(self):
16        # Simulate long task
17        import time
18        time.sleep(5)
```

Introduction to Machine Learning



Definition

Machine learning is the study of algorithms that improve automatically through experience, building models from sample data.

Types of ML:

- ▶ Supervised Learning
 - ▶ Classification
 - ▶ Regression
- ▶ Unsupervised Learning
 - ▶ Clustering
 - ▶ Dimensionality Reduction

Typical Workflow:

1. Data collection
2. Preprocessing Training
3. Evaluation
4. Deployment

Python Ecosystem

Scikit-learn, TensorFlow, PyTorch, Keras, XGBoost, LightGBM

```
1 from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score
6
7 # Load dataset
8 iris = load_iris()
9 X, y = iris.data, iris.target
10
11 # Split data
12 X_train, X_test, y_train, y_test = train_test_split(
13     X, y, test_size=0.2, random_state=42)
14
15 # Preprocessing
16 scaler = StandardScaler()
17 X_train = scaler.fit_transform(X_train)
18 X_test = scaler.transform(X_test)
```

```
1 from sklearn.preprocessing import (StandardScaler, MinMaxScaler,
2                                     OneHotEncoder, LabelEncoder)
3 from sklearn.impute import SimpleImputer
4 from sklearn.pipeline import Pipeline
5 from sklearn.compose import ColumnTransformer
6
7 # Numeric features
8 numeric_features = ['age', 'income']
9 numeric_transformer = Pipeline(steps=[
10     ('imputer', SimpleImputer(strategy='median')),
11     ('scaler', StandardScaler())])
12
13 # Categorical features
14 categorical_features = ['gender', 'country']
15 categorical_transformer = Pipeline(steps=[
16     ('imputer', SimpleImputer(strategy='constant', fill_value='missing'))
17     ,
18     ('onehot', OneHotEncoder(handle_unknown='ignore'))])
```

```
1 from sklearn.metrics import (accuracy_score, precision_score,
2                               recall_score, f1_score,
3                               confusion_matrix, classification_report)
4
5 # For classification
6 y_pred = model.predict(X_test)
7 print("Accuracy:", accuracy_score(y_test, y_pred))
8 print("Precision:", precision_score(y_test, y_pred, average='macro'))
9 print("Recall:", recall_score(y_test, y_pred, average='macro'))
10 print("F1:", f1_score(y_test, y_pred, average='macro'))
11 print(confusion_matrix(y_test, y_pred))
12 print(classification_report(y_test, y_pred))
13
14
```

```
1 from sklearn.metrics import (accuracy_score, precision_score,
2                               recall_score, f1_score,
3                               confusion_matrix, classification_report)
4
5
6 # For regression
7 from sklearn.metrics import mean_squared_error, r2_score
8 print("MSE:", mean_squared_error(y_test, y_pred))
9 print("R2:", r2_score(y_test, y_pred))
10
11 # Cross-validation
12 from sklearn.model_selection import cross_val_score
13 scores = cross_val_score(model, X, y, cv=5)
14 print("CV Accuracy: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std() *
15           2))
```

```
1 from sklearn.model_selection import GridSearchCV
2 # Define parameter grid
3 param_grid = {
4     'n_estimators': [50, 100, 200],
5     'max_depth': [None, 10, 20],
6     'min_samples_split': [2, 5, 10]}
7 # Create grid search
8 grid_search = GridSearchCV(
9     estimator=RandomForestClassifier(),
10    param_grid=param_grid,
11    cv=5,
12    n_jobs=-1, verbose=2)
13 # Fit grid search
14 grid_search.fit(X_train, y_train)
15 # Best parameters
16 print("Best params:", grid_search.best_params_)
17 print("Best score:", grid_search.best_score_)
18 # Use best model
19 best_model = grid_search.best_estimator_
```

```
1 from sklearn.feature_selection import SelectKBest, chi2
2 from sklearn.decomposition import PCA
3 from sklearn.feature_extraction.text import TfidfVectorizer
4
5 # Feature selection
6 selector = SelectKBest(chi2, k=10)
7 X_new = selector.fit_transform(X, y)
8
9 # Dimensionality reduction
10 pca = PCA(n_components=2)
11 X_pca = pca.fit_transform(X)
12
13 # Text feature extraction
14 corpus = ["This is good", "This is bad", "Excellent"]
15 vectorizer = TfidfVectorizer()
16 X_text = vectorizer.fit_transform(corpus)
```

