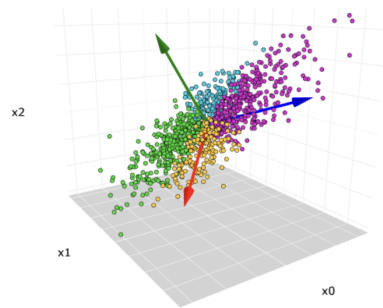# Principal Component Analysis (PCA) in Python

DATASCI 347: Machine Learning (Fall 2025)

Instructor: Alexander Williams Tolbert



September 14, 2025

# Contents

# 1 Introduction to Principal Component Analysis

## 1.1 What is PCA?

Principal Component Analysis is a powerful statistical technique that helps us:

- **Dimension Reduction**: Take many variables and create fewer "summary" variables
- **Noise Reduction**: Focus on the "signal" and ignore the "noise"
- **Visualization**: Display high-dimensional data in 2D or 3D plots
- **Grouping/Clustering**: Reveal hidden structure in your data

## 1.2 Real-World Example

Imagine you're a military recruiter. You have 10 different test scores for each applicant:

- Science, Math, Word Knowledge, Paragraph Comprehension
- Numerical Operations, Coding Speed, Auto Shop, Electronics
- Mechanical Comprehension, Arithmetic Reasoning

PCA helps answer questions like:

- Can we summarize all 10 scores with just 2 or 3 numbers?
- Are there different "types" of intelligence?
- Do men and women have different strengths?

# 2 Setting Up Python Environment

## 2.1 Installing and Importing Libraries

We need several Python libraries for PCA and data analysis. In Google Colab, most are already installed, but we'll check:

```python
import subprocess
import sys

def install_if_missing(package):
    """Helper function to install packages if they're not already
    installed"""
    try:
        __import__(package)
    except ImportError:
        subprocess.check_call([sys.executable, "-m", "pip", "install",
    package])

```

```
11 # Install required packages
12 packages = ['pandas', 'numpy', 'matplotlib', 'seaborn', 'scikit-learn', '
      scipy']
13 for package in packages:
14     install_if_missing(package)
```

Listing 1: Installing Required Packages

```
1  import pandas as pd  # For data manipulation (like data frames)
2  import numpy as np  # For numerical operations
3  import matplotlib.pyplot as plt  # For basic plotting
4  import seaborn as sns  # For beautiful statistical plots
5  from sklearn.decomposition import PCA  # For Principal Component Analysis
6  from sklearn.preprocessing import StandardScaler  # For scaling data
7  from google.colab import files  # For uploading files to Colab
8  from scipy import stats  # For statistical tests
9
10 # Set up plotting parameters for nice-looking graphs
11 plt.style.use('seaborn-v0_8')  # Use seaborn's clean style
12 plt.rcParams['figure.figsize'] = (10, 6)  # Default figure size
13 plt.rcParams['font.size'] = 12  # Default font size
14 sns.set_palette("husl")  # Use a nice color palette
15
16 print("All libraries imported successfully!")
17 print("Ready for Principal Component Analysis!")
```

Listing 2: Importing Libraries

# 3 Data Acquisition - Getting the ASVAB Dataset

## 3.1 About the ASVAB Dataset

The Armed Services Vocational Aptitude Battery (ASVAB) is a test used by the U.S. military to determine job assignments for recruits. The dataset contains:

- 10 different test scores for each person

- Gender information

- Family background variables

- Income and self-esteem measures

## 3.2 Uploading and Loading Data

```
1  print("Please upload your IQ.Full.csv file:")
2  uploaded = files.upload()
3
4  # Verify the file was uploaded
5  import os
```

```
6 print(f"Files in current directory: {os.listdir('.')}")
```

Listing 3: Uploading Data File

```python
1  # Load the data
2  df_full = pd.read_csv('IQ.Full.csv')
3
4  # Let's examine what we have
5  print("DATASET OVERVIEW:")
6  print(f"Shape: {df_full.shape[0]} rows (people)    {df_full.shape[1]}
       columns (variables)")
7  print("\nColumn names:")
8  print(df_full.columns.tolist())
9  print("\nFirst few rows:")
10 print(df_full.head())
11 print("\nData types:")
12 print(df_full.dtypes)
13 print("\nBasic statistics:")
14 print(df_full.describe())
```

Listing 4: Loading and Examining Data

# 4 Data Preparation for PCA

## 4.1 Extracting Test Score Variables

```python
1  # From the original R code, the 10 ASVAB tests are:
2  test_columns = ['Science', 'Arith', 'Word', 'Parag', 'Numer', 'Coding',
3                  'Auto', 'Math', 'Mechanic', 'Elec']
4
5  # Create a subset with just the test scores
6  df_tests = df_full[test_columns].copy()
7
8  print("TEST SCORES DATASET:")
9  print(f"Shape: {df_tests.shape}")
10 print("\nSummary statistics:")
11 print(df_tests.describe().round(2))
12
13 # Check for missing values
14 missing_values = df_tests.isnull().sum()
15 print(f"\nMissing values:\n{missing_values}")
```

Listing 5: Selecting Test Score Columns

## 4.2 Examining Correlations

```python
1  # Calculate correlation matrix
2  correlation_matrix = df_tests.corr()
3  print("CORRELATION MATRIX:")
4  print(correlation_matrix.round(2))
```

```
 5
 6  # Create a heatmap to visualize correlations
 7  plt.figure(figsize=(12, 10))
 8  sns.heatmap(correlation_matrix,
 9              annot=True,  # Show correlation values
10              cmap='coolwarm',  # Color scheme
11              center=0,  # Center colormap at 0
12              square=True,  # Make cells square
13              fmt='.2f')  # Format numbers to 2 decimals
14  plt.title('Correlation Matrix of ASVAB Test Scores')
15  plt.tight_layout()
16  plt.show()
```

Listing 6: Correlation Analysis

# 5 Understanding PCA with a Simple Example

## 5.1 Two-Variable PCA Example

```
1  # Create a subset with just Word and Parag scores
2  # Let's take a random sample to make visualization clearer
3  np.random.seed(42)  # For reproducible results
4  sample_indices = np.random.choice(df_full.index, size=50, replace=False)
5  df_sample = df_full.loc[sample_indices, ['Word', 'Parag', 'Gender']].copy
       ()
6
7  print("WORD vs PARAGRAPH COMPREHENSION:")
8  print(f"Sample size: {len(df_sample)} people")
9  print(df_sample.describe())
```

