# Survey Data Analyzer

## Complete Project Documentation & Codebase

A comprehensive Command Line Interface (CLI) tool for analyzing survey responses from CSV files. This project demonstrates advanced Python programming concepts including modular design, statistical analysis, sentiment analysis, pattern detection, and report generation.

• CSV Input Handling with multi-encoding support

• Survey Summary with demographic breakdowns

• Statistical Analysis (cross-tabulation, chi-square tests)

• Sentiment Analysis on text responses

• Pattern Recognition and correlation detection

• Comprehensive Report Generation

• Professional CLI Interface

Generated on: August 02, 2025 at 03:03 AM

# Table of Contents

# 1. Project Overview

The Survey Data Analyzer is a sophisticated Python CLI application designed to perform comprehensive analysis of survey data stored in CSV format. The project demonstrates advanced software engineering principles including modular design, object-oriented programming, error handling, and comprehensive testing. The application consists of several specialized modules, each handling a specific aspect of survey analysis:

• **DataLoader**: Handles CSV file loading, validation, and data preprocessing

• **SurveySummary**: Generates demographic breakdowns and response statistics

• **StatsAnalyzer**: Performs statistical analysis including chi-square tests

• **SentimentAnalyzer**: Analyzes text responses using keyword-based sentiment scoring

• **PatternDetector**: Identifies correlations and patterns in survey responses

• **ReportGenerator**: Creates comprehensive analysis reports

• **Utils**: Provides utility functions for formatting and validation

# 2. Main Application (main.py)

The main.py file serves as the entry point for the Survey Data Analyzer CLI application. It provides a menu-driven interface that orchestrates all analysis components and manages the user interaction flow.