```
1 from sklearn.feature_selection import SelectKBest, chi2
2 from sklearn.decomposition import PCA
3 from sklearn.feature_extraction.text import TfidfVectorizer
4
5
6
7 # Custom feature engineering
8 from sklearn.base import BaseEstimator, TransformerMixin
9
10 class CustomTransformer(BaseEstimator, TransformerMixin):
11     def fit(self, X, y=None):
12         return self
13     def transform(self, X):
14         return np.c_[X, X[:, 0]**2] # Add squared feature
15
```

```
1 import joblib
2 import pickle
3
4 # Save model with joblib (preferred)
5 joblib.dump(model, 'model.joblib')
6 loaded_model = joblib.load('model.joblib')
7
8 # Save model with pickle
9 with open('model.pkl', 'wb') as f:
10     pickle.dump(model, f)
11 with open('model.pkl', 'rb') as f:
12     loaded_model = pickle.load(f)
13
14 # Save entire pipeline
15 from sklearn.pipeline import Pipeline
16 pipe = Pipeline([
17     ('preprocessor', preprocessor),
18     ('classifier', model)
19 ])
```

```
1 import pandas as pd
2 from sklearn.ensemble import RandomForestRegressor
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import mean_absolute_error
5
6 # Load data
7 data = pd.read_csv('housing.csv')
8 X = data.drop('price', axis=1)
9 y = data['price']
10
11 # Split data
12 X_train, X_test, y_train, y_test = train_test_split(
13     X, y, test_size=0.2, random_state=42)
14
15 # Preprocessing (simplified)
16 X_train = pd.get_dummies(X_train)
17 X_test = pd.get_dummies(X_test)
18
19 # Align columns
```

Do's

- ▶ Split data before any processing
- ▶ Use pipelines for reproducibility
- ▶ Track experiments and parameters
- ▶ Monitor model performance over time
- ▶ Start with simple models

Don'ts

- ▶ Don't test on training data
- ▶ Don't ignore feature scaling
- ▶ Don't leak information from test set
- ▶ Don't use accuracy for imbalanced data

```
1 # Deep Learning with Keras
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense
4 model = Sequential([
5     Dense(64, activation='relu', input_shape=(10,)),
6     Dense(64, activation='relu'),
7     Dense(1)
8 ])
9 model.compile(optimizer='adam', loss='mse')
10 model.fit(X_train, y_train, epochs=10)
11 # Ensemble Methods
12 from sklearn.ensemble import VotingClassifier
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.svm import SVC
15 ensemble = VotingClassifier(estimators=[
16     ('lr', LogisticRegression()),
17     ('svm', SVC()),
18     ('rf', RandomForestClassifier())
19 ])
```

```
1 # Feature importance
2
3 importances = model.feature_importances_
4 indices = np.argsort(importances)[::-1]
5 for f in range(X_train.shape[1]):
6     print(f"{X_train.columns[indices[f]]}: {importances[indices[f]]}")
7
8 # SHAP values
9 import shap
10
11 explainer = shap.TreeExplainer(model)
12 shap_values = explainer.shap_values(X_test)
13 shap.summary_plot(shap_values, X_test)
14
15
```

```
1 # Partial dependence plots
2 from sklearn.inspection import plot_partial_dependence
3 plot_partial_dependence(model, X_train, features=[0, 1])
4
5
6 # LIME for local explanations
7 import lime
8 import lime.lime_tabular
9
10 explainer = lime.lime_tabular.LimeTabularExplainer(
11     X_train.values, feature_names=X_train.columns)
12 exp = explainer.explain_instance(X_test.iloc[0], model.predict_proba)
13 exp.show_in_notebook()
14
```

