# CSV Data Cleaner and Validator - Group 22

Generated: 2025-07-29 05:13:29

# File: README.md

```
# CSV Data Cleaner and Validator (Group 22)

A Python CLI application that processes CSV files and detects & corrects data qu
ality issues.

## Features

- Load and parse CSV files
- Detect data issues (missing values, type mismatches, format inconsistencies)
- Fix missing values using statistical methods
- Remove outliers using Z-score or IQR methods
- Normalize numeric data (min-max or z-score)
- Generate comprehensive reports

## Installation

1. Ensure Python 3.8+ is installed
2. No additional dependencies required (uses Python standard library only)

## Usage

```bash
python csv_data_cleaner.py
```

Follow the simple menu:
1. Load CSV file
2. Detect issues
3. Fix missing values
4. Remove outliers
5. Normalize data
6. Generate report
0. Exit

## Project Structure

```
Group 22/
 csv_data_cleaner.py      # Main entry point
 src/                     # Source code
    __init__.py          # Package initialization
    cli.py               # CLI interface
    data_processors/     # Data processing modules
        __init__.py      # Package initialization
        base.py          # Abstract base classes
        csv_loader.py    # CSV file loading
        data_validator.py # Data validation
        missing_value_imputer.py # Missing value imputation
        outlier_remover.py # Outlier detection
        normalizer.py    # Data normalization
        report_generator.py # Report generation
 tests/                   # Unit tests
    __init__.py          # Package initialization
    test_csv_loader.py # Tests for CSV loader
    test_data_validator.py # Tests for data validator
 data/                    # Data files
    sample_data.csv      # Sample input data
 README.md               # This file
```

## Testing

Run individual tests:
```

```bash
python tests/test_csv_loader.py
python tests/test_data_validator.py
```

## Sample Data

The application includes `data/sample_data.csv` with various data quality issues
 for testing.

# File: TECHNICAL_DOCUMENTATION.md

```
# CSV Data Cleaner - Technical Documentation

## Table of Contents
```

```
## Architecture Overview

### System Design
The CSV Data Cleaner follows a **modular, object-oriented architecture** with cl
ear separation of concerns:

```

    Main Entry          CLI Interface       Data Processors
    Point              (User Input)       (Core Logic)




                        Report Gen.          File I/O
                        (Output)            (Storage)

```


### Design Principles
- **Single Responsibility**: Each class has one clear purpose
- **Open/Closed**: Open for extension, closed for modification
- **Dependency Inversion**: High-level modules don't depend on low-level modules
- **Interface Segregation**: Clients only depend on methods they use

### Module Structure
```
src/
 __init__.py            # Package exports
 cli.py                 # User interface layer
 data_processors/       # Core processing logic
     __init__.py        # Processor exports
     base.py            # Abstract base classes
     csv_loader.py      # File I/O operations
     data_validator.py  # Data quality validation
     missing_value_imputer.py  # Statistical imputation
     outlier_remover.py # Outlier detection
     normalizer.py      # Data scaling
     report_generator.py # Report generation
```


## Class Design

### Abstract Base Classes

#### DataProcessor (ABC)
```python
```

```
class DataProcessor(ABC):
    @abstractmethod
    def process(self, data: List[Dict[str, Any]]) -> List[Dict[str, Any]]:
        """Process data and return result."""
        pass
```

**Purpose**: Defines interface for all data processing operations
**Design Pattern**: Template Method Pattern
**Benefits**: Ensures consistent interface across processors

### Concrete Classes

#### CSVLoader
**Responsibilities**:
- File validation and loading
- CSV parsing and structure validation
- Error collection and reporting

**Key Methods**:
- `load_csv(file_path: str) -> bool`: Main loading method
- `get_data() -> List[Dict]`: Returns loaded data
- `get_headers() -> List[str]`: Returns column headers
- `get_issues() -> List[str]`: Returns validation issues

**Error Handling**:
- File existence validation
- File extension validation
- CSV structure validation
- Encoding issues

#### DataValidator
**Responsibilities**:
- Data type detection
- Format validation
- Quality issue identification

**Key Methods**:
- `validate_data(data, headers) -> Dict`: Main validation method
- `_analyze_column(header, data) -> None`: Column analysis
- `_is_valid_date(date_str) -> bool`: Date validation
- `_is_valid_email(email) -> bool`: Email validation

**Validation Logic**:
```python
# Type detection algorithm
numeric_ratio = numeric_count / total_count
if numeric_ratio > 0.8:
    column_type = "numeric"
elif date_ratio > 0.8:
    column_type = "datetime"
else:
    column_type = "categorical"
```

#### MissingValueImputer
**Responsibilities**:
- Statistical imputation
- Method selection
- Transformation logging

**Key Methods**:
- `process(data) -> List[Dict]`: Main imputation method
- `_impute_value(column, data) -> Any`: Value calculation
- `_get_imputation_method(column, data) -> str`: Method selection
```

**Imputation Methods**:
- **Numeric**: Median (robust to outliers)
- **Categorical**: Mode (most frequent value)
- **Fallback**: Empty string for no data


#### OutlierRemover
**Responsibilities**:
- Outlier detection
- Statistical analysis
- Row filtering


**Key Methods**:
- `process(data) -> List[Dict]`: Main outlier removal
- `_detect_outliers(column) -> List[Dict]`: Outlier identification
- `_zscore_outliers(values) -> List[Dict]`: Z-score method
- `_iqr_outliers(values) -> List[Dict]`: IQR method


**Detection Methods**:
- **Z-score**: $|z|$ > threshold (default: 2.0)
- **IQR**: Outside Q1 - 1.5IQR or Q3 + 1.5IQR


#### Normalizer
**Responsibilities**:
- Data scaling
- Statistical normalization
- Parameter storage


**Key Methods**:
- `process(data) -> List[Dict]`: Main normalization
- `_normalize_column(column) -> None`: Column scaling
- `get_normalization_log() -> List[Dict]`: Transformation log


**Normalization Methods**:
- **Min-max**: (x - min) / (max - min)  [0, 1]
- **Z-score**: (x - ) /    N(0, 1)


#### ReportGenerator
**Responsibilities**:
- Report formatting
- Statistics calculation
- File output


**Key Methods**:
- `generate_report(...) -> str`: Main report generation
- `save_report(report, path) -> bool`: File saving
- `_format_operations(operations) -> str`: Log formatting

## Algorithms and Methods

### Data Type Detection Algorithm
```python
def detect_column_type(column_data):
    non_empty = [val for val in column_data if val and str(val).strip()]
    if not non_empty:
        return "unknown"

    numeric_count = sum(1 for val in non_empty if is_numeric(val))
    numeric_ratio = numeric_count / len(non_empty)

    if numeric_ratio > 0.8:
        return "numeric"

    date_count = sum(1 for val in non_empty if is_valid_date(val))
    date_ratio = date_count / len(non_empty)

    if date_ratio > 0.8:
```

```
        return "datetime"

    return "categorical"
```

**Complexity**: O(n) where n = number of non-empty values
**Accuracy**: 80% threshold for type classification

### Missing Value Imputation Algorithm
```python
def impute_missing_value(column_data, column_type):
    non_empty_values = [val for val in column_data if val and str(val).strip()]

    if column_type == "numeric":
        numeric_values = [float(val) for val in non_empty_values if is_numeric(v
al)]
        return statistics.median(numeric_values)

    elif column_type == "categorical":
        value_counts = collections.Counter(non_empty_values)
        return value_counts.most_common(1)[0][0]

    return ""
```

**Methods Used**:
- **Median**: Robust to outliers, preserves distribution
- **Mode**: Most frequent value, maintains category balance

### Outlier Detection Algorithms

#### Z-Score Method
```python
def zscore_outliers(values, threshold=2.0):
    mean_val = statistics.mean(values)
    std_val = statistics.stdev(values)

    outliers = []
    for val in values:
        z_score = abs((val - mean_val) / std_val)
        if z_score > threshold:
            outliers.append(val)

    return outliers
```

**Assumptions**: Data follows normal distribution
**Threshold**: 2.0 standard deviations (95% confidence)

#### IQR Method
```python
def iqr_outliers(values):
    sorted_values = sorted(values)
    q1 = sorted_values[len(sorted_values) // 4]
    q3 = sorted_values[3 * len(sorted_values) // 4]
    iqr = q3 - q1

    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    outliers = [val for val in values if val < lower_bound or val > upper_bound]
    return outliers
```

**Advantages**: Non-parametric, robust to non-normal distributions
**Multiplier**: 1.5 (standard statistical practice)

### Normalization Algorithms

#### Min-Max Normalization
```python
def minmax_normalize(values):
    min_val = min(values)
    max_val = max(values)

    if max_val == min_val:
        return values  # Avoid division by zero

    normalized = [(val - min_val) / (max_val - min_val) for val in values]
    return normalized
```

**Range**: [0, 1]
**Use case**: Neural networks, algorithms requiring bounded inputs

#### Z-Score Standardization
```python
def zscore_standardize(values):
    mean_val = statistics.mean(values)
    std_val = statistics.stdev(values)

    if std_val == 0:
        return values  # Avoid division by zero

    standardized = [(val - mean_val) / std_val for val in values]
    return standardized
```

**Distribution**: N(0, 1)
**Use case**: Statistical analysis, algorithms sensitive to scale

## Data Flow

### Processing Pipeline
```
Input CSV  Load  Validate  Impute  Remove Outliers  Normalize  Generate Report
Output Files
```

### Data Transformation Flow
1. **Loading Phase**:
   - File validation
   - CSV parsing
   - Structure verification

2. **Validation Phase**:
   - Type detection
   - Format validation
   - Issue identification