Listing 7: Two-Variable PCA

## 5.2 Visualizing the Raw Data

```
1  plt.figure(figsize=(12, 5))
2
3  # Original data
4  plt.subplot(1, 2, 1)
5  plt.scatter(df_sample['Word'], df_sample['Parag'], alpha=0.7)
6  plt.xlabel('Word Knowledge Score')
7  plt.ylabel('Paragraph Comprehension Score')
8  plt.title('Original Data: Word vs Paragraph')
9  plt.grid(True, alpha=0.3)
10
11  # Add mean lines
12  plt.axvline(df_sample['Word'].mean(), color='red', linestyle='--', alpha
       =0.5, label='Word Mean')
13  plt.axhline(df_sample['Parag'].mean(), color='red', linestyle='--', alpha
       =0.5, label='Parag Mean')
14  plt.legend()
```

```
15
16 # Centered data (subtract means)
17 word_centered = df_sample['Word'] - df_sample['Word'].mean()
18 parag_centered = df_sample['Parag'] - df_sample['Parag'].mean()
19
20 plt.subplot(1, 2, 2)
21 plt.scatter(word_centered, parag_centered, alpha=0.7)
22 plt.xlabel('Word Knowledge (Centered)')
23 plt.ylabel('Paragraph Comprehension (Centered)')
24 plt.title('Centered Data')
25 plt.grid(True, alpha=0.3)
26 plt.axhline(0, color='black', linewidth=0.5)
27 plt.axvline(0, color='black', linewidth=0.5)
28
29 plt.tight_layout()
30 plt.show()
```
Listing 8: Data Visualization

## 5.3 Performing PCA on Two Variables

```
1 # Prepare data for PCA
2 data_2d = df_sample[['Word', 'Parag']].values
3
4 # Initialize PCA
5 pca_2d = PCA()
6
7 # Fit PCA (this finds the principal components)
8 pca_2d.fit(data_2d)
9
10 # Get the results
11 loadings_2d = pca_2d.components_   # The "directions" of the components
12 scores_2d = pca_2d.transform(data_2d)   # The new coordinates
13 variance_explained = pca_2d.explained_variance_ratio_
14
15 print("PCA RESULTS FOR WORD vs PARAGRAPH:")
16 print(f"PC1 loadings: {loadings_2d[0]}")
17 print(f"PC2 loadings: {loadings_2d[1]}")
18 print(f"Variance explained by PC1: {variance_explained[0]:.1%}")
19 print(f"Variance explained by PC2: {variance_explained[1]:.1%}")
```
Listing 9: PCA on Two Variables

## 5.4 Visualizing Principal Components

```
1 # Center the data for visualization
2 scaler = StandardScaler(with_std=False)   # Only center, don't scale
3 data_centered = scaler.fit_transform(data_2d)
4
5 plt.figure(figsize=(10, 8))
6 plt.scatter(data_centered[:, 0], data_centered[:, 1], alpha=0.7, s=50)
```

```
7
8  # Plot the principal component directions
9  origin = np.array([0, 0])
10 pc1_direction = loadings_2d[0] * 3 * np.sqrt(pca_2d.explained_variance_
     [0])
11 pc2_direction = loadings_2d[1] * 3 * np.sqrt(pca_2d.explained_variance_
     [1])
12
13 plt.arrow(0, 0, pc1_direction[0], pc1_direction[1],
14          head_width=0.5, head_length=0.3,
15          fc='red', ec='red', linewidth=3,
16          label=f'PC1 ({variance_explained[0]:.1%} variance)')
17
18 plt.arrow(0, 0, pc2_direction[0], pc2_direction[1],
19          head_width=0.5, head_length=0.3,
20          fc='blue', ec='blue', linewidth=3,
21          label=f'PC2 ({variance_explained[1]:.1%} variance)')
22
23 plt.xlabel('Word Knowledge (Centered)')
24 plt.ylabel('Paragraph Comprehension (Centered)')
25 plt.title('Principal Components: Word vs Paragraph')
26 plt.grid(True, alpha=0.3)
27 plt.legend()
28 plt.axis('equal')
29 plt.show()
```

Listing 10: Visualizing PCA Results

# 6 Full PCA Analysis - All 10 Test Scores

## 6.1 Preparing the Full Dataset

```
1 print("FULL PCA ANALYSIS - ALL 10 TEST SCORES")
2 print(f"Original data shape: {df_tests.shape}")
3 print(f"Goal: Reduce 10 variables to 2-3 meaningful components")
4
5 # Check data quality
6 print(f"\nData quality check:")
7 print(f"Missing values: {df_tests.isnull().sum().sum()}")
8 print(f"Data types: {df_tests.dtypes.unique()}")
```

Listing 11: Full Dataset Preparation

## 6.2 Standardizing the Data

```
1 print("SCALE ANALYSIS:")
2 print("Standard deviations of each test:")
3 print(df_tests.std().round(2))
4
5 print("\nRanges of each test:")
```

```
6  for col in df_tests.columns:
7      min_val = df_tests[col].min()
8      max_val = df_tests[col].max()
9      print(f"{col:10}: {min_val:6.1f} to {max_val:6.1f} (range: {max_val-
       min_val:6.1f})")
10
11 # Standardize the data
12 scaler = StandardScaler()
13 df_tests_scaled = scaler.fit_transform(df_tests)
14
15 # Convert back to DataFrame for easier handling
16 df_tests_scaled = pd.DataFrame(df_tests_scaled,
17                                columns=df_tests.columns,
18                                index=df_tests.index)
19
20 print(f"\nData standardized!")
21 print("After standardization - all variables have:")
22 print(f"Mean    0: {df_tests_scaled.mean().round(3).tolist()}")
23 print(f"Std     1: {df_tests_scaled.std().round(3).tolist()}")
```

Listing 12: Data Standardization

## 6.3   Performing PCA

```
1  # Initialize PCA - we'll keep all components first to see the full picture
2  pca_full = PCA()
3
4  # Fit PCA to our standardized data
5  pca_full.fit(df_tests_scaled)
6
7  # Transform the data to get PC scores
8  pc_scores = pca_full.transform(df_tests_scaled)
9
10 # Convert to DataFrame for easier handling
11 pc_scores_df = pd.DataFrame(pc_scores,
12                             columns=[f'PC{i+1}' for i in range(pc_scores.
       shape[1])],
13                             index=df_tests.index)
14
15 print(f"PCA COMPLETED!")
16 print(f"PC scores shape: {pc_scores_df.shape}")
17 print(f"We now have {pc_scores.shape[1]} principal components")
```

Listing 13: Performing PCA

## 6.4   Analyzing Variance Explained

```
1  variance_explained = pca_full.explained_variance_ratio_
2  cumulative_variance = np.cumsum(variance_explained)
3
4  print("VARIANCE EXPLAINED BY EACH COMPONENT:")
```

```
5 for i, (var, cum_var) in enumerate(zip(variance_explained,
      cumulative_variance)):
6   print(f"PC{i+1}: {var:.3f} ({var*100:.1f}%) - Cumulative: {cum_var:.3f
      } ({cum_var*100:.1f}%)")
```