```
1 # Resampling techniques
2 from imblearn.over_sampling import SMOTE
3 from imblearn.under_sampling import RandomUnderSampler
4 from imblearn.pipeline import Pipeline
5
6 over = SMOTE(sampling_strategy=0.1)
7 under = RandomUnderSampler(sampling_strategy=0.5)
8 steps = [('o', over), ('u', under)]
9 pipeline = Pipeline(steps=steps)
10 X_resampled, y_resampled = pipeline.fit_resample(X, y)
11
12 # Class weights
13 model = RandomForestClassifier(class_weight='balanced')
14
15 # Appropriate metrics
16 from sklearn.metrics import balanced_accuracy_score, roc_auc_score
17 print("Balanced accuracy:", balanced_accuracy_score(y_test, y_pred))
18 print("ROC AUC:", roc_auc_score(y_test, y_pred_proba))
```

Time Complexity



Key Concepts

- ▶ **Time Complexity:** How runtime grows with input size
- ▶ **Space Complexity:** How memory usage grows with input size
- ▶ **Big-O Notation:** Describes upper bound of growth
- ▶ **Best/Average/Worst Case:** Different scenarios

Why It Matters

Input Size	$O(n^2)$ Time
10	100 operations
100	10,000 operations
1000	1,000,000 operations

Practical Impact

A $O(n^2)$ algorithm with 1M items could take days vs seconds for $O(n \log n)$

O(1) - Constant:

```
1 def get_first(items):
2     return items[0]
3
```

O(log n) - Logarithmic:

```
1 def binary_search(sorted_list,
2                     target):
3     low, high = 0, len(sorted_list)
4     -1
5     while low <= high:
6         mid = (low + high) // 2
7         if sorted_list[mid] ==
8             target:
9             return mid
10            elif sorted_list[mid] <
11                target:
12                    low = mid + 1
13                    else:
14                        high = mid - 1
15                    return -1
16
```

O(n) - Linear:

```
1 def find_max(items):
2     max_val = items[0]
3     for item in items:
4         if item > max_val:
5             max_val = item
6     return max_val
7
```

O(n²) - Quadratic:

```
1 def bubble_sort(items):
2     n = len(items)
3     for i in range(n):
4         for j in range(n-i-1):
5             if items[j] > items[j
+1]:
6                 items[j], items[j+1]
7 = items[j+1], items[j]
```

List Operations

Operation	Time Complexity
Indexing (lst[i])	O(1)
Append	O(1)
Insert	O(n)
Delete	O(n)
Search (x in lst)	O(n)
Sort	O(n log n)

Dict Operations

Operation	Time Complexity
Get item	O(1)
Set item	O(1)
Delete item	O(1)
Search (key in d)	O(1)

```
1 # O(n) - Linear
2 def linear(items):
3     total = 0
4     for item in items:          # n iterations
5         total += item          # O(1)
6     return total                # Total: O(n)
7
8 # O(n2) - Quadratic
9 def quadratic(items):
10    pairs = []
11    for x in items:           # n iterations
12        for y in items:       # n iterations
13            pairs.append((x, y)) # O(1)
14    return pairs               # Total: O(n2)
15
16
```

```
1
2 # O(n log n) - Linearithmic
3 def linearithmic(items):
4     items.sort()                      # O(n log n)
5     result = []
6     for item in items:                # O(n)
7         result.append(item*2)        # O(1)
8     return result                    # Total: O(n log n)
9
```

```
1 # O(1) - Constant space
2 def constant_space(n):
3     total = 0
4     for i in range(n):
5         total += i
6     return total
7
8 # O(n) - Linear space
9 def linear_space(n):
10    numbers = []
11    for i in range(n):
12        numbers.append(i)
13    return numbers
```

```
1  
2  
3 # O(n2) - Quadratic space  
4 def quadratic_space(n):  
5     matrix = []  
6     for i in range(n):  
7         row = []  
8         for j in range(n):  
9             row.append(i*j)  
10        matrix.append(row)  
11    return matrix  
12
```

Before:

```
1 def contains_duplicate(nums):
2     for i in range(len(nums)):
3         for j in range(i+1, len(nums)):
4             if nums[i] == nums[j]:
5                 return True
6     return False
7 # O(n2) time, O(1) space
8
```

After:

```
1 def contains_duplicate(nums):
2     seen = set()
3     for num in nums:
4         if num in seen:
5             return True
6         seen.add(num)
7     return False
8 # O(n) time, O(n) space
9
```

Tradeoffs

- ▶ Time vs space
- ▶ Readability vs performance
- ▶ Best case vs worst case