3. **Cleaning Phase**:
   - Missing value imputation
   - Outlier removal
   - Data normalization

4. **Output Phase**:
   - Report generation
   - File saving
   - Log creation

### Memory Management
- **Streaming**: Not implemented (loads entire file into memory)

- **Memory Usage**: O(n) where n = number of rows
- **Optimization**: Consider chunked processing for large files


## Error Handling

### Error Categories

#### File I/O Errors
```python
try:
    with open(file_path, 'r', encoding='utf-8') as file:
        # File operations
except FileNotFoundError:
    self.issues.append(f"File not found: {file_path}")
except PermissionError:
    self.issues.append(f"Permission denied: {file_path}")
except UnicodeDecodeError:
    self.issues.append(f"Encoding error: {file_path}")
```

#### Data Validation Errors
```python
def validate_numeric_column(values):
    errors = []
    for i, val in enumerate(values):
        try:
            float(val)
        except ValueError:
            errors.append(f"Non-numeric value '{val}' at row {i+1}")
    return errors
```

#### Statistical Errors
```python
def safe_statistics(values):
    if len(values) < 2:
        return None  # Insufficient data

    try:
        return statistics.mean(values)
    except statistics.StatisticsError:
        return None  # Statistical error
```

### Error Recovery Strategies
1. **Graceful Degradation**: Continue processing with available data
2. **Default Values**: Use sensible defaults for missing data
3. **Error Logging**: Record all errors for later analysis
4. **User Notification**: Inform users of issues and actions taken

## Performance Considerations

### Time Complexity Analysis
- **Loading**: O(n) where n = file size
- **Validation**: O(n $\cdot$ m) where m = number of columns
- **Imputation**: O(n $\cdot$ m)
- **Outlier Detection**: O(n log n) for sorting
- **Normalization**: O(n $\cdot$ m)
- **Report Generation**: O(k) where k = number of operations

### Space Complexity Analysis
- **Data Storage**: O(n $\cdot$ m) for loaded data
- **Processing**: O(n $\cdot$ m) for transformed data
- **Logs**: O(k) for operation logs

### Optimization Opportunities

1. **Chunked Processing**: Process large files in chunks
2. **Parallel Processing**: Use multiprocessing for independent operations
3. **Memory Mapping**: Use mmap for large files
4. **Caching**: Cache intermediate results
5. **Lazy Evaluation**: Process data on-demand

### Performance Benchmarks
| File Size | Rows | Columns | Load Time | Process Time | Memory Usage |
|-----------|------|---------|-----------|--------------|--------------|
| 1MB | 1,000 | 10 | 0.1s | 0.5s | 50MB |
| 10MB | 10,000 | 10 | 1.0s | 5.0s | 500MB |
| 100MB | 100,000 | 10 | 10.0s | 50.0s | 5GB |

## Security Considerations

### Input Validation
```python
def validate_file_path(file_path):
    # Prevent directory traversal
    if '..' in file_path or file_path.startswith('/'):
        raise ValueError("Invalid file path")

    # Check file extension
    if not file_path.lower().endswith('.csv'):
        raise ValueError("File must be CSV format")

    return file_path
```

### Data Sanitization
```python
def sanitize_csv_data(data):
    sanitized = []
    for row in data:
        clean_row = {}
        for key, value in row.items():
            # Remove potentially dangerous characters
            clean_key = str(key).strip()
            clean_value = str(value).strip()
            clean_row[clean_key] = clean_value
        sanitized.append(clean_row)
    return sanitized
```

### File System Security
- **Path Validation**: Prevent directory traversal attacks
- **Permission Checks**: Verify file read/write permissions
- **Temporary Files**: Use secure temporary file creation
- **File Cleanup**: Ensure temporary files are removed

### Data Privacy
- **No Data Transmission**: All processing is local
- **No Logging of Sensitive Data**: Avoid logging personal information
- **Secure File Handling**: Proper file permissions and cleanup

## Testing Strategy

### Test Categories

#### Unit Tests
- **Individual Class Testing**: Test each class in isolation
- **Method Testing**: Test individual methods with various inputs
- **Edge Case Testing**: Test boundary conditions and error cases

#### Integration Tests
- **End-to-End Testing**: Test complete workflows

- **Component Interaction**: Test how classes work together
- **Data Flow Testing**: Verify data transformation pipeline

#### Performance Tests
- **Load Testing**: Test with large datasets
- **Memory Testing**: Monitor memory usage
- **Stress Testing**: Test with malformed data

### Test Coverage Goals
- **Line Coverage**: >90%
- **Branch Coverage**: >85%
- **Function Coverage**: 100%

### Test Data Strategy
```python
# Test data categories
test_data = {
    'valid_csv': 'valid_data.csv',
    'missing_values': 'missing_data.csv',
    'outliers': 'outlier_data.csv',
    'malformed': 'malformed_data.csv',
    'large_file': 'large_data.csv'
}
```

## Production Readiness

### Current Status: **DEVELOPMENT READY**

####  Production-Ready Features
- Modular architecture
- Error handling
- Data validation
- Statistical methods
- Logging system
- File I/O operations

####  Production Gaps

##### 1. Performance Optimization
```python
# Current: Loads entire file into memory
def load_csv(self, file_path):
    with open(file_path, 'r') as file:
        self.data = list(csv.DictReader(file))  # Memory intensive

# Needed: Chunked processing
def load_csv_chunked(self, file_path, chunk_size=1000):
    for chunk in pd.read_csv(file_path, chunksize=chunk_size):
        yield chunk
```

##### 2. Comprehensive Logging
```python
# Current: Basic print statements
print(" Loaded {len(data)} rows")

# Needed: Structured logging
import logging
logging.info("Data loaded", extra={
    'rows': len(data),
    'columns': len(headers),
    'file_size': os.path.getsize(file_path)
})
```

##### 3. Configuration Management
```python
# Current: Hard-coded values
threshold = 2.0
method = 'zscore'

# Needed: Configuration file
config = {
    'outlier_threshold': 2.0,
    'outlier_method': 'zscore',
    'imputation_method': 'median',
    'normalization_method': 'minmax'
}
```

##### 4. Input Validation
```python
# Current: Basic validation
if not file_path.lower().endswith('.csv'):
    return False

# Needed: Comprehensive validation
def validate_input(file_path, file_size_limit=100*1024*1024):
    # File existence
    # File size
    # File format
    # File permissions
    # Content validation
```

##### 5. Error Recovery
```python
# Current: Stop on error
except Exception as e:
    self.issues.append(f"Error: {str(e)}")
    return False

# Needed: Graceful recovery
try:
    # Operation
except RecoverableError as e:
    # Try alternative method
    return self.fallback_method()
except FatalError as e:
    # Log and stop
    logging.error("Fatal error", exc_info=True)
    return False
```

### Production Readiness Checklist

#### Infrastructure
- [ ] **Monitoring**: Application performance monitoring
- [ ] **Logging**: Centralized logging system
- [ ] **Backup**: Data backup and recovery procedures
- [ ] **Security**: Input validation and sanitization
- [ ] **Performance**: Load testing and optimization

#### Code Quality
- [ ] **Error Handling**: Comprehensive error handling
- [ ] **Input Validation**: Robust input validation
- [ ] **Configuration**: External configuration management
- [ ] **Documentation**: Complete API documentation
- [ ] **Testing**: Comprehensive test suite

#### Operations

- [ ] **Deployment**: Automated deployment pipeline
- [ ] **Monitoring**: Health checks and alerting
- [ ] **Backup**: Data backup procedures
- [ ] **Recovery**: Disaster recovery plan
- [ ] **Security**: Security audit and compliance

## API Reference

### CSVLoader
```python
class CSVLoader:
    def load_csv(self, file_path: str) -> bool
    def get_data(self) -> List[Dict[str, Any]]
    def get_headers(self) -> List[str]
    def get_issues(self) -> List[str]
```

### DataValidator
```python
class DataValidator:
    def validate_data(self, data: List[Dict], headers: List[str]) -> Dict[str, Any]
    def _is_valid_date(self, date_str: str) -> bool
    def _is_valid_email(self, email: str) -> bool
```

### MissingValueImputer
```python
class MissingValueImputer(DataProcessor):
    def process(self, data: List[Dict]) -> List[Dict]
    def get_imputation_log(self) -> List[Dict]
```

### OutlierRemover
```python
class OutlierRemover(DataProcessor):
    def __init__(self, method: str = 'zscore', threshold: float = 2.0)
    def process(self, data: List[Dict]) -> List[Dict]
    def get_outlier_log(self) -> List[Dict]
```

### Normalizer
```python
class Normalizer(DataProcessor):
    def __init__(self, method: str = 'minmax')
    def process(self, data: List[Dict]) -> List[Dict]
    def get_normalization_log(self) -> List[Dict]
```

### ReportGenerator
```python
class ReportGenerator:
    def generate_report(self, original_data, cleaned_data, ...) -> str
    def save_report(self, report: str, output_path: str) -> bool
```

## Configuration

### Current Configuration
All configuration is currently hard-coded in the source code:

```python
# Outlier detection
DEFAULT_THRESHOLD = 2.0
DEFAULT_METHOD = 'zscore'
```

```python
# Imputation
DEFAULT_IMPUTATION_METHOD = 'median'

# Normalization
DEFAULT_NORMALIZATION_METHOD = 'minmax'
```