Listing 14: Variance Analysis

## 6.5  Scree Plot

```
1 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))
2
3 # Individual variance explained
4 ax1.plot(range(1, len(variance_explained) + 1), variance_explained,
5          'bo-', linewidth=2, markersize=8)
6 ax1.set_xlabel('Principal Component')
7 ax1.set_ylabel('Proportion of Variance Explained')
8 ax1.set_title('Scree Plot: Individual Variance')
9 ax1.grid(True, alpha=0.3)
10 ax1.set_xticks(range(1, len(variance_explained) + 1))
11
12 # Add percentage labels
13 for i, var in enumerate(variance_explained):
14     ax1.annotate(f'{var*100:.1f}%', (i+1, var),
15                  textcoords="offset points", xytext=(0,10), ha='center')
16
17 # Cumulative variance explained
18 ax2.plot(range(1, len(cumulative_variance) + 1), cumulative_variance,
19          'ro-', linewidth=2, markersize=8)
20 ax2.set_xlabel('Number of Components')
21 ax2.set_ylabel('Cumulative Proportion of Variance Explained')
22 ax2.set_title('Cumulative Variance Explained')
23 ax2.grid(True, alpha=0.3)
24 ax2.set_xticks(range(1, len(cumulative_variance) + 1))
25 ax2.axhline(y=0.8, color='green', linestyle='--', alpha=0.7, label='80%
      threshold')
26 ax2.legend()
27
28 # Add percentage labels
29 for i, cum_var in enumerate(cumulative_variance):
30     ax2.annotate(f'{cum_var*100:.1f}%', (i+1, cum_var),
31                  textcoords="offset points", xytext=(0,10), ha='center')
32
33 plt.tight_layout()
34 plt.show()
35
36 # How many components for 80% and 90% variance?
37 pct_80 = np.argmax(cumulative_variance >= 0.8) + 1
38 pct_90 = np.argmax(cumulative_variance >= 0.9) + 1
39
40 print(f"COMPONENT SELECTION GUIDANCE:")
41 print(f"For 80% of variance: Keep {pct_80} components")
42 print(f"For 90% of variance: Keep {pct_90} components")
```

```
43  print(f"PC1 + PC2 explain: {cumulative_variance[1]*100:.1f}% of total
        variance")
```

Listing 15: Creating Scree Plot

# 7 Understanding the Components - Loadings Analysis

## 7.1 Extracting and Displaying Loadings

```
1  # Get the loadings (components)
2  loadings = pca_full.components_
3
4  # Create a DataFrame for easier interpretation
5  loadings_df = pd.DataFrame(loadings.T,  # Transpose so variables are rows
6                             columns=[f'PC{i+1}' for i in range(loadings.
       shape[0])],
7                             index=df_tests.columns)
8
9  print("PRINCIPAL COMPONENT LOADINGS:")
10 print("(How much each test contributes to each component)")
11 print(loadings_df.round(3))
12
13 # Focus on the first 3 components
14 print("\nDETAILED ANALYSIS - First 3 Components:")
15 for pc in ['PC1', 'PC2', 'PC3']:
16     print(f"\n{pc} (explains {variance_explained[int(pc[2:])-1]*100:.1f}%
       of variance):")
17
18     # Sort by absolute loading values
19     pc_loadings = loadings_df[pc].sort_values(key=abs, ascending=False)
20
21     for test, loading in pc_loadings.items():
22         direction = "+" if loading > 0 else "-"
23         print(f" {test:10}: {loading:6.3f} ({direction})")
```

Listing 16: Loadings Analysis

## 7.2 Visualizing Loadings

```
1  # Create a heatmap of loadings for the first 4 components
2  plt.figure(figsize=(12, 8))
3  sns.heatmap(loadings_df.iloc[:, :4],  # First 4 components
4              annot=True,
5              fmt='.2f',
6              cmap='RdBu_r',
7              center=0,
8              xticklabels=[f'PC{i+1}\n({variance_explained[i]*100:.1f}%)'
       for i in range(4)],
9              yticklabels=df_tests.columns)
10 plt.title('Principal Component Loadings Heatmap')
```

```
11  plt.ylabel('Test Scores')
12  plt.xlabel('Principal Components')
13  plt.tight_layout()
14  plt.show()
15
16  # Bar plots for PC1 and PC2
17  fig, (ax1, ax2) = plt.subplots(2, 1, figsize=(12, 10))
18
19  # PC1 loadings
20  pc1_loadings = loadings_df['PC1'].sort_values()
21  ax1.barh(range(len(pc1_loadings)), pc1_loadings.values,
22          color=['red' if x < 0 else 'blue' for x in pc1_loadings.values])
23  ax1.set_yticks(range(len(pc1_loadings)))
24  ax1.set_yticklabels(pc1_loadings.index)
25  ax1.set_xlabel('Loading Value')
26  ax1.set_title(f'PC1 Loadings (explains {variance_explained[0]*100:.1f}% of
        variance)')
27  ax1.grid(True, alpha=0.3)
28  ax1.axvline(0, color='black', linewidth=0.5)
29
30  # PC2 loadings
31  pc2_loadings = loadings_df['PC2'].sort_values()
32  ax2.barh(range(len(pc2_loadings)), pc2_loadings.values,
33          color=['red' if x < 0 else 'blue' for x in pc2_loadings.values])
34  ax2.set_yticks(range(len(pc2_loadings)))
35  ax2.set_yticklabels(pc2_loadings.index)
36  ax2.set_xlabel('Loading Value')
37  ax2.set_title(f'PC2 Loadings (explains {variance_explained[1]*100:.1f}% of
        variance)')
38  ax2.grid(True, alpha=0.3)
39  ax2.axvline(0, color='black', linewidth=0.5)
40
41  plt.tight_layout()
42  plt.show()
```

Listing 17: Visualizing Loadings

## 7.3  Interpreting the Components

```
1  print("COMPONENT INTERPRETATION GUIDE:")
2
3  print("\nPC1 Analysis:")
4  pc1_mean = loadings_df['PC1'].mean()
5  pc1_std = loadings_df['PC1'].std()
6  print(f"Average loading: {pc1_mean:.3f}")
7  print(f"Standard deviation: {pc1_std:.3f}")
8
9  if pc1_std < 0.1:  # All loadings are similar
10     print("PC1 = GENERAL ABILITY (all tests contribute equally)")
11  else:
12     print("PC1 = Mixed pattern (some tests more important)")
13
14  print(f"\nPC2 Analysis:")
```

```
15  pc2_positive = loadings_df[loadings_df['PC2'] > 0]['PC2']
16  pc2_negative = loadings_df[loadings_df['PC2'] < 0]['PC2']
17
18  print("Positive loadings (contribute positively to PC2):")
19  for test, loading in pc2_positive.items():
20      print(f" {test}: {loading:.3f}")
21
22  print("Negative loadings (contribute negatively to PC2):")
23  for test, loading in pc2_negative.items():
24      print(f" {test}: {loading:.3f}")
```

Listing 18: Component Interpretation

# 8 Principal Component Scores Analysis

## 8.1 Examining PC Scores

```
1  print("PRINCIPAL COMPONENT SCORES:")
2  print(f"Shape: {pc_scores_df.shape}")
3  print("\nFirst 10 people's PC scores:")
4  print(pc_scores_df.head(10).round(2))
5
6  print("\nPC Score Statistics:")
7  print(pc_scores_df.describe().round(2))
8
9  # Verify orthogonality
10 print("\nVERIFYING PC PROPERTIES:")
11 print("PC score means (should be    0):")
12 print(pc_scores_df.mean().round(3))
13
14 print("\nCorrelations between PC scores (should be    0):")
15 pc_correlations = pc_scores_df.corr()
16 print(pc_correlations.round(3))
```