```
1 import timeit
2 import random
3 from functools import lru_cache
4 # Time measurement
5 def test_func():
6     return sum(range(1000000))
7
8 time = timeit.timeit(test_func, number=100)
9 print(f"Average time: {time/100:.6f} seconds")
10 # Profiling
11 import cProfile
12 cProfile.run('test_func()')
13 # Memoization
14 @lru_cache(maxsize=None)
15 def fibonacci(n):
16     if n < 2:
17         return n
18     return fibonacci(n-1) + fibonacci(n-2)
19 # Reduces from O(2^n) to O(n)
```

Techniques

- ▶ Memoization (caching)
- ▶ Using more efficient data structures
- ▶ Early termination
- ▶ Precomputing values
- ▶ Vectorization (NumPy)

```
1 # Vectorized vs loop
2 import numpy as np
3
4 def slow_dot(a, b):
5     result = 0
6     for x, y in zip(a, b):
7         result += x * y
8     return result
9
```

Naive Approach:

```
1 def two_sum(nums, target):
2     for i in range(len(nums)):
3         for j in range(i+1, len(nums)):
4             if nums[i] + nums[j] == target:
5                 return [i, j]
6     return []
7 # O(n2) time, O(1) space
8
```

Optimized:

```
1 def two_sum(nums, target):
2     seen = {}
3     for i, num in enumerate(nums):
4         complement = target - num
5         if complement in seen:
6             return [seen[complement], i]
7         seen[num] = i
8     return []
9 # O(n) time, O(n) space
10
```

Performance

For n=10,000:

- ▶ Naive: 50 million operations
- ▶ Optimized: 10,000 operations

Common Algorithms

Algorithm	Time	Space
Linear Search	$O(n)$	$O(1)$
Binary Search	$O(\log n)$	$O(1)$
Bubble Sort	$O(n^2)$	$O(1)$
Merge Sort	$O(n \log n)$	$O(n)$
Quick Sort	$O(n \log n)$	$O(\log n)$
Dijkstra's	$O((V+E) \log V)$	$O(V)$

Data Structures

Operation	List	Dict
Access	$O(1)$	$O(1)$
Search	$O(n)$	$O(1)$
Insert	$O(n)$	$O(1)$
Delete	$O(n)$	$O(1)$

Optimization Guidelines

- ▶ First make it work, then make it fast
- ▶ Profile before optimizing
- ▶ Focus on bottlenecks (90/10 rule)
- ▶ Consider readability/maintainability

```
1 # Example: Choosing data structures
2 # Need fast lookups? Use dict/set
3 # Need ordered data? Use list
4 # Need both? Consider OrderedDict
5 from collections import OrderedDict, defaultdict
6 # Example: Using built-ins
7 # Instead of:
8 result = []
9 for item in items:
10     result.append(func(item))
11 # Use:
```

```
1 # Common pitfalls
2 # 1. Hidden loops
3 def contains_duplicate(nums):
4     return len(nums) != len(set(nums))    # O(n)
5
6 # 2. String concatenation
7 result = ""
8 for s in strings:  # O(n2) - strings are immutable
9     result += s
10 # Better:
11 result = "".join(strings)  # O(n)
12
13
14
```

```
1
2
3 # 3. Unnecessary computations
4 def sum_of_squares(n):
5     total = 0
6     for i in range(n):
7         for j in range(n): # O(n2) when O(n) possible
8             total += i*j
9     return total
10 # Better:
11 def sum_of_squares(n):
12     sum_i = sum(range(n))
13     return sum_i * sum_i # O(n)
```

Conclusion



Covered Topics:

- ▶ Exception handling
- ▶ NumPy arrays
- ▶ Pandas data analysis
- ▶ Matplotlib visualization
- ▶ File I/O operations
- ▶ GUI development
- ▶ ML fundamentals
- ▶ Time complexity

Key Skills:

- ▶ Writing robust Python code
- ▶ Efficient data processing
- ▶ Creating visualizations
- ▶ Building applications
- ▶ Algorithm analysis

Further Learning

- ▶ Advanced Python features
- ▶ Web frameworks (Django, Flask)
- ▶ Parallel programming
- ▶ Cloud deployment

Thank You!

Questions and Discussion