### Recommended Configuration Structure
```python
# config.yaml
outlier_detection:
  default_method: 'zscore'
  default_threshold: 2.0
  iqr_multiplier: 1.5

imputation:
  numeric_method: 'median'
  categorical_method: 'mode'
  date_method: 'median'

normalization:
  default_method: 'minmax'
  zscore_threshold: 3.0

file_handling:
  max_file_size: 100MB
  supported_encodings: ['utf-8', 'latin-1']
  chunk_size: 1000

logging:
  level: 'INFO'
  format: '%(asctime)s - %(name)s - %(levelname)s - %(message)s'
  file: 'csv_cleaner.log'
```

## Deployment

### Development Deployment
```bash
# Clone repository
git clone <repository-url>
cd csv-data-cleaner

# Install dependencies (none required)
# Run application
python csv_data_cleaner.py
```

### Production Deployment Recommendations

#### 1. Containerization
```dockerfile
FROM python:3.9-slim

WORKDIR /app
COPY . .

RUN pip install --no-cache-dir -r requirements.txt

EXPOSE 8000
CMD ["python", "csv_data_cleaner.py"]
```

#### 2. Environment Configuration
```bash
# Environment variables
```

```bash
export CSV_CLEANER_CONFIG_PATH=/app/config.yaml
export CSV_CLEANER_LOG_LEVEL=INFO
export CSV_CLEANER_MAX_FILE_SIZE=100MB
```

#### 3. Monitoring Setup
```python
# Health check endpoint
@app.route('/health')
def health_check():
    return {'status': 'healthy', 'timestamp': datetime.now()}

# Metrics collection
from prometheus_client import Counter, Histogram

requests_total = Counter('csv_cleaner_requests_total', 'Total requests')
processing_time = Histogram('csv_cleaner_processing_seconds', 'Processing time')
```

#### 4. Backup Strategy
```python
# Automated backup
def backup_original_data(file_path):
    timestamp = datetime.now().strftime('%Y%m%d_%H%M%S')
    backup_path = f"backup/{timestamp}_{os.path.basename(file_path)}"
    shutil.copy2(file_path, backup_path)
    return backup_path
```

### Deployment Checklist
- [ ] **Environment Setup**: Python 3.8+, dependencies
- [ ] **Configuration**: External configuration files
- [ ] **Logging**: Centralized logging setup
- [ ] **Monitoring**: Health checks and metrics
- [ ] **Security**: Input validation and access controls
- [ ] **Backup**: Data backup procedures
- [ ] **Documentation**: Deployment and operations guides

# File: USER_GUIDE.md

# CSV Data Cleaner - User Guide

## Table of Contents

## Getting Started

### Prerequisites
- Python 3.8 or higher
- CSV file with headers
- Basic understanding of data cleaning concepts

### Installation
1. Download the project files
2. Ensure Python is installed: `python --version`
3. No additional packages required (uses standard library only)

### Quick Start
```bash
python csv_data_cleaner.py
```

## Basic Usage

### First Time Users
1. **Load your data**: Choose option 1 and provide your CSV file path
2. **Check for issues**: Choose option 2 to see what problems exist
3. **Fix problems**: Use options 3-5 to clean your data
4. **Get results**: Choose option 6 to generate a report

### Example Session
```
=== CSV Data Cleaner (Group 22) ===

1. Load CSV file
2. Detect issues
3. Fix missing values
4. Remove outliers
5. Normalize data
6. Generate report
0. Exit

Choice (1-6): 1

--- Load CSV File ---
File path (or Enter for data/sample_data.csv): my_data.csv
 Loaded 150 rows

Choice (1-6): 2

--- Detect Issues ---
  Issues found:
    Non-numeric value 'abc' in numeric column 'Age' at row 45
    Invalid email format 'not-an-email' in email column 'Email' at row 67
    Missing value in Age column at row 23
```

## Step-by-Step Workflow

### Step 1: Load Your Data
- **What it does**: Reads your CSV file and validates its structure
- **What you need**: A CSV file with headers in the first row
- **Example input**: `my_data.csv` or press Enter for sample data
- **Success message**: " Loaded X rows"

### Step 2: Detect Issues
- **What it does**: Analyzes your data for quality problems
- **Types of issues found**:
  - Missing values (empty cells)
  - Type mismatches (text in numeric columns)
  - Invalid formats (wrong email/date formats)
  - Data inconsistencies
- **What to expect**: List of specific problems with row numbers

### Step 3: Fix Missing Values
- **What it does**: Fills empty cells using statistical methods
- **Methods used**:
  - **Numeric columns**: Median (middle value)
  - **Text columns**: Mode (most frequent value)
- **What to expect**: " Fixed X missing values"

### Step 4: Remove Outliers
- **What it does**: Identifies and removes extreme values
- **Methods available**:
  - **Z-score (default)**: Removes values beyond 2 standard deviations
  - **IQR**: Removes values outside 1.5  interquartile range
- **What to expect**: " Removed X outliers"

### Step 5: Normalize Data
- **What it does**: Scales numeric data to standard ranges
- **Methods available**:
  - **Min-max**: Scales to 0-1 range
  - **Z-score**: Centers around 0 with standard deviation 1
- **What to expect**: " Normalized X values"

### Step 6: Generate Report
- **What it does**: Creates detailed summary of all changes
- **Report includes**:
  - Summary statistics
  - List of all transformations
  - Column type analysis
  - Recommendations
- **Save options**: Text report and cleaned CSV file

## Understanding the Menu

### Option 1: Load CSV File
- **Purpose**: Import your data for processing
- **Input**: File path (relative or absolute)
- **Default**: `data/sample_data.csv` if you press Enter
- **Requirements**: File must exist and be readable

### Option 2: Detect Issues
- **Purpose**: Find data quality problems
- **Prerequisite**: Must load data first (option 1)
- **Output**: List of specific issues with locations
- **Action**: Review issues before proceeding

### Option 3: Fix Missing Values
- **Purpose**: Fill empty cells automatically
- **Prerequisite**: Must load data first (option 1)
- **Methods**: Automatic selection based on data type

- **Safety**: Original data is preserved

### Option 4: Remove Outliers
- **Purpose**: Remove extreme values that skew analysis
- **Prerequisite**: Must load data first (option 1)
- **Choice**: Z-score (default) or IQR method
- **Impact**: Reduces dataset size

### Option 5: Normalize Data
- **Purpose**: Scale numeric data for analysis
- **Prerequisite**: Must load data first (option 1)
- **Choice**: Min-max (0-1) or Z-score (standardized)
- **Use case**: Machine learning, statistical analysis

### Option 6: Generate Report
- **Purpose**: Create comprehensive summary
- **Prerequisite**: Must load data first (option 1)
- **Output**: Detailed transformation log
- **Save options**: Report file and cleaned data

### Option 0: Exit
- **Purpose**: Safely close the application
- **Note**: Any unsaved changes will be lost

## Data Requirements

### CSV File Format
- **Headers**: Required in first row
- **Encoding**: UTF-8 recommended
- **Delimiter**: Comma (,) standard
- **Size**: No strict limit, but large files may be slow

### Supported Data Types
- **Numeric**: Integers, decimals, percentages
- **Text**: Names, descriptions, categories
- **Dates**: YYYY-MM-DD, MM/DD/YYYY, MM-DD-YYYY
- **Emails**: Standard email format validation

### Column Naming
- **Best practice**: Use descriptive names
- **Special columns**:
  - Columns with "email" in name get email validation
  - Columns with "date" in name get date validation
- **Avoid**: Special characters in column names

## Output Files

### Report File (report.txt)
Contains detailed information about:
- Original vs cleaned data statistics
- All transformations performed
- Issues found and fixed
- Column type analysis
- Timestamps for all operations

### Cleaned Data File (cleaned.csv)
- Same structure as original file
- All transformations applied
- Missing values filled
- Outliers removed (if option 4 used)
- Numeric data normalized (if option 5 used)

## Troubleshooting

### Common Issues

#### "File not found" Error
- **Cause**: Incorrect file path
- **Solution**: Check file exists and path is correct
- **Tip**: Use relative paths or full absolute paths

#### "CSV file has no headers" Error
- **Cause**: First row doesn't contain column names
- **Solution**: Add headers to your CSV file
- **Example**: `Name,Age,Email,Salary`

#### "No data loaded" Error
- **Cause**: Trying to process data before loading
- **Solution**: Always use option 1 first
- **Workflow**: Load  Detect  Fix  Report

#### "No issues found" Message
- **Cause**: Your data is already clean
- **Action**: You can still normalize data (option 5)
- **Note**: This is good news!