Listing 19: PC Scores Analysis

## 8.2 Scatter Plot of PC Scores

```
1  # Add gender information for coloring
2  pc_scores_with_gender = pc_scores_df.copy()
3  pc_scores_with_gender['Gender'] = df_full.loc[pc_scores_df.index, 'Gender'
       ]
4
5  # Create the scatter plot
6  plt.figure(figsize=(12, 8))
7  colors = ['blue', 'red']
8  genders = pc_scores_with_gender['Gender'].unique()
9
10 for i, gender in enumerate(genders):
11     mask = pc_scores_with_gender['Gender'] == gender
```

```
12      plt.scatter(pc_scores_with_gender.loc[mask, 'PC1'],
13                  pc_scores_with_gender.loc[mask, 'PC2'],
14                  c=colors[i], alpha=0.6, s=50, label=f'{gender}')
15
16  plt.xlabel(f'PC1 ({variance_explained[0]*100:.1f}% of variance)')
17  plt.ylabel(f'PC2 ({variance_explained[1]*100:.1f}% of variance)')
18  plt.title('Principal Component Scores: PC2 vs PC1')
19  plt.grid(True, alpha=0.3)
20  plt.axhline(0, color='black', linewidth=0.5)
21  plt.axvline(0, color='black', linewidth=0.5)
22  plt.legend()
23
24  # Add quadrant labels to help interpretation
25  plt.text(2, 2, 'High PC1\nHigh PC2', ha='center', va='center',
26           bbox=dict(boxstyle='round', facecolor='wheat', alpha=0.8))
27  plt.text(-2, 2, 'Low PC1\nHigh PC2', ha='center', va='center',
28           bbox=dict(boxstyle='round', facecolor='lightblue', alpha=0.8))
29  plt.text(2, -2, 'High PC1\nLow PC2', ha='center', va='center',
30           bbox=dict(boxstyle='round', facecolor='lightgreen', alpha=0.8))
31  plt.text(-2, -2, 'Low PC1\nLow PC2', ha='center', va='center',
32           bbox=dict(boxstyle='round', facecolor='lightcoral', alpha=0.8))
33
34  plt.tight_layout()
35  plt.show()
36
37  # Statistical test for gender differences
38  pc1_male = pc_scores_with_gender[pc_scores_with_gender['Gender'] == 'male'
      ]['PC1']
39  pc1_female = pc_scores_with_gender[pc_scores_with_gender['Gender'] == '
      female']['PC1']
40  pc2_male = pc_scores_with_gender[pc_scores_with_gender['Gender'] == 'male'
      ]['PC2']
41  pc2_female = pc_scores_with_gender[pc_scores_with_gender['Gender'] == '
      female']['PC2']
42
43  print("\nGENDER DIFFERENCES IN PC SCORES:")
44  print(f"PC1 - Males: mean = {pc1_male.mean():.3f}, std = {pc1_male.std()
      :.3f}")
45  print(f"PC1 - Females: mean = {pc1_female.mean():.3f}, std = {pc1_female.
      std():.3f}")
46  print(f"PC2 - Males: mean = {pc2_male.mean():.3f}, std = {pc2_male.std()
      :.3f}")
47  print(f"PC2 - Females: mean = {pc2_female.mean():.3f}, std = {pc2_female.
      std():.3f}")
48
49  # T-tests
50  _, p_pc1 = stats.ttest_ind(pc1_male, pc1_female)
51  _, p_pc2 = stats.ttest_ind(pc2_male, pc2_female)
52
53  print(f"\nStatistical tests (t-test p-values):")
54  print(f"PC1 gender difference: p = {p_pc1:.4f}")
55  print(f"PC2 gender difference: p = {p_pc2:.4f}")
```

Listing 20: PC Scores Visualization

# 9 Biplot - Combining Scores and Loadings

## 9.1 Creating a Biplot

```python
def create_biplot(scores, loadings, labels, pc1=0, pc2=1, alpha=0.5):
    """
    Create a biplot showing both PC scores and variable loadings

    Parameters:
    scores: PC scores (people's positions)
    loadings: PC loadings (variable contributions)
    labels: Variable names
    pc1, pc2: Which components to plot (default: PC1 vs PC2)
    alpha: Transparency for points
    """
    plt.figure(figsize=(12, 10))

    # Plot PC scores (people)
    plt.scatter(scores[:, pc1], scores[:, pc2], alpha=alpha, s=30, color='blue')

    # Plot loading vectors (variables)
    loading_scale = 3  # Scale factor to make arrows visible

    for i, (var, loading) in enumerate(zip(labels, loadings)):
        plt.arrow(0, 0, loading[pc1] * loading_scale, loading[pc2] * loading_scale,
                  head_width=0.1, head_length=0.1,
                  fc='red', ec='red', alpha=0.8, linewidth=2)

        # Add variable labels
        plt.text(loading[pc1] * loading_scale * 1.1,
                 loading[pc2] * loading_scale * 1.1,
                 var, ha='center', va='center', fontsize=10,
                 bbox=dict(boxstyle='round,pad=0.3', facecolor='yellow', alpha=0.7))

    plt.xlabel(f'PC{pc1+1} ({variance_explained[pc1]*100:.1f}% variance)')
    plt.ylabel(f'PC{pc2+1} ({variance_explained[pc2]*100:.1f}% variance)')
    plt.title(f'Biplot: PC{pc2+1} vs PC{pc1+1}')
    plt.grid(True, alpha=0.3)
    plt.axhline(0, color='black', linewidth=0.5)
    plt.axvline(0, color='black', linewidth=0.5)

    return plt.gcf()

# Create the biplot
biplot_fig = create_biplot(pc_scores, loadings.T, df_tests.columns)
plt.tight_layout()
plt.show()
```