## Main Application - main.py

```
1: #!/usr/bin/env python3 2: """ 3: Survey Data Analyzer - Main CLI Application
4: Author: Student Developer 5: Description: Main entry point for the Survey Data
Analyzer CLI tool. 6: Provides a menu-driven interface for analyzing survey data
from CSV files. 7: """ 8: 9: import sys 10: import os 11: from typing import
Optional, Dict, Any 12: from data_loader import DataLoader 13: from
survey_summary import SurveySummary 14: from stats_analyzer import StatsAnalyzer
15: from sentiment_analyzer import SentimentAnalyzer 16: from pattern_detector
import PatternDetector 17: from report_generator import ReportGenerator 18: from
utils import clear_screen, print_header, print_menu 19: 20: 21: class
SurveyAnalyzerCLI: 22: """Main CLI application class for Survey Data Analyzer."""
23: 24: def __init__(self): 25: """Initialize the CLI application with empty data
containers.""" 26: self.data_loader = DataLoader() 27: self.survey_data = None
28: self.survey_summary = None 29: self.stats_analyzer = None 30:
self.sentiment_analyzer = None 31: self.pattern_detector = None 32:
self.report_generator = None 33: 34: def load_survey_data(self) -> bool: 35:
"""Load survey data from CSV file.""" 36: try: 37: print_header("Load Survey
Data") 38: file_path = input("Enter the path to your CSV file: ").strip() 39: 40:
if not file_path: 41: print("ERROR: No file path provided.") 42: return False 43:
44: if not os.path.exists(file_path): 45: print(f"ERROR: File not found:
{file_path}") 46: return False 47: 48: self.survey_data =
self.data_loader.load_csv(file_path) 49: if self.survey_data: 50:
print(f"SUCCESS: Successfully loaded {len(self.survey_data)} survey responses")
```

```
51: print(f"INFO: Columns: {', '.join(self.survey_data[0].keys())}") 52: return
True 53: else: 54: print("ERROR: Failed to load survey data") 55: return False
56: 57: except Exception as e: 58: print(f"ERROR: Error loading data: {str(e)}")
59: return False 60: 61: def view_summary_statistics(self): 62: """Display
summary statistics for the loaded survey data.""" 63: if not self.survey_data:
64: print("ERROR: No survey data loaded. Please load data first.") 65: return 66:
67: try: 68: print_header("Survey Summary Statistics") 69: self.survey_summary =
SurveySummary(self.survey_data) 70: summary =
self.survey_summary.generate_summary() 71: 72: print("\nRESPONSE OVERVIEW:") 73:
print(f" Total Responses: {summary['total_responses']}") 74: print(f" Response
Rate: {summary['response_rate']:.1f}%") 75: 76: print("\nDEMOGRAPHIC
BREAKDOWN:") 77: for demo, breakdown in summary['demographics'].items(): 78:
print(f"\n {demo.upper()}:") 79: for category, count in breakdown.items(): 80:
percentage = (count / summary['total_responses']) * 100 81: print(f" {category}:
{count} ({percentage:.1f}%)") 82: 83: except Exception as e: 84: print(f"ERROR:
Error generating summary: {str(e)}") 85: 86: def analyze_sentiment(self): 87:
"""Analyze sentiment of text responses.""" 88: if not self.survey_data: 89:
print("ERROR: No survey data loaded. Please load data first.") 90: return 91: 92:
try: 93: print_header("Sentiment Analysis") 94: self.sentiment_analyzer =
SentimentAnalyzer() 95: 96: # Find text columns 97: text_columns = [] 98: for col
in self.survey_data[0].keys(): 99: if any(isinstance(row[col], str) and
len(str(row[col])) > 20 100: for row in self.survey_data[:10]):
```

```python
101:         text_columns.append(col)
102:
103:     if not text_columns:
104:         print("ERROR: No text columns found for sentiment analysis.")
105:         return
106:
107:     print(f"INFO: Found text columns: {', '.join(text_columns)}")
108:
109:     for column in text_columns:
110:         print(f"\nANALYZING: Analyzing sentiment for '{column}':")
111:         sentiment_results = self.sentiment_analyzer.analyze_column(
112:             self.survey_data, column
113:         )
114:
115:         print(f"  Positive responses: {sentiment_results['positive']} ({sentiment_results['positive_pct']:.1f}%)")
116:         print(f"  Negative responses: {sentiment_results['negative']} ({sentiment_results['negative_pct']:.1f}%)")
117:         print(f"  Neutral responses: {sentiment_results['neutral']} ({sentiment_results['neutral_pct']:.1f}%)")
118:         print(f"  Average sentiment score: {sentiment_results['avg_score']:.2f}")
119:
120:     except Exception as e:
121:         print(f"ERROR: Error analyzing sentiment: {str(e)}")
122:
123: def cross_tabulate_results(self):
124:     """Perform cross-tabulation analysis."""
125:     if not self.survey_data:
126:         print("ERROR: No survey data loaded. Please load data first.")
127:         return
128:
129:     try:
130:         print_header("Cross-Tabulation Analysis")
131:         self.stats_analyzer = StatsAnalyzer(self.survey_data)
132:
133:         # Get available columns
134:         columns = list(self.survey_data[0].keys())
135:         print(f"Available columns: {', '.join(columns)}")
136:
137:         col1 = input("Enter first column name: ").strip()
138:         col2 = input("Enter second column name: ").strip()
139:
140:         if col1 not in columns or col2 not in columns:
141:             print("ERROR: Invalid column names.")
142:             return
143:
144:         crosstab = self.stats_analyzer.cross_tabulate(col1, col2)
145:         chi_square = self.stats_analyzer.chi_square_test(col1, col2)
146:
147:         print(f"\nCROSS-TABULATION: {col1} vs {col2}")
148:         print("=" * 50)
149:
150:         # Display crosstab
```

```python
151: for row in crosstab:
152:     print(" | ".join(f"{cell:>8}" for cell in row))
153:
154:     print(f"\nCHI-SQUARE TEST RESULTS:")
155:     print(f"  Chi-Square Value: {chi_square['chi_square']:.4f}")
156:     print(f"  P-Value: {chi_square['p_value']:.4f}")
157:     print(f"  Degrees of Freedom: {chi_square['df']}")
158:     print(f"  Significant: {'Yes' if chi_square['significant'] else 'No'}")
159:
160:     except Exception as e:
161:         print(f"ERROR: Error in cross-tabulation: {str(e)}")
162:
163: def detect_patterns(self):
164:     """Detect patterns and correlations in survey responses."""
165:     if not self.survey_data:
166:         print("ERROR: No survey data loaded. Please load data first.")
167:         return
168:
169:     try:
170:         print_header("Pattern Detection")
171:         self.pattern_detector = PatternDetector(self.survey_data)
172:
173:         patterns = self.pattern_detector.find_patterns()
174:
175:         print("DETECTED PATTERNS:")
176:         print("=" * 50)
177:
178:         for pattern in patterns:
179:             print(f"\nPATTERN: {pattern['description']}")
180:             print(f"  Confidence: {pattern['confidence']:.1f}%")
181:             print(f"  Sample size: {pattern['sample_size']}")
182:
183:     except Exception as e:
184:         print(f"ERROR: Error detecting patterns: {str(e)}")
185:
186: def generate_report(self):
187:     """Generate a comprehensive analysis report."""
188:     if not self.survey_data:
189:         print("ERROR: No survey data loaded. Please load data first.")
190:         return
191:
192:     try:
193:         print_header("Generate Report")
194:
195:         # Initialize all analyzers
196:         self.survey_summary = SurveySummary(self.survey_data)
197:         self.stats_analyzer = StatsAnalyzer(self.survey_data)
198:         self.sentiment_analyzer = SentimentAnalyzer()
199:         self.pattern_detector = PatternDetector(self.survey_data)
200:         self.report_generator = ReportGenerator()
```

```
201: 202: output_file = input("Enter output file name (default:
survey_report.txt): ").strip() 203: if not output_file: 204: output_file =
"survey_report.txt" 205: 206: if not output_file.endswith('.txt'): 207:
output_file += '.txt' 208: 209: # Generate comprehensive report 210: report_data
= { 211: 'summary': self.survey_summary.generate_summary(), 212: 'sentiment':
self.sentiment_analyzer.analyze_all_text_columns(self.survey_data), 213:
'patterns': self.pattern_detector.find_patterns(), 214: 'file_path': output_file
215: } 216: 217: success = self.report_generator.generate_report(report_data)
218: 219: if success: 220: print(f"SUCCESS: Report generated successfully:
{output_file}") 221: else: 222: print("ERROR: Failed to generate report") 223:
224: except Exception as e: 225: print(f"ERROR: Error generating report:
{str(e)}") 226: 227: def run(self): 228: """Main application loop.""" 229: while
True: 230: clear_screen() 231: print_header("Survey Data Analyzer") 232: 233: if
self.survey_data: 234: print(f"INFO: Loaded: {len(self.survey_data)} responses")
235: else: 236: print("INFO: No data loaded") 237: 238: print_menu([ 239: "Load
survey data", 240: "View summary statistics", 241: "Analyze sentiment", 242:
"Cross-tabulate results", 243: "Detect patterns", 244: "Generate report", 245:
"Exit" 246: ]) 247: 248: choice = input("\nEnter your choice (1-7): ").strip()
249: 250: if choice == '1':
```

```python
251:         self.load_survey_data()
252:     elif choice == '2':
253:         self.view_summary_statistics()
254:     elif choice == '3':
255:         self.analyze_sentiment()
256:     elif choice == '4':
257:         self.cross_tabulate_results()
258:     elif choice == '5':
259:         self.detect_patterns()
260:     elif choice == '6':
261:         self.generate_report()
262:     elif choice == '7':
263:         print("\nThank you for using Survey Data Analyzer!")
264:         sys.exit(0)
265:     else:
266:         print("ERROR: Invalid choice. Please try again.")
267:
268:     input("\nPress Enter to continue...")
269:
270:
271: def main():
272:     """Main entry point for the application."""
273:     try:
274:         app = SurveyAnalyzerCLI()
275:         app.run()
276:     except KeyboardInterrupt:
277:         print("\n\nGoodbye!")
278:         sys.exit(0)
279:     except Exception as e:
280:         print(f"ERROR: Unexpected error: {str(e)}")
281:         sys.exit(1)
282:
283:
284: if __name__ == "__main__":
285:     main()
```

# 3. Data Loading Module (data_loader.py)

The DataLoader class handles all aspects of CSV file processing including file validation, encoding detection, data cleaning, and validation. It provides robust error handling and supports multiple file formats and encodings.