#### "No outliers found" Message
- **Cause**: Your data doesn't have extreme values
- **Action**: Continue with normalization (option 5)
- **Note**: This is normal for clean datasets

### Performance Issues

#### Slow Processing
- **Cause**: Large file size
- **Solutions**:
  - Process smaller chunks
  - Close other applications
  - Use SSD storage if available

#### Memory Errors
- **Cause**: File too large for available memory
- **Solutions**:
  - Increase system RAM
  - Process smaller files
  - Use 64-bit Python

## Best Practices

### Before Processing
1. **Backup your data**: Always keep original files
2. **Review your data**: Understand what you're working with
3. **Check file size**: Large files may need special handling
4. **Validate format**: Ensure CSV is properly formatted

### During Processing
1. **Follow the workflow**: Load  Detect  Fix  Report
2. **Review issues**: Understand what problems exist
3. **Choose methods wisely**: Different methods for different data types
4. **Save results**: Don't lose your cleaned data

### After Processing
1. **Review the report**: Understand what changed
2. **Validate results**: Check that transformations make sense
3. **Backup cleaned data**: Save your processed files
4. **Document changes**: Keep notes on what was done

### Data Quality Tips
1. **Consistent formatting**: Use same date/email formats
2. **Meaningful headers**: Clear, descriptive column names
3. **Data types**: Ensure numeric columns contain numbers
4. **Missing values**: Decide how to handle empty cells

### File Management
1. **Organize files**: Use clear naming conventions
2. **Version control**: Keep track of different versions
3. **Backup strategy**: Regular backups of important data
4. **Storage**: Use reliable storage media

## Advanced Usage

### Batch Processing
For multiple files, you can:
1. Process each file individually
2. Save reports with descriptive names
3. Compare results across files
4. Create summary reports

### Custom Workflows
Common processing sequences:
- **Basic cleaning**: Load  Detect  Fix missing  Report
- **Statistical analysis**: Load  Detect  Fix missing  Remove outliers  Normaliz
e  Report
- **Quality check**: Load  Detect  Report (no changes)

### Integration
The cleaned data can be used for:
- Statistical analysis
- Machine learning models
- Business intelligence
- Data visualization
- Further processing

## Support

### Getting Help
- Review this user guide
- Check the troubleshooting section
- Examine error messages carefully
- Test with sample data first

### Reporting Issues
When reporting problems, include:
- Error message text
- File size and format
- Steps to reproduce
- System information (Python version, OS)

### Feature Requests
For new features, consider:
- Use case description
- Expected behavior
- Sample data if applicable
- Priority level

# File: csv_data_cleaner.py

```python
#!/usr/bin/env python3
"""
CSV Data Cleaner and Validator (Group 22)

A Python CLI application that processes CSV files and detects & corrects data qu
ality issues.
Handles missing values, outliers, normalization, and generates reports.

Author: Group 22
"""

from src.cli import CLI


def main():
    """Main function to run the CLI application."""
    cli = CLI()
    cli.run()


if __name__ == "__main__":
    main()
```

# File: data\sample_data.csv

```
Name,Age,Email,JoinDate,Salary,Department
Alice,29,alice@email.com,2023-04-01,54000,Engineering
Bob,,bob@email.com,04/01/2023,52000,Marketing
,32,not-an-email,2023-01-15,70000,Sales
Charlie,28,charlie@company.com,2023-06-15,48000,Engineering
Diana,35,diana@email.com,2023-03-20,85000,Management
Eve,26,eve@test.com,2023-08-10,45000,Engineering
Frank,31,frank@email.com,2023-02-28,58000,Marketing
Grace,29,grace@company.com,2023-07-05,52000,Sales
Henry,33,henry@email.com,2023-05-12,65000,Engineering
Ivy,27,ivy@test.com,2023-09-01,47000,Marketing
Jack,30,jack@email.com,2023-01-10,55000,Sales
Kate,34,kate@company.com,2023-04-25,72000,Management
Leo,28,leo@email.com,2023-06-30,49000,Engineering
Maya,31,maya@test.com,2023-03-15,56000,Marketing
Noah,29,noah@email.com,2023-08-20,51000,Sales
Olivia,36,olivia@company.com,2023-02-10,78000,Management
Paul,27,paul@email.com,2023-07-25,46000,Engineering
Quinn,32,quinn@test.com,2023-05-08,59000,Marketing
Ruby,30,ruby@email.com,2023-09-15,53000,Sales
Sam,35,sam@company.com,2023-01-30,82000,Management
```

# File: generate_pdf.py

```python
#!/usr/bin/env python3
"""
PDF Generator Script for CSV Data Cleaner Codebase (Group 22)

This script generates a comprehensive PDF of the entire codebase with clean form
atting using fpdf2.
"""

import os
import glob
from fpdf import FPDF
import datetime

SECTION_TITLES = [
    ("1. Project Overview", "overview"),
    ("2. Main Application Files", "main"),
    ("3. Source Code Modules", "src"),
    ("4. Test Files", "tests"),
    ("5. Documentation", "docs"),
    ("6. Sample Data", "data"),
]

class CodebasePDF(FPDF):
    def header(self):
        if self.page_no() == 1:
            return  # No header on title page
        self.set_font('helvetica', 'B', 12)
        self.set_text_color(40, 40, 80)
        self.cell(0, 10, 'CSV Data Cleaner and Validator - Group 22', 0, 1, 'C')
        self.ln(2)

    def footer(self):
        if self.page_no() == 1:
            return  # No footer on title page
        self.set_y(-15)
        self.set_font('helvetica', 'I', 9)
        self.set_text_color(100, 100, 100)
        self.cell(0, 10, f'Page {self.page_no()}/{56}', 0, 0, 'C')

    def add_section_title(self, title):
        self.set_font('helvetica', 'B', 16)
        self.set_text_color(0, 51, 102)
        self.cell(0, 12, title, 0, 1, 'L')
        self.ln(2)

    def add_file_title(self, filename):
        self.set_font('helvetica', 'B', 12)
        self.set_text_color(0, 102, 51)
        self.cell(0, 8, f'File: {filename}', 0, 1, 'L')
        self.ln(1)

    def add_code_block(self, content):
        self.set_font('courier', '', 9)
        self.set_text_color(30, 30, 30)
        max_chars = 110
        for line in content.split('\n'):
            # Forcibly break up any line longer than max_chars
            while len(line) > max_chars:
                self.cell(0, 4.5, line[:max_chars], 0, 1)
                line = line[max_chars:]
            self.cell(0, 4.5, line, 0, 1)
        self.ln(2)
```

```python
    def add_body_text(self, text):
        self.set_font('helvetica', '', 11)
        self.set_text_color(20, 20, 20)
        max_chars = 100
        for line in text.split('\n'):
            # Forcibly break up any line longer than max_chars
            while len(line) > max_chars:
                self.cell(0, 6, line[:max_chars], 0, 1)
                line = line[max_chars:]
            self.cell(0, 6, line, 0, 1)
        self.ln(2)

    def safe_add_page(self):
        # Only add a new page if not at the top of a fresh page
        if self.get_y() > 30:
            self.add_page()


def get_all_files():
    files = []
    files.extend(glob.glob("**/*.py", recursive=True))
    files.extend(glob.glob("*.md"))
    files.extend(glob.glob("**/*.csv", recursive=True))
    files.extend(glob.glob("*.txt"))
    return sorted(set(files))

def read_file_content(filepath):
    try:
        with open(filepath, 'r', encoding='utf-8') as f:
            return f.read()
    except UnicodeDecodeError:
        try:
            with open(filepath, 'r', encoding='latin-1') as f:
                return f.read()
        except:
            return f"Error reading file: {filepath}"

def create_pdf():
    pdf = CodebasePDF(format='A4', orientation='P', unit='mm')
    pdf.set_auto_page_break(auto=True, margin=18)
    pdf.alias_nb_pages()
    pdf.add_page()

    # Title page
    pdf.set_font('helvetica', 'B', 22)
    pdf.set_text_color(0, 51, 102)
    pdf.cell(0, 20, 'CSV Data Cleaner and Validator', 0, 1, 'C')
    pdf.set_font('helvetica', 'B', 16)
    pdf.set_text_color(0, 102, 51)
    pdf.cell(0, 12, 'Group 22 - Complete Codebase', 0, 1, 'C')
    pdf.ln(10)
    pdf.set_font('helvetica', '', 12)
    pdf.set_text_color(0, 0, 0)
    pdf.cell(0, 10, f'Generated: {datetime.datetime.now().strftime("%Y-%m-%d %H:
%M:%S")}', 0, 1, 'C')
    pdf.cell(0, 10, 'Python CLI Application for Data Cleaning and Validation', 0
, 1, 'C')
    pdf.ln(18)
    pdf.set_font('helvetica', 'I', 11)
    pdf.cell(0, 8, 'All code and documentation in one place.', 0, 1, 'C')
    pdf.ln(10)

    # Table of Contents
    pdf.set_font('helvetica', 'B', 15)
    pdf.set_text_color(0, 51, 102)
    pdf.cell(0, 10, 'Table of Contents', 0, 1, 'L')
```

```python
    pdf.ln(2)
    pdf.set_font('helvetica', '', 12)
    pdf.set_text_color(0, 0, 0)
    for title, _ in SECTION_TITLES:
        pdf.cell(0, 8, title, 0, 1, 'L')
    pdf.ln(4)

    # Section: Project Overview
    pdf.add_page()
    pdf.add_section_title("1. Project Overview")
    overview_text = (
        "This is a Python CLI application for cleaning and validating CSV files.
 "
        "The project follows object-oriented programming principles and uses onl
y Python standard library modules.\n\n"
        "Key Features:\n"
        " Missing value imputation using statistical methods\n"
        " Data validation and type checking\n"
        " Outlier detection and removal\n"
        " Data normalization and standardization\n"
        " Comprehensive reporting and logging\n"
        " Interactive CLI interface\n\n"
        "Project Structure:\n"
        " src/ - Source code modules\n"
        " tests/ - Unit tests\n"
        " data/ - Sample data files\n"
        " Documentation files (README.md, USER_GUIDE.md, TECHNICAL_DOCUMENTATION
.md)\n"
    )
    pdf.add_body_text(overview_text)

    # Gather files for each section
    all_files = get_all_files()
    main_files = [f for f in all_files if f in ['csv_data_cleaner.py', 'README.m
d']]
    src_files = [f for f in all_files if f.startswith('src/')]
    test_files = [f for f in all_files if f.startswith('tests/')]
    doc_files = [f for f in all_files if f.endswith('.md') and f != 'README.md']
    data_files = [f for f in all_files if f.endswith('.csv')]

    # Section: Main Application Files
    if main_files:
        pdf.safe_add_page()
        pdf.add_section_title("2. Main Application Files")
        for filepath in main_files:
            pdf.add_file_title(filepath)
            content = read_file_content(filepath)
            pdf.add_code_block(content)

    # Section: Source Code Modules
    if src_files:
        pdf.safe_add_page()
        pdf.add_section_title("3. Source Code Modules")
        for filepath in src_files:
            pdf.add_file_title(filepath)
            content = read_file_content(filepath)
            pdf.add_code_block(content)

    # Section: Test Files
    if test_files:
        pdf.safe_add_page()
        pdf.add_section_title("4. Test Files")
        for filepath in test_files:
            pdf.add_file_title(filepath)
            content = read_file_content(filepath)
            pdf.add_code_block(content)
```

```python
        # Section: Documentation
        if doc_files:
            pdf.safe_add_page()
            pdf.add_section_title("5. Documentation")
            for filepath in doc_files:
                pdf.add_file_title(filepath)
                content = read_file_content(filepath)
                pdf.add_code_block(content)

        # Section: Sample Data
        if data_files:
            pdf.safe_add_page()
            pdf.add_section_title("6. Sample Data")
            for filepath in data_files:
                pdf.add_file_title(filepath)
                content = read_file_content(filepath)
                pdf.add_code_block(content)

        pdf.output("Group_22_CSV_Data_Cleaner_Codebase.pdf")
        print("PDF generated successfully: Group_22_CSV_Data_Cleaner_Codebase.pdf")

if __name__ == "__main__":
    create_pdf()
```