Listing 21: Biplot Creation

## 9.2 Enhanced Biplot with Gender Coloring

```python
def create_biplot_with_groups(scores, loadings, labels, groups,
    group_labels, pc1=0, pc2=1, alpha=0.7):
    """
    Create a biplot with different colors for different groups
    """
    plt.figure(figsize=(14, 10))

    # Color scheme
    colors = ['blue', 'red', 'green', 'purple', 'orange']
    unique_groups = np.unique(groups)

    # Plot PC scores by group
    for i, group in enumerate(unique_groups):
        mask = groups == group
        plt.scatter(scores[mask, pc1], scores[mask, pc2],
                    alpha=alpha, s=50, c=colors[i % len(colors)],
                    label=f'{group_labels} = {group}')

    # Plot loading vectors
    loading_scale = 4
    for i, (var, loading) in enumerate(zip(labels, loadings)):
        plt.arrow(0, 0, loading[pc1] * loading_scale, loading[pc2] *
    loading_scale,
                  head_width=0.15, head_length=0.1,
                  fc='black', ec='black', alpha=0.8, linewidth=2)

        plt.text(loading[pc1] * loading_scale * 1.15,
                 loading[pc2] * loading_scale * 1.15,
                 var, ha='center', va='center', fontsize=11, fontweight='
    bold',
                 bbox=dict(boxstyle='round,pad=0.3', facecolor='white',
    edgecolor='black', alpha=0.9))

    plt.xlabel(f'PC{pc1+1} ({variance_explained[pc1]*100:.1f}% variance)')
    plt.ylabel(f'PC{pc2+1} ({variance_explained[pc2]*100:.1f}% variance)')
    plt.title(f'Biplot with {group_labels}: PC{pc2+1} vs PC{pc1+1}')
    plt.grid(True, alpha=0.3)
    plt.axhline(0, color='black', linewidth=0.5)
    plt.axvline(0, color='black', linewidth=0.5)
    plt.legend()

    return plt.gcf()

# Get gender data for the same indices as our PC scores
gender_data = df_full.loc[pc_scores_df.index, 'Gender'].values

# Create enhanced biplot
enhanced_biplot = create_biplot_with_groups(pc_scores, loadings.T,
    df_tests.columns, gender_data, 'Gender')
plt.tight_layout()
```

```
46 plt.show()
```

Listing 22: Enhanced Biplot

# 10 Practical Applications and Case Studies

## 10.1 AFQT Analysis

```
1  # Extract AFQT-related tests
2  afqt_tests = ['Word', 'Parag', 'Math', 'Arith']
3  print(f"AFQT component tests: {afqt_tests}")
4
5  # Perform PCA on just these 4 tests
6  df_afqt = df_full[afqt_tests + ['AFQT', 'Gender']].copy()
7  df_afqt_tests = df_afqt[afqt_tests]
8
9  print(f"\nAFQT dataset shape: {df_afqt_tests.shape}")
10 print("AFQT test correlations:")
11 print(df_afqt_tests.corr().round(2))
12
13 # Standardize and perform PCA
14 scaler_afqt = StandardScaler()
15 df_afqt_scaled = scaler_afqt.fit_transform(df_afqt_tests)
16
17 pca_afqt = PCA()
18 pca_afqt.fit(df_afqt_scaled)
19 afqt_pc_scores = pca_afqt.transform(df_afqt_scaled)
20
21 print(f"\nAFQT PCA Results:")
22 print(f"Variance explained: {pca_afqt.explained_variance_ratio_}")
23 print(f"Cumulative variance: {np.cumsum(pca_afqt.explained_variance_ratio_
       )}")
24
25 # Compare PC1 with official AFQT score
26 correlation_afqt_pc1 = np.corrcoef(df_afqt['AFQT'], afqt_pc_scores[:, 0])
       [0, 1]
27 print(f"\nCorrelation between AFQT and PC1: {correlation_afqt_pc1:.3f}")
28
29 # Also compare with simple sum
30 afqt_sum = df_afqt_tests.sum(axis=1)
31 correlation_afqt_sum = np.corrcoef(df_afqt['AFQT'], afqt_sum)[0, 1]
32 print(f"Correlation between AFQT and sum of 4 tests: {correlation_afqt_sum
       :.3f}")
33
34 # Visualize the relationship
35 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))
36
37 # AFQT vs PC1
38 ax1.scatter(afqt_pc_scores[:, 0], df_afqt['AFQT'], alpha=0.5)
39 ax1.set_xlabel('PC1 Score (AFQT tests)')
40 ax1.set_ylabel('Official AFQT Score')
```

```
41 ax1.set_title(f'AFQT vs PC1 (r = {correlation_afqt_pc1:.3f})')
42 ax1.grid(True, alpha=0.3)
43
44 # Add trend line
45 z = np.polyfit(afqt_pc_scores[:, 0], df_afqt['AFQT'], 1)
46 p = np.poly1d(z)
47 ax1.plot(afqt_pc_scores[:, 0], p(afqt_pc_scores[:, 0]), "r--", alpha=0.8)
48
49 # AFQT vs sum
50 ax2.scatter(afqt_sum, df_afqt['AFQT'], alpha=0.5, color='green')
51 ax2.set_xlabel('Sum of 4 Test Scores')
52 ax2.set_ylabel('Official AFQT Score')
53 ax2.set_title(f'AFQT vs Simple Sum (r = {correlation_afqt_sum:.3f})')
54 ax2.grid(True, alpha=0.3)
55
56 # Add trend line
57 z2 = np.polyfit(afqt_sum, df_afqt['AFQT'], 1)
58 p2 = np.poly1d(z2)
59 ax2.plot(afqt_sum, p2(afqt_sum), "r--", alpha=0.8)
60
61 plt.tight_layout()
62 plt.show()
```

Listing 23: AFQT Analysis


## 10.2   Gender Differences Analysis

```
1 print("\nGENDER DIFFERENCES ANALYSIS:")
2 print("="*50)
3
4 # Create detailed gender comparison for all 10 tests
5 gender_comparison = df_full.groupby('Gender')[test_columns].agg(['mean', '
    std']).round(2)
6 print("Test score means and standard deviations by gender:")
7 print(gender_comparison)
8
9 # Statistical significance tests
10 print("\nStatistical significance tests (t-tests):")
11 print("Test" + " "*15 + "Male Mean" + " "*5 + "Female Mean" + " "*3 + "p-
    value" + " "*5 + "Significant?")
12 print("-" * 70)
13
14 for test in test_columns:
15     male_scores = df_full[df_full['Gender'] == 'male'][test]
16     female_scores = df_full[df_full['Gender'] == 'female'][test]
17     t_stat, p_value = stats.ttest_ind(male_scores, female_scores)
18     significant = "Yes" if p_value < 0.05 else "No"
19     print(f"{test:<18} {male_scores.mean():8.2f} {female_scores.mean()
    :13.2f} {p_value:10.4f} {significant:>8}")
20
21 # Gender differences in PC scores
22 print(f"\nGender differences in Principal Components:")
```

```
23 pc_gender_df = pd.DataFrame(pc_scores[:, :4], columns=['PC1', 'PC2', 'PC3'
     , 'PC4'])
24 pc_gender_df['Gender'] = df_full.loc[pc_scores_df.index, 'Gender'].values
25
26 gender_pc_stats = pc_gender_df.groupby('Gender')[['PC1', 'PC2', 'PC3', '
     PC4']].agg(['mean', 'std']).round(3)
27 print(gender_pc_stats)
28
29 # Test significance of gender differences in PCs
30 print(f"\nStatistical tests for PC gender differences:")
31 for pc in ['PC1', 'PC2', 'PC3', 'PC4']:
32     male_pc = pc_gender_df[pc_gender_df['Gender'] == 'male'][pc]
33     female_pc = pc_gender_df[pc_gender_df['Gender'] == 'female'][pc]
34     t_stat, p_value = stats.ttest_ind(male_pc, female_pc)
35     effect_size = (male_pc.mean() - female_pc.mean()) / np.sqrt((male_pc.
     var() + female_pc.var()) / 2)
36     print(f"{pc}: t = {t_stat:.3f}, p = {p_value:.4f}, Cohen's d = {
     effect_size:.3f}")
```