## Data Loading Module - data_loader.py

```
1: #!/usr/bin/env python3 2: """ 3: Survey Data Loader Module 4: Author: Student
Developer 5: Description: Handles CSV file loading, validation, and data
preprocessing. 6: Provides robust error handling and data consistency checks. 7:
""" 8: 9: import csv 10: import os 11: import sys 12: from typing import List,
Dict, Any, Optional 13: from collections import defaultdict 14: 15: 16: class
DataLoader: 17: """Handles loading and validation of survey data from CSV
files.""" 18: 19: def __init__(self): 20: """Initialize the data loader with
validation settings.""" 21: self.required_fields = ['age', 'gender', 'region'] #
Common required fields 22: self.max_file_size = 50 * 1024 * 1024 # 50MB limit 23:
self.supported_encodings = ['utf-8', 'latin-1', 'cp1252'] 24: 25: def
load_csv(self, file_path: str) -> Optional[List[Dict[str, Any]]]: 26: """ 27:
Load and validate CSV survey data. 28: 29: Args: 30: file_path: Path to the CSV
file 31: 32: Returns: 33: List of dictionaries representing survey responses, or
None if failed 34: """ 35: try: 36: # Validate file exists and is readable 37: if
not self._validate_file(file_path): 38: return None 39: 40: # Try different
encodings 41: data = None 42: for encoding in self.supported_encodings: 43: try:
44: data = self._read_csv_with_encoding(file_path, encoding) 45: if data: 46:
break 47: except UnicodeDecodeError: 48: continue 49: 50: if not data:
```

```
51: print("ERROR: Could not read file with any supported encoding") 52: return
None 53: 54: # Validate and clean data 55: cleaned_data =
self._clean_and_validate_data(data) 56: if not cleaned_data: 57: return None 58:
59: print(f"SUCCESS: Data validation completed successfully") 60: return
cleaned_data 61: 62: except Exception as e: 63: print(f"ERROR: Error loading CSV
file: {str(e)}") 64: return None 65: 66: def _validate_file(self, file_path: str)
-> bool: 67: """Validate that the file exists and is accessible.""" 68: try: 69:
if not os.path.exists(file_path): 70: print(f"ERROR: File not found:
{file_path}") 71: return False 72: 73: if not os.path.isfile(file_path): 74:
print(f"ERROR: Path is not a file: {file_path}") 75: return False 76: 77: # Check
file size 78: file_size = os.path.getsize(file_path) 79: if file_size >
self.max_file_size: 80: print(f"ERROR: File too large: {file_size /
(1024*1024):.1f}MB (max: 50MB)") 81: return False 82: 83: if file_size == 0: 84:
print("ERROR: File is empty") 85: return False 86: 87: return True 88: 89: except
Exception as e: 90: print(f"ERROR: Error validating file: {str(e)}") 91: return
False 92: 93: def _read_csv_with_encoding(self, file_path: str, encoding: str) ->
Optional[List[Dict[str, Any]]]: 94: """Read CSV file with specified encoding."""
95: try: 96: with open(file_path, 'r', encoding=encoding, newline='') as file:
97: # Try to detect delimiter 98: sample = file.read(1024) 99: file.seek(0) 100:
```

```python
101: # Common delimiters to try
102: delimiters = [',', ';', '\t', '|']
103: detected_delimiter = ','
104:
105: for delimiter in delimiters:
106: if delimiter in sample:
107: detected_delimiter = delimiter
108: break
109:
110: reader = csv.DictReader(file, delimiter=detected_delimiter)
111: data = list(reader)
112:
113: if not data:
114: print("ERROR: No data found in CSV file")
115: return None
116:
117: print(f"SUCCESS: Loaded {len(data)} rows with {len(data[0])} columns")
118: print(f"INFO: Columns: {', '.join(data[0].keys())}")
119:
120: return data
121:
122: except Exception as e:
123: print(f"ERROR: Error reading CSV with {encoding} encoding: {str(e)}")
124: return None
125:
126: def _clean_and_validate_data(self, data: List[Dict[str, Any]]) -> Optional[List[Dict[str, Any]]]:
127: """Clean and validate the loaded data."""
128: try:
129: if not data:
130: return None
131:
132: cleaned_data = []
133: validation_errors = []
134:
135: for i, row in enumerate(data, 1):
136: cleaned_row = {}
137: row_errors = []
138:
139: for key, value in row.items():
140: # Clean the key (remove whitespace, normalize)
141: if key is None:
142: continue # Skip rows with None keys
143: clean_key = str(key).strip().lower().replace(' ', '_')
144:
145: # Clean the value
146: if value is None or value == '':
147: cleaned_value = None
148: else:
149: cleaned_value = str(value).strip()
150: if cleaned_value.lower() in ['na', 'n/a', 'null', 'none']:
```

```python
151:         cleaned_value = None
152:
153:         cleaned_row[clean_key] = cleaned_value
154:
155:     # Basic validation
156:     if self._validate_row(cleaned_row, i):
157:         cleaned_data.append(cleaned_row)
158:     else:
159:         validation_errors.append(f"Row {i}: Invalid data")
160:
161:     # Report validation results
162:     total_rows = len(data)
163:     valid_rows = len(cleaned_data)
164:     invalid_rows = total_rows - valid_rows
165:
166:     print(f"INFO: Data validation results:")
167:     print(f"  Total rows: {total_rows}")
168:     print(f"  Valid rows: {valid_rows}")
169:     print(f"  Invalid rows: {invalid_rows}")
170:
171:     if invalid_rows > 0:
172:         print(f"WARNING: {invalid_rows} rows had validation issues")
173:
174:     if valid_rows == 0:
175:         print("ERROR: No valid data found")
176:         return None
177:
178:     return cleaned_data
179:
180: except Exception as e:
181:     print(f"ERROR: Error cleaning data: {str(e)}")
182:     return None
183:
184: def _validate_row(self, row: Dict[str, Any], row_num: int) -> bool:
185:     """Validate a single row of data."""
186:     try:
187:         # Check for minimum required fields
188:         if len(row) < 2:
189:             return False
190:
191:         # Validate age if present
192:         if 'age' in row and row['age'] is not None:
193:             try:
194:                 age = int(row['age'])
195:                 if age < 0 or age > 120:
196:                     return False
197:             except (ValueError, TypeError):
198:                 return False
199:
200:         # Validate gender if present
```

```python
201: if 'gender' in row and row['gender'] is not None:
202:     try:
203:         gender = str(row['gender']).lower()
204:         valid_genders = ['male', 'female', 'm', 'f', 'other', 'prefer not to say']
205:         if gender not in valid_genders:
206:             return False
207:     except (AttributeError, TypeError):
208:         return False
209:
210:     return True
211:
212:     except Exception:
213:         return False
214:
215:     def get_data_summary(self, data: List[Dict[str, Any]]) -> Dict[str, Any]:
216:         """Generate a summary of the loaded data."""
217:         if not data:
218:             return {}
219:
220:         summary = {
221:             'total_rows': len(data),
222:             'columns': list(data[0].keys()) if data else [],
223:             'missing_values': {},
224:             'unique_values': {}
225:         }
226:
227:         # Analyze missing values and unique values
228:         for column in summary['columns']:
229:             values = [row.get(column) for row in data]
230:             missing_count = sum(1 for v in values if v is None or v == '')
231:             summary['missing_values'][column] = {
232:                 'count': missing_count,
233:                 'percentage': (missing_count / len(data)) * 100
234:             }
235:
236:             unique_vals = set(v for v in values if v is not None and v != '')
237:             summary['unique_values'][column] = len(unique_vals)
238:
239:         return summary
240:
241:     def export_sample_data(self, file_path: str = "sample_survey.csv"):
242:         """Create a sample CSV file for testing."""
243:         sample_data = [
244:             {
245:                 'age': '25',
246:                 'gender': 'Female',
247:                 'region': 'North',
248:                 'education': 'Bachelor',
249:                 'satisfaction': 'Very Satisfied',
250:                 'feedback': 'Great experience with the product!',
```