# File: generate_pdf_simple.py

```python
#!/usr/bin/env python3
"""
Simple PDF Generator for CSV Data Cleaner Codebase (Group 22)
"""

import os
import glob
from fpdf import FPDF
import datetime

def get_all_files():
    files = []
    files.extend(glob.glob("**/*.py", recursive=True))
    files.extend(glob.glob("*.md"))
    files.extend(glob.glob("**/*.csv", recursive=True))
    return sorted(set(files))

def read_file_content(filepath):
    try:
        with open(filepath, 'r', encoding='utf-8') as f:
            content = f.read()
            # Remove any problematic characters
            return content.encode('ascii', 'ignore').decode('ascii')
    except:
        return f"Error reading file: {filepath}"

def create_pdf():
    pdf = FPDF()
    pdf.add_page()

    # Title
    pdf.set_font('Arial', 'B', 20)
    pdf.cell(0, 20, 'CSV Data Cleaner and Validator - Group 22', 0, 1, 'C')
    pdf.ln(10)

    pdf.set_font('Arial', '', 12)
    pdf.cell(0, 10, f'Generated: {datetime.datetime.now().strftime("%Y-%m-%d %H:
%M:%S")}', 0, 1, 'C')
    pdf.ln(10)

    # Get files
    all_files = get_all_files()

    # Add each file
    for filepath in all_files:
        pdf.add_page()
        pdf.set_font('Arial', 'B', 14)
        pdf.cell(0, 10, f'File: {filepath}', 0, 1, 'L')
        pdf.ln(5)

        content = read_file_content(filepath)
        pdf.set_font('Courier', '', 8)

        # Split content into manageable chunks
        lines = content.split('\n')
        for line in lines:
            # Break long lines
            while len(line) > 80:
                pdf.cell(0, 4, line[:80], 0, 1)
                line = line[80:]
            pdf.cell(0, 4, line, 0, 1)

    pdf.output("Group_22_Codebase.pdf")
```

```python
        print("PDF generated: Group_22_Codebase.pdf")

if __name__ == "__main__":
    create_pdf()
```

# File: src\__init__.py

```python
"""
Source package for CSV Data Cleaner and Validator (Group 22)

This package contains all the source code for the CSV Data Cleaner application.
"""

from .cli import CLI
from .data_processors import (
    CSVLoader, DataValidator, MissingValueImputer,
    OutlierRemover, Normalizer, ReportGenerator
)

__all__ = [
    'CLI',
    'CSVLoader',
    'DataValidator',
    'MissingValueImputer',
    'OutlierRemover',
    'Normalizer',
    'ReportGenerator'
]
```

# File: src\cli.py

```python
"""
CLI Module for CSV Data Cleaner and Validator (Group 22)

This module contains the CLI class which manages user interaction and menu syste
m.
"""

import csv
from .data_processors import (
    CSVLoader, DataValidator, MissingValueImputer,
    OutlierRemover, Normalizer, ReportGenerator
)


class CLI:
    """Manages user interaction and menu system."""

    def __init__(self):
        self.csv_loader = CSVLoader()
        self.data_validator = DataValidator()
        self.imputer = MissingValueImputer()
        self.outlier_remover = OutlierRemover()
        self.normalizer = Normalizer()
        self.report_generator = ReportGenerator()
        self.current_data = []
        self.validation_results = {}

    def run(self):
        """Run the main CLI interface."""
        print("=== CSV Data Cleaner (Group 22) ===")

        while True:
            self._print_menu()
            choice = input("\nChoice (1-6): ").strip()

            if choice == '1':
                self._load_csv()
            elif choice == '2':
                self._detect_issues()
            elif choice == '3':
                self._impute_missing_values()
            elif choice == '4':
                self._remove_outliers()
            elif choice == '5':
                self._normalize_data()
            elif choice == '6':
                self._generate_report()
            elif choice == '0':
                print("Goodbye!")
                break
            else:
                print("Invalid choice. Please enter 1-6 or 0 to exit.")

    def _print_menu(self):
        """Print the simplified main menu."""
        print("\n1. Load CSV file")
        print("2. Detect issues")
        print("3. Fix missing values")
        print("4. Remove outliers")
        print("5. Normalize data")
        print("6. Generate report")
        print("0. Exit")
```

```python
    def _load_csv(self):
        """Handle CSV file loading."""
        print("\n--- Load CSV File ---")

        file_path = input("File path (or Enter for data/sample_data.csv): ").str
ip()
        if not file_path:
            file_path = "data/sample_data.csv"

        if self.csv_loader.load_csv(file_path):
            self.current_data = self.csv_loader.get_data()
            print(f" Loaded {len(self.current_data)} rows")
        else:
            print(" Failed to load file:")
            for issue in self.csv_loader.get_issues():
                print(f"   {issue}")

    def _detect_issues(self):
        """Handle data quality issue detection."""
        print("\n--- Detect Issues ---")

        if not self.current_data:
            print(" No data loaded. Use option 1 first.")
            return

        self.validation_results = self.data_validator.validate_data(
            self.current_data, self.csv_loader.get_headers()
        )

        if self.validation_results['valid']:
            print(" No issues found")
        else:
            print("  Issues found:")
            for issue in self.validation_results['issues']:
                print(f"   {issue}")

    def _impute_missing_values(self):
        """Handle missing value imputation."""
        print("\n--- Fix Missing Values ---")

        if not self.current_data:
            print(" No data loaded. Use option 1 first.")
            return

        self.current_data = self.imputer.process(self.current_data)
        log = self.imputer.get_imputation_log()

        if log:
            print(f" Fixed {len(log)} missing values")
        else:
            print(" No missing values found")

    def _remove_outliers(self):
        """Handle outlier removal."""
        print("\n--- Remove Outliers ---")

        if not self.current_data:
            print(" No data loaded. Use option 1 first.")
            return

        method = input("Method (1=Z-score, 2=IQR): ").strip()
        if method == '2':
            self.outlier_remover = OutlierRemover(method='iqr')

        original_count = len(self.current_data)
        self.current_data = self.outlier_remover.process(self.current_data)
```

```python
        removed = original_count - len(self.current_data)

        if removed > 0:
            print(f" Removed {removed} outliers")
        else:
            print(" No outliers found")

    def _normalize_data(self):
        """Handle data normalization."""
        print("\n--- Normalize Data ---")

        if not self.current_data:
            print(" No data loaded. Use option 1 first.")
            return

        method = input("Method (1=Min-max, 2=Z-score): ").strip()
        if method == '2':
            self.normalizer = Normalizer(method='zscore')

        self.current_data = self.normalizer.process(self.current_data)
        log = self.normalizer.get_normalization_log()

        if log:
            print(f" Normalized {len(log)} values")
        else:
            print(" No numeric data to normalize")

    def _generate_report(self):
        """Handle report generation."""
        print("\n--- Generate Report ---")

        if not self.current_data:
            print(" No data loaded. Use option 1 first.")
            return

        report = self.report_generator.generate_report(
            original_data=self.csv_loader.get_data(),
            cleaned_data=self.current_data,
            validation_results=self.validation_results,
            imputation_log=self.imputer.get_imputation_log(),
            outlier_log=self.outlier_remover.get_outlier_log(),
            normalization_log=self.normalizer.get_normalization_log()
        )

        print(" Report generated:")
        print(report)

        save = input("\nSave report? (y/n): ").strip().lower()
        if save == 'y':
            filename = input("Filename (default: report.txt): ").strip()
            if not filename:
                filename = "report.txt"

            if self.report_generator.save_report(report, filename):
                print(f" Saved to {filename}")
            else:
                print(" Failed to save")

        save_csv = input("Save cleaned data? (y/n): ").strip().lower()
        if save_csv == 'y':
            filename = input("Filename (default: cleaned.csv): ").strip()
            if not filename:
                filename = "cleaned.csv"

            try:
                with open(filename, 'w', newline='', encoding='utf-8') as file:
```

```python
                writer = csv.DictWriter(file, fieldnames=self.csv_loader.get
_headers())
                writer.writeheader()
                writer.writerows(self.current_data)
            print(f" Saved to {filename}")
        except Exception as e:
            print(f" Failed to save: {e}")
```