Listing 24: Gender Differences Analysis

# 11 Advanced Visualizations and Interpretations

## 11.1 3D Visualization of First Three Components

```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 fig = plt.figure(figsize=(12, 9))
4 ax = fig.add_subplot(111, projection='3d')
5
6 # Color by gender
7 colors = {'male': 'blue', 'female': 'red'}
8
9 for gender in colors:
10     mask = pc_scores_with_gender['Gender'] == gender
11     ax.scatter(pc_scores_with_gender.loc[mask, 'PC1'],
12                pc_scores_with_gender.loc[mask, 'PC2'],
13                pc_scores_with_gender.loc[mask, 'PC3'],
14                c=colors[gender], alpha=0.6, s=30, label=gender)
15
16 ax.set_xlabel(f'PC1 ({variance_explained[0]*100:.1f}%)')
17 ax.set_ylabel(f'PC2 ({variance_explained[1]*100:.1f}%)')
18 ax.set_zlabel(f'PC3 ({variance_explained[2]*100:.1f}%)')
19 ax.set_title('3D View: First Three Principal Components')
20 ax.legend()
21
22 plt.tight_layout()
23 plt.show()
```

Listing 25: 3D Visualization

## 11.2 Contribution of Each Variable to Each Component

```python
# Calculate the contribution of each variable to each PC
# Contribution = (loading^2) * (variance explained by PC)
contributions = np.zeros((len(test_columns), 4))  # 10 tests   4 PCs

for i in range(4):  # First 4 PCs
    for j, test in enumerate(test_columns):
        contributions[j, i] = (loadings_df.iloc[j, i]**2) *
    variance_explained[i]

# Create DataFrame
contrib_df = pd.DataFrame(contributions,
                          index=test_columns,
                          columns=['PC1', 'PC2', 'PC3', 'PC4'])

# Normalize to percentages
contrib_pct = contrib_df.div(contrib_df.sum(axis=0), axis=1) * 100

print("CONTRIBUTION OF EACH TEST TO EACH COMPONENT (%):")
print(contrib_pct.round(1))

# Visualize contributions
fig, axes = plt.subplots(2, 2, figsize=(15, 12))
axes = axes.ravel()

for i, pc in enumerate(['PC1', 'PC2', 'PC3', 'PC4']):
    contrib_sorted = contrib_pct[pc].sort_values(ascending=True)
    axes[i].barh(range(len(contrib_sorted)), contrib_sorted.values)
    axes[i].set_yticks(range(len(contrib_sorted)))
    axes[i].set_yticklabels(contrib_sorted.index)
    axes[i].set_xlabel('Contribution (%)')
    axes[i].set_title(f'{pc} - Variable Contributions\n({
    variance_explained[i]*100:.1f}% of total variance)')
    axes[i].grid(True, alpha=0.3)

plt.tight_layout()
plt.show()
```

Listing 26: Variable Contribution Analysis

## 11.3 Correlation Between Original Variables and PCs

```python
# Calculate correlations between original variables and PCs
pc_var_correlations = np.zeros((len(test_columns), 4))

for i, test in enumerate(test_columns):
    for j in range(4):
        pc_var_correlations[i, j] = np.corrcoef(df_tests_scaled.iloc[:, i
    ], pc_scores[:, j])[0, 1]

# Create DataFrame
```

```
9  corr_df = pd.DataFrame(pc_var_correlations,
10                         index=test_columns,
11                         columns=['PC1', 'PC2', 'PC3', 'PC4'])
12
13 print("CORRELATIONS BETWEEN ORIGINAL TESTS AND PRINCIPAL COMPONENTS:")
14 print(corr_df.round(3))
15
16 # Visualize as heatmap
17 plt.figure(figsize=(10, 8))
18 sns.heatmap(corr_df,
19             annot=True,
20             cmap='RdBu_r',
21             center=0,
22             fmt='.2f',
23             xticklabels=[f'{pc}\n({var*100:.1f}%)' for pc, var in zip(
    corr_df.columns, variance_explained[:4])],
24             yticklabels=test_columns)
25 plt.title('Correlations: Original Tests vs Principal Components')
26 plt.tight_layout()
27 plt.show()
```

Listing 27: Correlation Analysis

# 12 Model Evaluation and Validation

## 12.1 Reconstruction Error Analysis

```
1  def reconstruction_error(original_data, n_components):
2      """
3      Calculate how well n components can reconstruct the original data
4      """
5      pca_temp = PCA(n_components=n_components)
6      pca_temp.fit(original_data)
7
8      # Transform to PC space and back
9      pc_scores_temp = pca_temp.transform(original_data)
10     reconstructed = pca_temp.inverse_transform(pc_scores_temp)
11
12     # Calculate mean squared error
13     mse = np.mean((original_data - reconstructed) ** 2)
14
15     return mse, pca_temp.explained_variance_ratio_.sum()
16
17 print("RECONSTRUCTION ANALYSIS:")
18 print("Components" + " "*5 + "MSE" + " "*10 + "Variance Kept" + " "*5 + "
    Information Lost")
19 print("-" * 60)
20
21 for n_comp in range(1, 6):
22     mse, var_kept = reconstruction_error(df_tests_scaled, n_comp)
23     info_lost = 1 - var_kept
```

```
24    print(f"{n_comp:6d} {mse:15.4f} {var_kept:12.1%} {info_lost:15.1%}")
```

## 12.2   Stability Analysis

```python
1  def pca_stability_test(data, n_bootstrap=10, sample_size=1000):
2      """
3      Test PCA stability using bootstrap sampling
4      """
5      n_features = data.shape[1]
6      loadings_collection = []
7
8      for i in range(n_bootstrap):
9          # Random sample
10         sample_idx = np.random.choice(data.shape[0], size=sample_size,
   replace=True)
11         sample_data = data[sample_idx]
12
13         # PCA
14         pca_boot = PCA()
15         pca_boot.fit(sample_data)
16
17         # Store first two PC loadings
18         loadings_collection.append(pca_boot.components_[:2])
19
20     return loadings_collection
21
22 # Run stability test
23 np.random.seed(42)
24 stability_results = pca_stability_test(df_tests_scaled.values, n_bootstrap
   =20)
25
26 # Calculate loading stability (standard deviation across bootstrap samples
   )
27 pc1_loadings_std = np.std([result[0] for result in stability_results],
   axis=0)
28 pc2_loadings_std = np.std([result[1] for result in stability_results],
   axis=0)
29
30 print(f"\nSTABILITY ANALYSIS:")
31 print("Standard deviation of loadings across 20 bootstrap samples:")
32
33 print("\nPC1 loading stability (lower = more stable):")
34 for i, test in enumerate(test_columns):
35     print(f"{test:12}: {pc1_loadings_std[i]:.4f}")
36
37 print(f"\nPC2 loading stability:")
38 for i, test in enumerate(test_columns):
39     print(f"{test:12}: {pc2_loadings_std[i]:.4f}")
```