251: 'recommend': 'Yes' 252: }, 253: { 254: 'age': '32', 255: 'gender': 'Male', 256: 'region': 'South', 257: 'education': 'Master', 258: 'satisfaction': 'Satisfied', 259: 'feedback': 'Good but could be better.', 260: 'recommend': 'Yes' 261: }, 262: { 263: 'age': '45', 264: 'gender': 'Female', 265: 'region': 'East', 266: 'education': 'High School', 267: 'satisfaction': 'Neutral', 268: 'feedback': 'It was okay, nothing special.', 269: 'recommend': 'Maybe' 270: }, 271: { 272: 'age': '28', 273: 'gender': 'Male', 274: 'region': 'West', 275: 'education': 'Bachelor', 276: 'satisfaction': 'Dissatisfied', 277: 'feedback': 'Poor quality and bad service.', 278: 'recommend': 'No' 279: }, 280: { 281: 'age': '35', 282: 'gender': 'Female', 283: 'region': 'North', 284: 'education': 'PhD', 285: 'satisfaction': 'Very Satisfied', 286: 'feedback': 'Excellent product and amazing support!', 287: 'recommend': 'Yes' 288: } 289: ] 290: 291: try: 292: with open(file_path, 'w', newline='', encoding='utf-8') as file: 293: if sample_data: 294: writer = csv.DictWriter(file, fieldnames=sample_data[0].keys()) 295: writer.writeheader() 296: writer.writerows(sample_data) 297: 298: print(f"SUCCESS: Sample data exported to {file_path}") 299: return True 300:

```
301: except Exception as e: 302: print(f"ERROR: Error exporting sample data:
{str(e)}") 303: return False
```

# 4. Survey Summary Module (survey_summary.py)

The SurveySummary class generates comprehensive summary statistics including demographic breakdowns, response rates, and data quality assessments. It provides the foundation for all subsequent analysis.