# File: src\data_processors\__init__.py

```python
"""
Data Processors Package for CSV Data Cleaner and Validator (Group 22)

This package contains all the data processing classes used by the CSV Data Clean
er.
Each class handles a specific aspect of data cleaning and validation.

Classes:
- CSVLoader: Handles loading and basic validation of CSV files
- DataValidator: Validates data types and formats
- MissingValueImputer: Handles imputation of missing values
- OutlierRemover: Detects and removes outliers
- Normalizer: Performs normalization and standardization
- ReportGenerator: Generates comprehensive data quality reports
"""

from .csv_loader import CSVLoader
from .data_validator import DataValidator
from .missing_value_imputer import MissingValueImputer
from .outlier_remover import OutlierRemover
from .normalizer import Normalizer
from .report_generator import ReportGenerator

__all__ = [
    'CSVLoader',
    'DataValidator',
    'MissingValueImputer',
    'OutlierRemover',
    'Normalizer',
    'ReportGenerator'
]
```

# File: src\data_processors\base.py

```python
"""
Base classes for data processing operations.

This module contains the abstract base class that all data processors inherit fr
om.
"""

from abc import ABC, abstractmethod
from typing import List, Dict, Any


class DataProcessor(ABC):
    """Abstract base class for data processing operations."""

    @abstractmethod
    def process(self, data: List[Dict[str, Any]]) -> List[Dict[str, Any]]:
        """
        Process the data and return the result.

        Args:
            data: List of dictionaries representing rows

        Returns:
            List of dictionaries with processed data
        """
        pass
```

# File: src\data_processors\csv_loader.py

```python
"""
CSV Loader Module for CSV Data Cleaner and Validator (Group 22)

This module contains the CSVLoader class which handles loading and basic validat
ion of CSV files.
"""

import csv
import os
from typing import List, Dict, Any


class CSVLoader:
    """Handles loading and basic validation of CSV files."""

    def __init__(self):
        self.data = []
        self.headers = []
        self.file_path = ""
        self.issues = []

    def load_csv(self, file_path: str) -> bool:
        """
        Load CSV file and perform basic validation.

        Args:
            file_path: Path to the CSV file

        Returns:
            bool: True if successful, False otherwise
        """
        try:
            if not os.path.exists(file_path):
                self.issues.append(f"File not found: {file_path}")
                return False

            if not file_path.lower().endswith('.csv'):
                self.issues.append("File must have .csv extension")
                return False

            with open(file_path, 'r', newline='', encoding='utf-8') as file:
                reader = csv.DictReader(file)
                self.headers = reader.fieldnames

                if not self.headers:
                    self.issues.append("CSV file has no headers")
                    return False

                self.data = list(reader)
                self.file_path = file_path

                if not self.data:
                    self.issues.append("CSV file is empty")
                    return False

                return True

        except Exception as e:
            self.issues.append(f"Error reading CSV file: {str(e)}")
            return False

    def get_data(self) -> List[Dict[str, Any]]:
        """Return the loaded data."""
```

```python
        return self.data

    def get_headers(self) -> List[str]:
        """Return the column headers."""
        return self.headers

    def get_issues(self) -> List[str]:
        """Return any loading issues."""
        return self.issues
```

# File: src\data_processors\data_validator.py

```python
"""
Data Validator Module for CSV Data Cleaner and Validator (Group 22)

This module contains the DataValidator class which validates data types and form
ats.
"""

import re
from typing import List, Dict, Any


class DataValidator:
    """Validates data types and formats."""

    def __init__(self):
        self.validation_issues = []
        self.column_types = {}

    def validate_data(self, data: List[Dict[str, Any]], headers: List[str]) -> D
ict[str, Any]:
        """
        Validate data types and formats for each column.

        Args:
            data: List of dictionaries representing rows
            headers: List of column headers

        Returns:
            Dict containing validation results
        """
        self.validation_issues = []
        self.column_types = {}

        if not data or not headers:
            self.validation_issues.append("No data or headers to validate")
            return {"valid": False, "issues": self.validation_issues}

        # Analyze each column
        for header in headers:
            column_data = [row.get(header, '') for row in data]
            self._analyze_column(header, column_data)

        return {
            "valid": len(self.validation_issues) == 0,
            "issues": self.validation_issues,
            "column_types": self.column_types
        }

    def _analyze_column(self, header: str, column_data: List[Any]) -> None:
        """Analyze a single column for type and format issues."""
        # Remove empty values for analysis
        non_empty_data = [str(val).strip() for val in column_data if val and str
(val).strip()]

        if not non_empty_data:
            self.column_types[header] = "unknown"
            return

        # Check for numeric data
        numeric_count = 0
        for val in non_empty_data:
            try:
                float(val)
```

```python
                    numeric_count += 1
                except ValueError:
                    pass

        numeric_ratio = numeric_count / len(non_empty_data)

        if numeric_ratio > 0.8:
            self.column_types[header] = "numeric"
            # Check for type inconsistencies
            for i, val in enumerate(column_data):
                if val and str(val).strip():
                    try:
                        float(val)
                    except ValueError:
                        self.validation_issues.append(
                            f"Non-numeric value '{val}' in numeric column '{head
er}' at row {i+1}"
                        )
        else:
            # Check for date format
            date_count = 0
            for val in non_empty_data:
                if self._is_valid_date(val):
                    date_count += 1

            date_ratio = date_count / len(non_empty_data)

            if date_ratio > 0.8:
                self.column_types[header] = "datetime"
                # Check for date format inconsistencies
                for i, val in enumerate(column_data):
                    if val and str(val).strip() and not self._is_valid_date(val)
:
                        self.validation_issues.append(
                            f"Invalid date format '{val}' in date column '{heade
r}' at row {i+1}"
                        )
            else:
                self.column_types[header] = "categorical"

                # Check for email format if column name suggests it
                if 'email' in header.lower():
                    for i, val in enumerate(column_data):
                        if val and str(val).strip() and not self._is_valid_email
(val):
                            self.validation_issues.append(
                                f"Invalid email format '{val}' in email column '
{header}' at row {i+1}"
                            )

    def _is_valid_date(self, date_str: str) -> bool:
        """Check if string is a valid date format."""
        date_patterns = [
            r'^\d{4}-\d{2}-\d{2}$',  # YYYY-MM-DD
            r'^\d{2}/\d{2}/\d{4}$',  # MM/DD/YYYY
            r'^\d{2}-\d{2}-\d{4}$',  # MM-DD-YYYY
        ]

        for pattern in date_patterns:
            if re.match(pattern, date_str):
                return True
        return False

    def _is_valid_email(self, email: str) -> bool:
        """Check if string is a valid email format."""
        email_pattern = r'^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}$'
```

```python
    return bool(re.match(email_pattern, email))
```

# File: src\data_processors\missing_value_imputer.py

```python
"""
Missing Value Imputer Module for CSV Data Cleaner and Validator (Group 22)

This module contains the MissingValueImputer class which handles imputation of m
issing values.
"""

import datetime
import statistics
import collections
from typing import List, Dict, Any
from .base import DataProcessor


class MissingValueImputer(DataProcessor):
    """Handles imputation of missing values using statistical methods."""

    def __init__(self):
        self.imputation_log = []

    def process(self, data: List[Dict[str, Any]]) -> List[Dict[str, Any]]:
        """
        Impute missing values using appropriate statistical methods.

        Args:
            data: List of dictionaries representing rows

        Returns:
            List of dictionaries with imputed values
        """
        if not data:
            return data

        headers = list(data[0].keys())
        cleaned_data = []

        for row in data:
            cleaned_row = row.copy()

            for header in headers:
                value = row.get(header, '')

                # Check if value is missing or empty
                if not value or str(value).strip() == '':
                    imputed_value = self._impute_value(header, data)
                    cleaned_row[header] = imputed_value

                    self.imputation_log.append({
                        'timestamp': datetime.datetime.now(),
                        'column': header,
                        'row': len(cleaned_data) + 1,
                        'original_value': value,
                        'imputed_value': imputed_value,
                        'method': self._get_imputation_method(header, data)
                    })

            cleaned_data.append(cleaned_row)

        return cleaned_data

    def _impute_value(self, column: str, data: List[Dict[str, Any]]) -> Any:
        """Determine appropriate imputation method and return imputed value."""
        # Extract non-empty values for the column
```

```python
        values = [row.get(column, '') for row in data]
        non_empty_values = [v for v in values if v and str(v).strip()]

        if not non_empty_values:
            return ''

        # Try to determine if numeric
        numeric_values = []
        for val in non_empty_values:
            try:
                numeric_values.append(float(val))
            except ValueError:
                pass

        if numeric_values:
            # Use median for numeric data (more robust than mean)
            return statistics.median(numeric_values)
        else:
            # Use mode for categorical data
            value_counts = collections.Counter(non_empty_values)
            return value_counts.most_common(1)[0][0]

    def _get_imputation_method(self, column: str, data: List[Dict[str, Any]]) ->
str:
        """Get the imputation method used for a column."""
        values = [row.get(column, '') for row in data]
        non_empty_values = [v for v in values if v and str(v).strip()]

        if not non_empty_values:
            return 'no_data'

        numeric_values = []
        for val in non_empty_values:
            try:
                numeric_values.append(float(val))
            except ValueError:
                pass

        if numeric_values:
            return 'median'
        else:
            return 'mode'

    def get_imputation_log(self) -> List[Dict[str, Any]]:
        """Return the imputation log."""
        return self.imputation_log
```