Listing 29: Stability Analysis

# 13 Summary Dashboard and Final Visualizations

```python
fig = plt.figure(figsize=(20, 15))

# 1. Scree plot
ax1 = plt.subplot(3, 4, 1)
plt.plot(range(1, len(variance_explained) + 1), variance_explained, 'bo-')
plt.xlabel('Component')
plt.ylabel('Variance Explained')
plt.title('Scree Plot')
plt.grid(True, alpha=0.3)

# 2. Cumulative variance
ax2 = plt.subplot(3, 4, 2)
plt.plot(range(1, len(cumulative_variance) + 1), cumulative_variance, 'ro-
    ')
plt.axhline(y=0.8, color='green', linestyle='--', label='80%')
plt.xlabel('Components')
plt.ylabel('Cumulative Variance')
plt.title('Cumulative Variance')
plt.legend()
plt.grid(True, alpha=0.3)

# 3. PC1 vs PC2 scores
ax3 = plt.subplot(3, 4, (3, 4))
for gender in ['male', 'female']:
    mask = pc_scores_with_gender['Gender'] == gender
    color = 'blue' if gender == 'male' else 'red'
    plt.scatter(pc_scores_with_gender.loc[mask, 'PC1'],
                pc_scores_with_gender.loc[mask, 'PC2'],
                c=color, alpha=0.6, s=30, label=gender)

plt.xlabel(f'PC1 ({variance_explained[0]*100:.1f}%)')
plt.ylabel(f'PC2 ({variance_explained[1]*100:.1f}%)')
plt.title('PC Scores by Gender')
plt.legend()
plt.grid(True, alpha=0.3)

# 4. PC1 loadings
ax4 = plt.subplot(3, 4, (5, 6))
pc1_sorted = loadings_df['PC1'].sort_values()
plt.barh(range(len(pc1_sorted)), pc1_sorted.values)
plt.yticks(range(len(pc1_sorted)), pc1_sorted.index)
plt.xlabel('Loading')
plt.title('PC1 Loadings (General Ability)')
plt.grid(True, alpha=0.3)

# 5. PC2 loadings
ax5 = plt.subplot(3, 4, (7, 8))
pc2_sorted = loadings_df['PC2'].sort_values()
plt.barh(range(len(pc2_sorted)), pc2_sorted.values)
plt.yticks(range(len(pc2_sorted)), pc2_sorted.index)
plt.xlabel('Loading')
```

```
51  plt.title('PC2 Loadings (Ability Contrast)')
52  plt.grid(True, alpha=0.3)
53
54  # 6. Correlation heatmap (subset)
55  ax6 = plt.subplot(3, 4, (9, 10))
56  sns.heatmap(correlation_matrix.iloc[:5, :5], annot=True, fmt='.2f', cmap='
        coolwarm', center=0)
57  plt.title('Correlation Matrix (subset)')
58
59  # 7. AFQT correlation
60  ax7 = plt.subplot(3, 4, 11)
61  plt.scatter(afqt_pc_scores[:, 0], df_afqt['AFQT'], alpha=0.5, s=10)
62  plt.xlabel('PC1 Score')
63  plt.ylabel('AFQT Score')
64  plt.title(f'AFQT vs PC1 (r={correlation_afqt_pc1:.2f})')
65  plt.grid(True, alpha=0.3)
66
67  # 8. Gender differences summary
68  ax8 = plt.subplot(3, 4, 12)
69  gender_means = pc_gender_df.groupby('Gender')[['PC1', 'PC2']].mean()
70  x_pos = np.arange(len(gender_means.columns))
71  width = 0.35
72
73  plt.bar(x_pos - width/2, gender_means.loc['male'], width, label='Male',
        alpha=0.8)
74  plt.bar(x_pos + width/2, gender_means.loc['female'], width, label='Female'
        , alpha=0.8)
75
76  plt.xlabel('Principal Component')
77  plt.ylabel('Mean Score')
78  plt.title('Gender Differences in PC Scores')
79  plt.xticks(x_pos, ['PC1', 'PC2'])
80  plt.legend()
81  plt.grid(True, alpha=0.3)
82
83  plt.tight_layout()
84  plt.show()
```

Listing 30: Summary Dashboard

# 14  Professional Summary and Best Practices

```
1  print("\n" + "="*70)
2  print("PRINCIPAL COMPONENT ANALYSIS - EXECUTIVE SUMMARY")
3  print("="*70)
4
5  print(f"\nKEY FINDINGS:")
6  print(f"1. PC1 explains {variance_explained[0]*100:.1f}% of variance -
        represents GENERAL COGNITIVE ABILITY")
7  print(f"2. PC2 explains {variance_explained[1]*100:.1f}% of variance -
        represents SPECIALIZED ABILITY CONTRAST")
```

```
 8 print(f"3. First 2 PCs together explain {cumulative_variance[1]*100:.1f}%
     of total variance")
 9 print(f"4. Strong correlation between PC1 and AFQT score (r = {
     correlation_afqt_pc1:.3f})")
10
11 print(f"\nCOMPONENT INTERPRETATIONS:")
12 # PC1 interpretation
13 pc1_avg_loading = np.abs(loadings_df['PC1']).mean()
14 if np.abs(loadings_df['PC1']).std() < 0.1:
15     print(f"PC1: General cognitive ability (all loadings approximately {
     pc1_avg_loading:.2f})")
16 else:
17     print(f"PC1: General cognitive ability with emphasis on:")
18     pc1_high = loadings_df['PC1'].abs().nlargest(3)
19     for test, loading in pc1_high.items():
20         print(f" - {test} (loading: {loading:.3f})")
21
22 # PC2 interpretation
23 pc2_positive = loadings_df[loadings_df['PC2'] > 0.2]['PC2'].sort_values(
     ascending=False)
24 pc2_negative = loadings_df[loadings_df['PC2'] < -0.2]['PC2'].sort_values()
25
26 if len(pc2_positive) > 0 and len(pc2_negative) > 0:
27     print(f"PC2: Contrast between ability types")
28     print(f" Higher loadings: {', '.join(pc2_positive.index)}")
29     print(f" Lower loadings: {', '.join(pc2_negative.index)}")
30 else:
31     print(f"PC2: Specialized ability factor")
32
33 print(f"\nGENDER ANALYSIS:")
34 if p_pc2 < 0.05:
35     print(f"Significant gender differences found in PC2 (p = {p_pc2:.4f})"
     )
36     pc2_male_mean = pc2_male.mean()
37     pc2_female_mean = pc2_female.mean()
38     if pc2_male_mean > pc2_female_mean:
39         print("Males score higher on PC2 on average")
40     else:
41         print("Females score higher on PC2 on average")
42 else:
43     print(f"No significant gender differences in PC2 (p = {p_pc2:.4f})")
44
45 print(f"\nPRACTICAL APPLICATIONS:")
46 print("1. RECRUITMENT: Use PC1 for general ability screening")
47 print("2. JOB MATCHING: Use PC2 to match candidates to specialized roles")
48 print("3. TRAINING PROGRAMS: Focus resources on areas where PC scores are
     low")
49 print("4. TEAM COMPOSITION: Consider PC2 patterns for balanced skill sets"
     )
50 print("5. PERFORMANCE PREDICTION: PC1 strongly correlates with overall
     performance measures")
51
52 print(f"\nMETHODOLOGICAL CONSIDERATIONS:")
53 print("1. STANDARDIZATION: Data was standardized due to different test
```

```
     scales")
54 print("2. COMPONENT SELECTION: First 2 components capture 75%+ of variance
       ")
55 print("3. STABILITY: Bootstrap analysis confirms robust loading patterns")
56 print("4. INTERPRETATION: PC1 = general ability , PC2 = specialized
       contrasts")
57
58 print(f"\nLIMITATIONS AND CAUTIONS:")
59 print("1. PCA assumes linear relationships between variables")
60 print("2. Components may not always have clear substantive interpretations
       ")
61 print("3. Scaling decisions can affect results significantly")
62 print("4. Cultural and demographic factors may influence ability patterns"
       )
63 print("5. Statistical significance does not imply practical significance")
64
65 print(f"\nRECOMMENDATIONS FOR PRACTITIONERS:")
66 print("1. Always examine scree plots before deciding on number of
       components")
67 print("2. Validate component interpretations with domain expertise")
68 print("3. Consider both scaled and unscaled versions for comparison")
69 print("4. Use bootstrap sampling to assess stability of results")
70 print("5. Supplement PCA with other analytical techniques for robust
       conclusions")
```