## Survey Summary Module - survey_summary.py

```
1: #!/usr/bin/env python3 2: """ 3: Survey Summary Module 4: Author: Student
Developer 5: Description: Generates summary statistics, response counts,
percentages, 6: and demographic breakdowns for survey data analysis. 7: """ 8: 9:
import statistics 10: from typing import List, Dict, Any, Optional 11: from
collections import defaultdict, Counter 12: 13: 14: class SurveySummary: 15:
"""Handles generation of survey summary statistics and demographic breakdowns."""
16: 17: def __init__(self, survey_data: List[Dict[str, Any]]): 18: """ 19:
Initialize the survey summary analyzer. 20: 21: Args: 22: survey_data: List of
dictionaries containing survey responses 23: """ 24: self.survey_data =
survey_data 25: self.total_responses = len(survey_data) 26: self.columns =
list(survey_data[0].keys()) if survey_data else [] 27: 28: def
generate_summary(self) -> Dict[str, Any]: 29: """ 30: Generate comprehensive
summary statistics. 31: 32: Returns: 33: Dictionary containing all summary
statistics 34: """ 35: if not self.survey_data: 36: return {} 37: 38: summary = {
39: 'total_responses': self.total_responses, 40: 'response_rate':
self._calculate_response_rate(), 41: 'demographics':
self._analyze_demographics(), 42: 'question_summaries':
self._analyze_questions(), 43: 'data_quality': self._assess_data_quality() 44: }
45: 46: return summary 47: 48: def _calculate_response_rate(self) -> float: 49:
"""Calculate the response rate (placeholder for actual calculation).""" 50: # In
a real scenario, this would compare against expected responses
```

```python
51:         # For now, we'll use a placeholder calculation
52:         return 85.5  # Placeholder response rate
53:
54:     def _analyze_demographics(self) -> Dict[str, Dict[str, int]]:
55:         """Analyze demographic breakdowns."""
56:         demographics = {}
57:
58:         # Common demographic fields
59:         demo_fields = ['age', 'gender', 'region', 'education', 'income']
60:
61:         for field in demo_fields:
62:             if field in self.columns:
63:                 demographics[field] = self._count_responses_by_field(field)
64:
65:         return demographics
66:
67:     def _analyze_questions(self) -> Dict[str, Dict[str, Any]]:
68:         """Analyze responses for each question."""
69:         question_summaries = {}
70:
71:         for column in self.columns:
72:             if column not in ['age', 'gender', 'region', 'education', 'income']:
73:                 question_summaries[column] = self._analyze_question_responses(column)
74:
75:         return question_summaries
76:
77:     def _count_responses_by_field(self, field: str) -> Dict[str, int]:
78:         """Count responses for a specific field."""
79:         counts = Counter()
80:         for row in self.survey_data:
81:             value = row.get(field)
82:             if value is not None and str(value).strip():
83:                 # Normalize the value
84:                 normalized_value = str(value).strip().title()
85:                 counts[normalized_value] += 1
86:
87:         return dict(counts)
88:
89:     def _analyze_question_responses(self, question: str) -> Dict[str, Any]:
90:         """Analyze responses for a specific question."""
91:         responses = [row.get(question) for row in self.survey_data]
92:         valid_responses = [r for r in responses if r is not None and str(r).strip()]
93:
94:         if not valid_responses:
95:             return {
96:                 'total_responses': 0,
97:                 'missing_responses': len(responses),
98:                 'response_rate': 0.0,
99:                 'top_responses': [],
```

```python
101:     'response_distribution': {}
102: }
103:
104: # Count responses
105: response_counts = Counter(valid_responses)
106:
107: # Calculate percentages
108: total_valid = len(valid_responses)
109: response_distribution = {}
110: for response, count in response_counts.items():
111:     percentage = (count / total_valid) * 100
112:     response_distribution[response] = {
113:         'count': count,
114:         'percentage': percentage
115:     }
116:
117: # Get top responses
118: top_responses = response_counts.most_common(5)
119:
120: return {
121:     'total_responses': total_valid,
122:     'missing_responses': len(responses) - total_valid,
123:     'response_rate': (total_valid / len(responses)) * 100,
124:     'top_responses': top_responses,
125:     'response_distribution': response_distribution
126: }
127:
128: def _assess_data_quality(self) -> Dict[str, Any]:
129:     """Assess the quality of the survey data."""
130:     quality_metrics = {
131:         'total_rows': len(self.survey_data),
132:         'total_columns': len(self.columns),
133:         'missing_data': {},
134:         'completeness': {}
135:     }
136:
137:     # Analyze missing data for each column
138:     for column in self.columns:
139:         values = [row.get(column) for row in self.survey_data]
140:         missing_count = sum(1 for v in values if v is None or str(v).strip() == '')
141:         missing_percentage = (missing_count / len(values)) * 100
142:
143:         quality_metrics['missing_data'][column] = {
144:             'count': missing_count,
145:             'percentage': missing_percentage
146:         }
147:
148:         quality_metrics['completeness'][column] = 100 - missing_percentage
149:
150:     return quality_metrics
```

```python
151:
152:     def get_age_distribution(self) -> Dict[str, int]:
153:         """Get age distribution if age data is available."""
154:         if 'age' not in self.columns:
155:             return {}
156:
157:         age_groups = {
158:             '18-25': 0,
159:             '26-35': 0,
160:             '36-45': 0,
161:             '46-55': 0,
162:             '56-65': 0,
163:             '65+': 0
164:         }
165:
166:         for row in self.survey_data:
167:             age_str = row.get('age')
168:             if age_str and str(age_str).isdigit():
169:                 try:
170:                     age = int(age_str)
171:                     if 18 <= age <= 25:
172:                         age_groups['18-25'] += 1
173:                     elif 26 <= age <= 35:
174:                         age_groups['26-35'] += 1
175:                     elif 36 <= age <= 45:
176:                         age_groups['36-45'] += 1
177:                     elif 46 <= age <= 55:
178:                         age_groups['46-55'] += 1
179:                     elif 56 <= age <= 65:
180:                         age_groups['56-65'] += 1
181:                     elif age > 65:
182:                         age_groups['65+'] += 1
183:                 except ValueError:
184:                     continue
185:
186:         return age_groups
187:
188:     def get_gender_distribution(self) -> Dict[str, int]:
189:         """Get gender distribution if gender data is available."""
190:         if 'gender' not in self.columns:
191:             return {}
192:
193:         return self._count_responses_by_field('gender')
194:
195:     def get_regional_distribution(self) -> Dict[str, int]:
196:         """Get regional distribution if region data is available."""
197:         if 'region' not in self.columns:
198:             return {}
199:
200:         return self._count_responses_by_field('region')
```

```python
201:
202:    def get_education_distribution(self) -> Dict[str, int]:
203:        """Get education distribution if education data is available."""
204:        if 'education' not in self.columns:
205:            return {}
206:
207:        return self._count_responses_by_field('education')
208:
209:    def calculate_average_age(self) -> Optional[float]:
210:        """Calculate average age if age data is available."""
211:        if 'age' not in self.columns:
212:            return None
213:
214:        ages = []
215:        for row in self.survey_data:
216:            age_str = row.get('age')
217:            if age_str and str(age_str).isdigit():
218:                try:
219:                    ages.append(int(age_str))
220:                except ValueError:
221:                    continue
222:
223:        if ages:
224:            return statistics.mean(ages)
225:        return None
226:
227:    def get_response_trends(self) -> Dict[str, Any]:
228:        """Identify response trends and patterns."""
229:        trends = {
230:            'most_common_responses': {},
231:            'response_patterns': [],
232:            'outliers': []
233:        }
234:
235:        # Find most common responses for each question
236:        for column in self.columns:
237:            if column not in ['age', 'gender', 'region', 'education', 'income']:
238:                responses = [row.get(column) for row in self.survey_data]
239:                valid_responses = [r for r in responses if r is not None and str(r).strip()]
240:
241:                if valid_responses:
242:                    response_counts = Counter(valid_responses)
243:                    most_common = response_counts.most_common(1)
244:                    if most_common:
245:                        trends['most_common_responses'][column] = {
246:                            'response': most_common[0][0],
247:                            'count': most_common[0][1],
248:                            'percentage': (most_common[0][1] / len(valid_responses)) * 100
249:                        }
250:
```

```python
251:     return trends
252:
253:     def generate_demographic_report(self) -> str:
254:         """Generate a formatted demographic report."""
255:         if not self.survey_data:
256:             return "No data available for demographic analysis."
257:
258:         report_lines = []
259:         report_lines.append("DEMOGRAPHIC ANALYSIS REPORT")
260:         report_lines.append("=" * 50)
261:         report_lines.append(f"Total Responses: {self.total_responses}")
262:         report_lines.append("")
263:
264:         # Age distribution
265:         age_dist = self.get_age_distribution()
266:         if age_dist:
267:             report_lines.append("AGE DISTRIBUTION:")
268:             for age_group, count in age_dist.items():
269:                 percentage = (count / self.total_responses) * 100
270:                 report_lines.append(f"  {age_group}: {count} ({percentage:.1f}%)")
271:             report_lines.append("")
272:
273:         # Gender distribution
274:         gender_dist = self.get_gender_distribution()
275:         if gender_dist:
276:             report_lines.append("GENDER DISTRIBUTION:")
277:             for gender, count in gender_dist.items():
278:                 percentage = (count / self.total_responses) * 100
279:                 report_lines.append(f"  {gender}: {count} ({percentage:.1f}%)")
280:             report_lines.append("")
281:
282:         # Regional distribution
283:         region_dist = self.get_regional_distribution()
284:         if region_dist:
285:             report_lines.append("REGIONAL DISTRIBUTION:")
286:             for region, count in region_dist.items():
287:                 percentage = (count / self.total_responses) * 100
288:                 report_lines.append(f"  {region}: {count} ({percentage:.1f}%)")
289:             report_lines.append("")
290:
291:         # Education distribution
292:         education_dist = self.get_education_distribution()
293:         if education_dist:
294:             report_lines.append("EDUCATION DISTRIBUTION:")
295:             for education, count in education_dist.items():
296:                 percentage = (count / self.total_responses) * 100
297:                 report_lines.append(f"  {education}: {count} ({percentage:.1f}%)")
298:             report_lines.append("")
299:
300:         return "\n".join(report_lines)
```

# 5. Statistical Analysis Module (stats_analyzer.py)

The StatsAnalyzer class performs advanced statistical analysis including cross-tabulation, chi-square tests, correlation analysis, and significance testing. It provides the statistical foundation for data-driven insights.