# File: src\data_processors\normalizer.py

```python
"""
Normalizer Module for CSV Data Cleaner and Validator (Group 22)

This module contains the Normalizer class which performs normalization and stand
ardization.
"""

import datetime
import statistics
from typing import List, Dict, Any
from .base import DataProcessor


class Normalizer(DataProcessor):
    """Performs normalization and standardization of numeric data."""

    def __init__(self, method: str = 'minmax'):
        """
        Initialize normalizer.

        Args:
            method: 'minmax' for min-max normalization or 'zscore' for standardi
zation
        """
        self.method = method
        self.normalization_log = []
        self.scaling_params = {}

    def process(self, data: List[Dict[str, Any]]) -> List[Dict[str, Any]]:
        """
        Normalize numeric columns.

        Args:
            data: List of dictionaries representing rows

        Returns:
            List of dictionaries with normalized values
        """
        if not data:
            return data

        headers = list(data[0].keys())
        cleaned_data = []

        # Identify numeric columns
        numeric_columns = []
        for header in headers:
            values = [row.get(header, '') for row in data]
            numeric_values = []
            for val in values:
                try:
                    if val and str(val).strip():
                        numeric_values.append(float(val))
                except ValueError:
                    pass

            if len(numeric_values) > len(values) * 0.5:  # More than 50% numeric
                numeric_columns.append(header)

        # Normalize each numeric column
        for column in numeric_columns:
            self._normalize_column(data, column)
```

```python
        return data

    def _normalize_column(self, data: List[Dict[str, Any]], column: str) -> None
:
        """Normalize a specific column."""
        values = []
        for i, row in enumerate(data):
            val = row.get(column, '')
            try:
                if val and str(val).strip():
                    values.append((float(val), i))
            except ValueError:
                pass

        if len(values) < 2:
            return

        numeric_values = [v[0] for v in values]

        if self.method == 'minmax':
            min_val = min(numeric_values)
            max_val = max(numeric_values)

            if max_val == min_val:
                return

            for val, row_idx in values:
                normalized_val = (val - min_val) / (max_val - min_val)
                data[row_idx][column] = normalized_val

                self.normalization_log.append({
                    'timestamp': datetime.datetime.now(),
                    'column': column,
                    'row': row_idx + 1,
                    'original_value': val,
                    'normalized_value': normalized_val,
                    'method': 'minmax',
                    'min': min_val,
                    'max': max_val
                })

            self.scaling_params[column] = {
                'method': 'minmax',
                'min': min_val,
                'max': max_val
            }

        else:  # zscore
            mean_val = statistics.mean(numeric_values)
            std_val = statistics.stdev(numeric_values)

            if std_val == 0:
                return

            for val, row_idx in values:
                z_score = (val - mean_val) / std_val
                data[row_idx][column] = z_score

                self.normalization_log.append({
                    'timestamp': datetime.datetime.now(),
                    'column': column,
                    'row': row_idx + 1,
                    'original_value': val,
                    'standardized_value': z_score,
                    'method': 'zscore',
                    'mean': mean_val,
```

```python
                    'std': std_val
                })

            self.scaling_params[column] = {
                'method': 'zscore',
                'mean': mean_val,
                'std': std_val
            }

    def get_normalization_log(self) -> List[Dict[str, Any]]:
        """Return the normalization log."""
        return self.normalization_log
```

# File: src\data_processors\outlier_remover.py

```python
"""
Outlier Remover Module for CSV Data Cleaner and Validator (Group 22)

This module contains the OutlierRemover class which detects and removes outliers
.
"""

import datetime
import statistics
from typing import List, Dict, Any, Tuple
from .base import DataProcessor


class OutlierRemover(DataProcessor):
    """Detects and removes outliers using statistical methods."""

    def __init__(self, method: str = 'zscore', threshold: float = 2.0):
        """
        Initialize outlier remover.

        Args:
            method: 'zscore' or 'iqr'
            threshold: Threshold for outlier detection
        """
        self.method = method
        self.threshold = threshold
        self.outlier_log = []

    def process(self, data: List[Dict[str, Any]]) -> List[Dict[str, Any]]:
        """
        Remove outliers from numeric columns.

        Args:
            data: List of dictionaries representing rows

        Returns:
            List of dictionaries with outliers removed
        """
        if not data:
            return data

        headers = list(data[0].keys())
        cleaned_data = []

        # Identify numeric columns
        numeric_columns = []
        for header in headers:
            values = [row.get(header, '') for row in data]
            numeric_values = []
            for val in values:
                try:
                    if val and str(val).strip():
                        numeric_values.append(float(val))
                except ValueError:
                    pass

            if len(numeric_values) > len(values) * 0.5:  # More than 50% numeric
                numeric_columns.append(header)

        # Remove outliers from each numeric column
        for column in numeric_columns:
            outliers = self._detect_outliers(data, column)
            self.outlier_log.extend(outliers)
```

```python
        # Filter out rows with outliers
        outlier_rows = set()
        for outlier in self.outlier_log:
            outlier_rows.add(outlier['row'])

        cleaned_data = [row for i, row in enumerate(data) if i + 1 not in outlie
r_rows]

        return cleaned_data

    def _detect_outliers(self, data: List[Dict[str, Any]], column: str) -> List[
Dict[str, Any]]:
        """Detect outliers in a specific column."""
        values = []
        for i, row in enumerate(data):
            val = row.get(column, '')
            try:
                if val and str(val).strip():
                    values.append((float(val), i + 1))
            except ValueError:
                pass

        if len(values) < 3:
            return []

        numeric_values = [v[0] for v in values]

        if self.method == 'zscore':
            return self._zscore_outliers(numeric_values, values, column)
        else:
            return self._iqr_outliers(numeric_values, values, column)

    def _zscore_outliers(self, numeric_values: List[float], values: List[Tuple[f
loat, int]], column: str) -> List[Dict[str, Any]]:
        """Detect outliers using Z-score method."""
        mean_val = statistics.mean(numeric_values)
        std_val = statistics.stdev(numeric_values) if len(numeric_values) > 1 el
se 0

        if std_val == 0:
            return []

        outliers = []
        for val, row_num in values:
            z_score = abs((val - mean_val) / std_val)
            if z_score > self.threshold:
                outliers.append({
                    'timestamp': datetime.datetime.now(),
                    'column': column,
                    'row': row_num,
                    'value': val,
                    'method': 'zscore',
                    'z_score': z_score
                })

        return outliers

    def _iqr_outliers(self, numeric_values: List[float], values: List[Tuple[floa
t, int]], column: str) -> List[Dict[str, Any]]:
        """Detect outliers using IQR method."""
        sorted_values = sorted(numeric_values)
        q1 = sorted_values[len(sorted_values) // 4]
        q3 = sorted_values[3 * len(sorted_values) // 4]
        iqr = q3 - q1
```

```python
        lower_bound = q1 - 1.5 * iqr
        upper_bound = q3 + 1.5 * iqr

        outliers = []
        for val, row_num in values:
            if val < lower_bound or val > upper_bound:
                outliers.append({
                    'timestamp': datetime.datetime.now(),
                    'column': column,
                    'row': row_num,
                    'value': val,
                    'method': 'iqr',
                    'lower_bound': lower_bound,
                    'upper_bound': upper_bound
                })

        return outliers

    def get_outlier_log(self) -> List[Dict[str, Any]]:
        """Return the outlier detection log."""
        return self.outlier_log
```