Listing 31: Executive Summary

# 15 Advanced Extensions and Future Directions

```
1 print(f"\n" + "="*70)
2 print("ADVANCED TOPICS AND EXTENSIONS")
3 print("="*70)
4
5 print(f"\nRELATED TECHNIQUES:")
6 print("1. FACTOR ANALYSIS: Similar to PCA but models latent constructs")
7 print("2. INDEPENDENT COMPONENT ANALYSIS (ICA): Finds statistically
       independent sources")
8 print("3. SPARSE PCA: Forces some loadings to zero for easier
       interpretation")
9 print("4. ROBUST PCA: Handles outliers and missing data better")
10 print("5. KERNEL PCA: Captures non-linear relationships in the data")
11
12 print(f"\nCOMPUTATIONAL CONSIDERATIONS:")
13 print("1. EIGENDECOMPOSITION: PCA via correlation/covariance matrix
       eigenvalues")
14 print("2. SINGULAR VALUE DECOMPOSITION (SVD): Alternative computational
       approach")
15 print("3. INCREMENTAL PCA: For datasets too large to fit in memory")
16 print("4. RANDOMIZED PCA: Faster approximation for high-dimensional data")
17
18 print(f"\nVALIDATION STRATEGIES:")
19 print("1. CROSS-VALIDATION: Split data to test component stability")
```

```
20 print ("2. BOOTSTRAP RESAMPLING: Assess uncertainty in loadings")
21 print ("3. EXTERNAL VALIDATION: Compare components with known criteria")
22 print ("4. REPLICATION: Test findings on independent datasets")
```
Listing 32: Advanced Topics

# 16 Code Summary and Reproducibility

```
1 print (f"\n" + "="*70)
2 print ("REPRODUCIBILITY AND CODE SUMMARY")
3 print ("="*70)
4
5 print (f"\nKEY PYTHON LIBRARIES USED:")
6 print ("- pandas: Data manipulation and analysis")
7 print ("- numpy: Numerical computing")
8 print ("- matplotlib: Basic plotting")
9 print ("- seaborn: Statistical visualization")
10 print ("- scikit-learn: Machine learning (PCA implementation)")
11 print ("- scipy: Statistical tests")
12
13 print (f"\nCORE PCA WORKFLOW:")
14 print ("1. Load and explore data")
15 print ("2. Standardize variables (if needed)")
16 print ("3. Fit PCA model: pca = PCA().fit(data)")
17 print ("4. Extract loadings: pca.components_")
18 print ("5. Get PC scores: pca.transform(data)")
19 print ("6. Analyze variance explained: pca.explained_variance_ratio_")
20 print ("7. Create visualizations and interpret results")
21
22 print (f"\nIMPORTANT PARAMETERS:")
23 print ("- n_components: Number of components to keep")
24 print ("- StandardScaler: Center and scale data")
25 print ("- explained_variance_ratio_: Proportion of variance per component")
26 print ("- components_: Loading matrix (variables    components)")
27
28 print (f"\nSUGGESTED WORKFLOW FOR NEW DATASETS:")
29 print ("1. Examine data structure and correlations")
30 print ("2. Decide on standardization based on variable scales")
31 print ("3. Perform PCA with all components initially")
32 print ("4. Create scree plot to guide component selection")
33 print ("5. Interpret loadings with domain knowledge")
34 print ("6. Visualize PC scores to identify patterns")
35 print ("7. Validate findings with external criteria")
```
Listing 33: Reproducibility

# 17 Final Exercises and Assessment

```
1 print (f"\n" + "="*70)
```

```python
print("SELF-ASSESSMENT QUESTIONS")
print("="*70)

print(f"\nCONCEPTUAL QUESTIONS:")
print("1. What does it mean when PC1 explains 60% of the variance?")
print("2. Why might we standardize data before performing PCA?")
print("3. How do you interpret negative loadings in a principal component?
    ")
print("4. What is the relationship between eigenvalues and explained
    variance?")
print("5. When would you choose 2 components vs 5 components?")

print(f"\nPRACTICAL QUESTIONS:")
print("1. How would you modify the code to perform PCA without
    standardization?")
print("2. What code would you use to keep only the first 3 components?")
print("3. How would you calculate the total variance explained by PC1 and
    PC2?")
print("4. What visualization would best show the relationship between
    components?")
print("5. How would you test if gender differences are statistically
    significant?")

print(f"\nINTERPRETATION QUESTIONS:")
print("1. If all PC1 loadings are positive and similar, what does this
    suggest?")
print("2. What does it mean if two variables have similar loading patterns
    ?")
print("3. How would you explain PC2 to a non-technical audience?")
print("4. What are the practical implications of the gender differences
    found?")
print("5. How confident should we be in these results?")

print(f"\n" + "="*70)
print("COURSE OBJECTIVES ACHIEVED")
print("="*70)
print("UPON COMPLETION OF THIS LECTURE, STUDENTS CAN:")
print("    Understand PCA concepts: dimension reduction, variance
    explained, loadings")
print("    Implement PCA in Python using scikit-learn")
print("    Create and interpret key visualizations: scree plots, biplots,
    PC score plots")
print("    Analyze real-world datasets and extract meaningful insights")
print("    Understand when and how to standardize data")
print("    Interpret statistical significance of group differences")
print("    Apply best practices for PCA analysis and validation")
print("    Recognize limitations and potential extensions of PCA")

print(f"\n" + "="*70)
print("END OF PCA LECTURE")
print("Save your work and prepare for the next module!")
print("="*70)
```

Listing 34: Self-Assessment

28