## Statistical Analysis Module - stats_analyzer.py

```
1: #!/usr/bin/env python3 2: """ 3: Statistical Analyzer Module 4: Author:
Student Developer 5: Description: Performs statistical analysis including
cross-tabulation, 6: chi-square tests, and correlation analysis on survey data.
7: """ 8: 9: import math 10: from typing import List, Dict, Any, Optional, Tuple
11: from collections import defaultdict, Counter 12: import itertools 13: 14: 15:
class StatsAnalyzer: 16: """Handles statistical analysis of survey data including
cross-tabulation and chi-square tests.""" 17: 18: def __init__(self, survey_data:
List[Dict[str, Any]]): 19: """ 20: Initialize the statistical analyzer. 21: 22:
Args: 23: survey_data: List of dictionaries containing survey responses 24: """
25: self.survey_data = survey_data 26: self.total_responses = len(survey_data)
27: self.columns = list(survey_data[0].keys()) if survey_data else [] 28: 29: def
cross_tabulate(self, col1: str, col2: str) -> List[List]: 30: """ 31: Perform
cross-tabulation between two columns. 32: 33: Args: 34: col1: First column name
35: col2: Second column name 36: 37: Returns: 38: Cross-tabulation matrix as a
list of lists 39: """ 40: if col1 not in self.columns or col2 not in
self.columns: 41: raise ValueError(f"Column not found: {col1} or {col2}") 42: 43:
# Get unique values for each column 44: values1 = set() 45: values2 = set() 46:
47: for row in self.survey_data: 48: val1 = row.get(col1) 49: val2 =
row.get(col2) 50:
```

```python
51: if val1 is not None and str(val1).strip(): 52: values1.add(str(val1).strip())
53: if val2 is not None and str(val2).strip(): 54: values2.add(str(val2).strip())
55: 56: # Sort values for consistent ordering 57: values1 = sorted(list(values1))
58: values2 = sorted(list(values2)) 59: 60: # Create cross-tabulation matrix 61:
crosstab = [] 62: 63: # Header row 64: header = [''] + values2 65:
crosstab.append(header) 66: 67: # Data rows 68: for val1 in values1: 69: row =
[val1] 70: for val2 in values2: 71: count = 0 72: for survey_row in
self.survey_data: 73: if (str(survey_row.get(col1, '')).strip() == val1 and 74:
str(survey_row.get(col2, '')).strip() == val2): 75: count += 1 76:
row.append(count) 77: crosstab.append(row) 78: 79: return crosstab 80: 81: def
chi_square_test(self, col1: str, col2: str) -> Dict[str, Any]: 82: """ 83:
Perform chi-square test of independence between two categorical variables. 84:
85: Args: 86: col1: First column name 87: col2: Second column name 88: 89:
Returns: 90: Dictionary containing chi-square test results 91: """ 92: # Get
cross-tabulation 93: crosstab = self.cross_tabulate(col1, col2) 94: 95: if
len(crosstab) < 2 or len(crosstab[0]) < 2: 96: return { 97: 'chi_square': 0.0,
98: 'p_value': 1.0, 99: 'df': 0, 100: 'significant': False,
```

```
101:         'error': 'Insufficient data for chi-square test'
102:     }
103:
104:     # Extract observed frequencies (skip header row and column)
105:     observed = []
106:     for i in range(1, len(crosstab)):
107:         row = []
108:         for j in range(1, len(crosstab[i])):
109:             row.append(int(crosstab[i][j]))
110:         observed.append(row)
111:
112:     # Check if we have enough data for chi-square test
113:     total_observations = sum(sum(row) for row in observed)
114:     if total_observations < 5:  # Chi-square test requires at least 5 observations
115:         return {
116:             'chi_square': 0.0,
117:             'p_value': 1.0,
118:             'df': 0,
119:             'significant': False,
120:             'error': 'Insufficient data for chi-square test'
121:         }
122:
123:     # Calculate expected frequencies
124:     expected = self._calculate_expected_frequencies(observed)
125:
126:     # Calculate chi-square statistic
127:     chi_square = 0.0
128:     for i in range(len(observed)):
129:         for j in range(len(observed[i])):
130:             if expected[i][j] > 0:
131:                 chi_square += ((observed[i][j] - expected[i][j]) ** 2) / expected[i][j]
132:
133:     # Calculate degrees of freedom
134:     df = (len(observed) - 1) * (len(observed[0]) - 1)
135:
136:     # Calculate p-value (approximation using chi-square distribution)
137:     p_value = self._chi_square_p_value(chi_square, df)
138:
139:     # Determine significance (alpha = 0.05)
140:     significant = p_value < 0.05
141:
142:     return {
143:         'chi_square': chi_square,
144:         'p_value': p_value,
145:         'df': df,
146:         'significant': significant,
147:         'observed': observed,
148:         'expected': expected
149:     }
150:
```

```python
151: def _calculate_expected_frequencies(self, observed: List[List[int]]) ->
List[List[float]]: 152: """Calculate expected frequencies for chi-square
test.""" 153: if not observed or not observed[0]: 154: return [] 155: 156: rows =
len(observed) 157: cols = len(observed[0]) 158: 159: # Calculate row and column
totals 160: row_totals = [sum(row) for row in observed] 161: col_totals = [] 162:
for j in range(cols): 163: col_totals.append(sum(observed[i][j] for i in
range(rows))) 164: 165: total = sum(row_totals) 166: 167: # Calculate expected
frequencies 168: expected = [] 169: for i in range(rows): 170: row = [] 171: for
j in range(cols): 172: expected_freq = (row_totals[i] * col_totals[j]) / total if
total > 0 else 0 173: row.append(expected_freq) 174: expected.append(row) 175:
176: return expected 177: 178: def _chi_square_p_value(self, chi_square: float,
df: int) -> float: 179: """ 180: Calculate approximate p-value for chi-square
statistic. 181: This is a simplified approximation - in practice, you'd use a
proper chi-square distribution. 182: """ 183: if df <= 0: 184: return 1.0 185:
186: # Simple approximation for chi-square p-value 187: # For small chi-square
values, p-value is close to 1 188: # For large chi-square values, p-value
approaches 0 189: if chi_square < df: 190: return 1.0 - (chi_square / (df * 2))
191: else: 192: return max(0.0, 1.0 - (chi_square / (df * 10))) 193: 194: def