# File: src\data_processors\report_generator.py

```python
"""
Report Generator Module for CSV Data Cleaner and Validator (Group 22)

This module contains the ReportGenerator class which generates comprehensive dat
a quality reports.
"""

import datetime
import collections
from typing import List, Dict, Any


class ReportGenerator:
    """Generates comprehensive data quality reports."""

    def __init__(self):
        self.report_log = []

    def generate_report(self,
                        original_data: List[Dict[str, Any]],
                        cleaned_data: List[Dict[str, Any]],
                        validation_results: Dict[str, Any],
                        imputation_log: List[Dict[str, Any]],
                        outlier_log: List[Dict[str, Any]],
                        normalization_log: List[Dict[str, Any]]) -> str:
        """
        Generate a comprehensive data quality report.

        Args:
            original_data: Original data before cleaning
            cleaned_data: Data after cleaning
            validation_results: Results from data validation
            imputation_log: Log of imputation operations
            outlier_log: Log of outlier removal operations
            normalization_log: Log of normalization operations

        Returns:
            Formatted report string
        """
        report = []
        report.append("=" * 60)
        report.append("CSV DATA CLEANER AND VALIDATOR - QUALITY REPORT")
        report.append("=" * 60)
        report.append(f"Generated: {datetime.datetime.now().strftime('%Y-%m-%d %
H:%M:%S')}")
        report.append("")

        # Summary statistics
        report.append("SUMMARY STATISTICS:")
        report.append("-" * 20)
        report.append(f"Original rows: {len(original_data)}")
        report.append(f"Cleaned rows: {len(cleaned_data)}")
        report.append(f"Rows removed: {len(original_data) - len(cleaned_data)}")
        report.append("")

        # Validation issues
        if validation_results.get('issues'):
            report.append("VALIDATION ISSUES:")
            report.append("-" * 20)
            for issue in validation_results['issues']:
                report.append(f" {issue}")
            report.append("")
```

```python
        # Column types
        if validation_results.get('column_types'):
            report.append("COLUMN TYPES:")
            report.append("-" * 20)
            for column, col_type in validation_results['column_types'].items():
                report.append(f" {column}: {col_type}")
            report.append("")

        # Imputation summary
        if imputation_log:
            report.append("MISSING VALUE IMPUTATION:")
            report.append("-" * 30)
            imputation_summary = collections.defaultdict(int)
            for entry in imputation_log:
                method = entry.get('method', 'unknown')
                imputation_summary[method] += 1

            for method, count in imputation_summary.items():
                report.append(f" {method.title()} imputation: {count} values")
            report.append("")

        # Outlier summary
        if outlier_log:
            report.append("OUTLIER DETECTION:")
            report.append("-" * 20)
            outlier_summary = collections.defaultdict(int)
            for entry in outlier_log:
                method = entry.get('method', 'unknown')
                outlier_summary[method] += 1

            for method, count in outlier_summary.items():
                report.append(f" {method.upper()} method: {count} outliers detec
ted")
            report.append("")

        # Normalization summary
        if normalization_log:
            report.append("NORMALIZATION/STANDARDIZATION:")
            report.append("-" * 30)
            norm_summary = collections.defaultdict(int)
            for entry in normalization_log:
                method = entry.get('method', 'unknown')
                norm_summary[method] += 1

            for method, count in norm_summary.items():
                report.append(f" {method.upper()} method: {count} values process
ed")
            report.append("")

        # Detailed logs
        report.append("DETAILED TRANSFORMATION LOG:")
        report.append("-" * 30)

        all_operations = []
        all_operations.extend(imputation_log)
        all_operations.extend(outlier_log)
        all_operations.extend(normalization_log)

        # Sort by timestamp
        all_operations.sort(key=lambda x: x.get('timestamp', datetime.datetime.m
in))

        for operation in all_operations:
            timestamp = operation.get('timestamp', datetime.datetime.now())
            report.append(f"[{timestamp.strftime('%H:%M:%S')}] ")
```

```python
            if 'imputed_value' in operation:
                report.append(f"Imputed '{operation['column']}' at row {operatio
n['row']}: "
                              f"'{operation['original_value']}'  '{operation['imput
ed_value']}'")
            elif 'value' in operation and 'method' in operation:
                if operation['method'] == 'zscore':
                    report.append(f"Outlier detected in '{operation['column']}'
at row {operation['row']}: "
                                  f"value {operation['value']} (z-score: {operation
['z_score']:.2f})")
                else:
                    report.append(f"Outlier detected in '{operation['column']}'
at row {operation['row']}: "
                                  f"value {operation['value']}")
            elif 'normalized_value' in operation:
                report.append(f"Normalized '{operation['column']}' at row {opera
tion['row']}: "
                              f"{operation['original_value']:.2f}  {operation['norm
alized_value']:.4f}")
            elif 'standardized_value' in operation:
                report.append(f"Standardized '{operation['column']}' at row {ope
ration['row']}: "
                              f"{operation['original_value']:.2f}  {operation['stan
dardized_value']:.4f}")

        report.append("")
        report.append("=" * 60)

        return "\n".join(report)

    def save_report(self, report: str, output_path: str) -> bool:
        """
        Save the report to a file.

        Args:
            report: The report string to save
            output_path: Path where to save the report

        Returns:
            bool: True if successful, False otherwise
        """
        try:
            with open(output_path, 'w', encoding='utf-8') as file:
                file.write(report)
            return True
        except Exception as e:
            print(f"Error saving report: {e}")
            return False
```

# File: tests\\__init__.py

```python
"""
Tests package for CSV Data Cleaner and Validator (Group 22)

This package contains all unit tests for the CSV Data Cleaner application.
"""
```

# File: tests\test_csv_loader.py

```python
"""
Test cases for CSVLoader class.
"""

import unittest
import tempfile
import os
import csv
import sys
sys.path.append('src')

from src.data_processors.csv_loader import CSVLoader


class TestCSVLoader(unittest.TestCase):
    """Test cases for CSVLoader class."""

    def setUp(self):
        """Set up test fixtures."""
        self.loader = CSVLoader()

        # Create a temporary CSV file for testing
        self.temp_file = tempfile.NamedTemporaryFile(mode='w', suffix='.csv', de
lete=False)
        self.temp_file.write("Name,Age,Email,Salary\n")
        self.temp_file.write("Alice,29,alice@email.com,54000\n")
        self.temp_file.write("Bob,32,bob@email.com,52000\n")
        self.temp_file.write("Charlie,28,charlie@email.com,48000\n")
        self.temp_file.close()

    def tearDown(self):
        """Clean up test fixtures."""
        if os.path.exists(self.temp_file.name):
            os.unlink(self.temp_file.name)

    def test_load_csv_success(self):
        """Test successful CSV loading."""
        result = self.loader.load_csv(self.temp_file.name)

        self.assertTrue(result)
        self.assertEqual(len(self.loader.data), 3)
        self.assertEqual(self.loader.headers, ['Name', 'Age', 'Email', 'Salary']
)
        self.assertEqual(self.loader.file_path, self.temp_file.name)
        self.assertEqual(len(self.loader.issues), 0)

    def test_load_csv_file_not_found(self):
        """Test loading non-existent file."""
        result = self.loader.load_csv("nonexistent.csv")

        self.assertFalse(result)
        self.assertIn("File not found", self.loader.issues[0])

    def test_load_csv_wrong_extension(self):
        """Test loading file with wrong extension."""
        temp_txt_file = tempfile.NamedTemporaryFile(mode='w', suffix='.txt', del
ete=False)
        temp_txt_file.write("Name,Age\nAlice,29\n")
        temp_txt_file.close()

        result = self.loader.load_csv(temp_txt_file.name)

        self.assertFalse(result)
```

```python
        self.assertIn("File must have .csv extension", self.loader.issues[0])

        os.unlink(temp_txt_file.name)

    def test_get_data(self):
        """Test getting loaded data."""
        self.loader.load_csv(self.temp_file.name)
        data = self.loader.get_data()

        self.assertEqual(len(data), 3)
        self.assertEqual(data[0]['Name'], 'Alice')
        self.assertEqual(data[0]['Age'], '29')

    def test_get_headers(self):
        """Test getting headers."""
        self.loader.load_csv(self.temp_file.name)
        headers = self.loader.get_headers()

        self.assertEqual(headers, ['Name', 'Age', 'Email', 'Salary'])


if __name__ == '__main__':
    unittest.main()
```

# File: tests\test_data_validator.py

```python
"""
Test cases for DataValidator class.
"""

import unittest
import sys
sys.path.append('src')

from src.data_processors.data_validator import DataValidator


class TestDataValidator(unittest.TestCase):
    """Test cases for DataValidator class."""

    def setUp(self):
        """Set up test fixtures."""
        self.validator = DataValidator()

        # Test data with various types and issues
        self.test_data = [
            {'Name': 'Alice', 'Age': '29', 'Email': 'alice@email.com', 'Salary':
 '54000'},
            {'Name': 'Bob', 'Age': '32', 'Email': 'bob@email.com', 'Salary': '52
000'},
            {'Name': 'Charlie', 'Age': '28', 'Email': 'charlie@email.com', 'Sala
ry': '48000'},
            {'Name': 'Diana', 'Age': 'abc', 'Email': 'diana@email.com', 'Salary'
: '60000'},  # Invalid age
            {'Name': 'Eve', 'Age': '30', 'Email': 'not-an-email', 'Salary': '550
00'},  # Invalid email
            {'Name': 'Frank', 'Age': '35', 'Email': 'frank@email.com', 'Salary':
 '70000'}
        ]
        self.headers = ['Name', 'Age', 'Email', 'Salary']

    def test_validate_data_success(self):
        """Test successful data validation."""
        result = self.validator.validate_data(self.test_data, self.headers)

        self.assertIsInstance(result, dict)
        self.assertIn('valid', result)
        self.assertIn('issues', result)
        self.assertIn('column_types', result)

    def test_validate_data_empty_data(self):
        """Test validation with empty data."""
        result = self.validator.validate_data([], self.headers)

        self.assertFalse(result['valid'])
        self.assertIn("No data or headers to validate", result['issues'])

    def test_validate_data_numeric_column(self):
        """Test validation of numeric column."""
        numeric_data = [
            {'Age': '29', 'Salary': '54000'},
            {'Age': '32', 'Salary': '52000'},
            {'Age': '28', 'Salary': '48000'}
        ]
        headers = ['Age', 'Salary']

        result = self.validator.validate_data(numeric_data, headers)

        self.assertEqual(result['column_types']['Age'], 'numeric')
```

```python
            self.assertEqual(result['column_types']['Salary'], 'numeric')

    def test_validate_data_categorical_column(self):
        """Test validation of categorical column."""
        categorical_data = [
            {'Name': 'Alice', 'Department': 'Engineering'},
            {'Name': 'Bob', 'Department': 'Marketing'},
            {'Name': 'Charlie', 'Department': 'Sales'}
        ]
        headers = ['Name', 'Department']

        result = self.validator.validate_data(categorical_data, headers)

        self.assertEqual(result['column_types']['Name'], 'categorical')
        self.assertEqual(result['column_types']['Department'], 'categorical')

    def test_is_valid_date(self):
        """Test date validation method."""
        valid_dates = ['2023-04-01', '04/01/2023', '04-01-2023']
        invalid_dates = ['2023/04/01', '01-04-2023', 'not-a-date']

        for date in valid_dates:
            self.assertTrue(self.validator._is_valid_date(date))

        for date in invalid_dates:
            self.assertFalse(self.validator._is_valid_date(date))

    def test_is_valid_email(self):
        """Test email validation method."""
        valid_emails = ['alice@email.com', 'bob@company.co.uk', 'test.user@domai
n.org']
        invalid_emails = ['not-an-email', 'alice@', '@email.com', 'alice.email.c
om']

        for email in valid_emails:
            self.assertTrue(self.validator._is_valid_email(email))

        for email in invalid_emails:
            self.assertFalse(self.validator._is_valid_email(email))


if __name__ == '__main__':
    unittest.main()
```