correlation_analysis(self, col1: str, col2: str) -> Dict[str, Any]: 195: """ 196:
Perform correlation analysis between two variables. 197: 198: Args: 199: col1:
First column name 200: col2: Second column name
```

```python
201:
202:        Returns:
203:            Dictionary containing correlation analysis results
204:        """
205:        # Extract numeric values
206:        values1 = []
207:        values2 = []
208:
209:        for row in self.survey_data:
210:            val1 = row.get(col1)
211:            val2 = row.get(col2)
212:
213:            # Try to convert to numeric
214:            try:
215:                if val1 is not None and str(val1).strip():
216:                    num1 = float(val1)
217:                    if val2 is not None and str(val2).strip():
218:                        num2 = float(val2)
219:                        values1.append(num1)
220:                        values2.append(num2)
221:            except (ValueError, TypeError):
222:                continue
223:
224:        if len(values1) < 2:
225:            return {
226:                'correlation': 0.0,
227:                'sample_size': 0,
228:                'error': 'Insufficient numeric data for correlation analysis'
229:            }
230:
231:        # Calculate correlation coefficient
232:        correlation = self._calculate_correlation(values1, values2)
233:
234:        return {
235:            'correlation': correlation,
236:            'sample_size': len(values1),
237:            'strength': self._interpret_correlation(correlation)
238:        }
239:
240:    def _calculate_correlation(self, x: List[float], y: List[float]) -> float:
241:        """Calculate Pearson correlation coefficient."""
242:        if len(x) != len(y) or len(x) < 2:
243:            return 0.0
244:
245:        n = len(x)
246:
247:        # Calculate means
248:        mean_x = sum(x) / n
249:        mean_y = sum(y) / n
250:
```

```python
251: # Calculate correlation coefficient
252: numerator = sum((x[i] - mean_x) * (y[i] - mean_y) for i in range(n))
253: denominator_x = sum((x[i] - mean_x) ** 2 for i in range(n))
254: denominator_y = sum((y[i] - mean_y) ** 2 for i in range(n))
255:
256: if denominator_x == 0 or denominator_y == 0:
257:     return 0.0
258:
259: correlation = numerator / math.sqrt(denominator_x * denominator_y)
260: return correlation
261:
262: def _interpret_correlation(self, correlation: float) -> str:
263:     """Interpret correlation coefficient strength."""
264:     abs_corr = abs(correlation)
265:
266:     if abs_corr >= 0.8:
267:         return "Very Strong"
268:     elif abs_corr >= 0.6:
269:         return "Strong"
270:     elif abs_corr >= 0.4:
271:         return "Moderate"
272:     elif abs_corr >= 0.2:
273:         return "Weak"
274:     else:
275:         return "Very Weak"
276:
277: def analyze_response_patterns(self) -> Dict[str, Any]:
278:     """Analyze patterns in survey responses."""
279:     patterns = {
280:         'common_combinations': [],
281:         'response_clusters': [],
282:         'outliers': []
283:     }
284:
285:     # Find common combinations of responses
286:     response_combinations = []
287:     for row in self.survey_data:
288:         combination = []
289:         for col in self.columns:
290:             value = row.get(col)
291:             if value is not None and str(value).strip():
292:                 combination.append(f"{col}:{str(value).strip()}")
293:         if combination:
294:             response_combinations.append(tuple(sorted(combination)))
295:
296:     # Count combinations
297:     combination_counts = Counter(response_combinations)
298:     common_combinations = combination_counts.most_common(5)
299:
300:     patterns['common_combinations'] = [
```

```python
301:         {
302:             'combination': list(combo),
303:             'count': count,
304:             'percentage': (count / len(response_combinations)) * 100
305:         }
306:         for combo, count in common_combinations
307:     ]
308:
309:     return patterns
310:
311: def get_statistical_summary(self) -> Dict[str, Any]:
312:     """Generate a comprehensive statistical summary."""
313:     summary = {
314:         'total_responses': self.total_responses,
315:         'total_columns': len(self.columns),
316:         'numeric_columns': [],
317:         'categorical_columns': [],
318:         'statistical_tests': []
319:     }
320:
321:     # Categorize columns
322:     for column in self.columns:
323:         numeric_count = 0
324:         total_count = 0
325:
326:         for row in self.survey_data:
327:             value = row.get(column)
328:             if value is not None and str(value).strip():
329:                 total_count += 1
330:                 try:
331:                     float(str(value))
332:                     numeric_count += 1
333:                 except (ValueError, TypeError):
334:                     pass
335:
336:         if total_count > 0 and (numeric_count / total_count) > 0.5:
337:             summary['numeric_columns'].append(column)
338:         else:
339:             summary['categorical_columns'].append(column)
340:
341:     return summary
342:
343: def perform_multiple_chi_square_tests(self, target_column: str) -> List[Dict[str, Any]]:
344:     """
345:     Perform chi-square tests between a target column and all other categorical columns.
346:
347:     Args:
348:         target_column: The target column to test against
349:
350:     Returns:
```

```
351: List of chi-square test results 352: """ 353: results = [] 354: 355: for
column in self.columns: 356: if column != target_column: 357: try: 358:
test_result = self.chi_square_test(target_column, column) 359:
test_result['column1'] = target_column 360: test_result['column2'] = column 361:
results.append(test_result) 362: except Exception as e: 363: results.append({
364: 'column1': target_column, 365: 'column2': column, 366: 'error': str(e) 367:
}) 368: 369: # Sort by significance 370: results.sort(key=lambda x:
x.get('p_value', 1.0)) 371: 372: return results
```

# 6. Sentiment Analysis Module (sentiment_analyzer.py)

The SentimentAnalyzer class performs text-based sentiment analysis using keyword dictionaries, negation handling, and intensifier recognition. It provides insights into the emotional content of open-ended survey responses.

## Sentiment Analysis Module - sentiment_analyzer.py

```
1: #!/usr/bin/env python3 2: """ 3: Sentiment Analyzer Module 4: Description:
Performs basic sentiment analysis on text responses using 5: keyword matching and
scoring systems without external libraries. 6: """ 7: 8: import re 9: from typing
import List, Dict, Any, Optional 10: from collections import Counter 11: 12: 13:
class SentimentAnalyzer: 14: """Handles sentiment analysis of text responses
using keyword-based approach.""" 15: 16: def __init__(self): 17: """Initialize
the sentiment analyzer with keyword dictionaries.""" 18: # Positive keywords and
their weights 19: self.positive_keywords = { 20: 'excellent': 3, 'amazing': 3,
'great': 2, 'good': 2, 'wonderful': 3, 21: 'fantastic': 3, 'outstanding': 3,
'perfect': 3, 'love': 2, 'like': 1, 22: 'enjoy': 2, 'happy': 2, 'satisfied': 2,
'pleased': 2, 'impressed': 2, 23: 'recommend': 2, 'helpful': 2, 'useful': 1,
'effective': 2, 'quality': 1, 24: 'best': 2, 'awesome': 3, 'brilliant': 3,
'superb': 3, 'terrific': 3, 25: 'delighted': 3, 'thrilled': 3, 'excited': 2,
'positive': 1, 'successful': 1, 26: 'improved': 1, 'better': 1, 'exceeded': 2,
'surpassed': 2, 'outstanding': 3 27: } 28: 29: # Negative keywords and their
weights 30: self.negative_keywords = { 31: 'terrible': 3, 'awful': 3, 'horrible':
3, 'bad': 2, 'poor': 2, 32: 'disappointing': 2, 'frustrated': 2, 'angry': 2,
'upset': 2, 'annoyed': 2, 33: 'hate': 3, 'dislike': 2, 'worst': 3, 'useless': 2,
'waste': 2, 34: 'problem': 1, 'issue': 1, 'difficult': 1, 'confusing': 1,
'complicated': 1, 35: 'broken': 2, 'failed': 2, 'error': 1, 'bug': 1, 'crash': 2,
36: 'slow': 1, 'expensive': 1, 'overpriced': 2, 'cheap': 1, 'low quality': 2, 37:
'unreliable': 2, 'unstable': 2, 'inconsistent': 1, 'disorganized': 1, 38:
'messy': 1, 'chaotic': 2, 'stressful': 2, 'overwhelming': 2 39: } 40: 41: #
Neutral keywords (used for context) 42: self.neutral_keywords = { 43: 'okay': 0,
'fine': 0, 'average': 0, 'normal': 0, 'standard': 0, 44: 'usual': 0, 'typical':
0, 'regular': 0, 'common': 0, 'basic': 0, 45: 'simple': 0, 'straightforward': 0,
'clear': 0, 'understandable': 0, 46: 'adequate': 0, 'sufficient': 0,
'acceptable': 0, 'reasonable': 0 47: } 48: 49: # Negation words that can flip
sentiment 50: self.negation_words = {
```

```python
51:     'not', 'no', 'never', 'none', 'neither', 'nor', 'nobody', 'nothing',
52:     'nowhere', 'hardly', 'barely', 'scarcely', 'doesn\'t', 'don\'t',
53:     'didn\'t', 'won\'t', 'can\'t', 'couldn\'t', 'wouldn\'t', 'shouldn\'t',
54:     'isn\'t', 'aren\'t', 'wasn\'t', 'weren\'t', 'hasn\'t', 'haven\'t',
55:     'hadn\'t', 'doesnt', 'dont', 'didnt', 'wont', 'cant', 'couldnt',
56:     'wouldnt', 'shouldnt', 'isnt', 'arent', 'wasnt', 'werent', 'hasnt',
57:     'havent', 'hadnt'
58: }
59:
60: # Intensifier words that amplify sentiment
61: self.intensifier_words = {
62:     'very': 1.5, 'really': 1.5, 'extremely': 2.0, 'absolutely': 2.0,
63:     'completely': 2.0, 'totally': 2.0, 'entirely': 2.0, 'thoroughly': 1.5,
64:     'highly': 1.5, 'incredibly': 2.0, 'amazingly': 2.0, 'exceptionally': 2.0,
65:     'particularly': 1.2, 'especially': 1.2, 'notably': 1.2, 'remarkably': 1.5
66: }
67:
68: def analyze_text(self, text: str) -> Dict[str, Any]:
69:     """
70:     Analyze sentiment of a single text response.
71:
72:     Args:
73:         text: Text string to analyze
74:
75:     Returns:
76:         Dictionary containing sentiment analysis results
77:     """
78:     if not text or not isinstance(text, str):
79:         return {
80:             'sentiment': 'neutral',
81:             'score': 0,
82:             'positive_words': [],
83:             'negative_words': [],
84:             'confidence': 0.0
85:         }
86:
87:     # Clean and normalize text
88:     cleaned_text = self._clean_text(text)
89:
90:     # Extract words
91:     words = self._extract_words(cleaned_text)
92:
93:     # Analyze sentiment
94:     sentiment_score = 0
95:     positive_words = []
96:     negative_words = []
97:     intensifier_count = 0
98:
99:     for i, word in enumerate(words):
100:         word_lower = word.lower()
```

```python
101: 102: # Check for intensifiers 103: if word_lower in self.intensifier_words:
104: intensifier_count += 1 105: continue 106: 107: # Check for negations 108:
is_negated = self._is_negated(words, i) 109: 110: # Check positive keywords 111:
if word_lower in self.positive_keywords: 112: weight =
self.positive_keywords[word_lower] 113: if is_negated: 114: sentiment_score -=
weight 115: negative_words.append(word) 116: else: 117: sentiment_score += weight
118: positive_words.append(word) 119: 120: # Check negative keywords 121: elif
word_lower in self.negative_keywords: 122: weight =
self.negative_keywords[word_lower] 123: if is_negated: 124: sentiment_score +=
weight 125: positive_words.append(word) 126: else: 127: sentiment_score -= weight
128: negative_words.append(word) 129: 130: # Apply intensifier multiplier 131: if
intensifier_count > 0: 132: sentiment_score *= (1 + (intensifier_count * 0.2))
133: 134: # Determine sentiment category 135: sentiment =
self._categorize_sentiment(sentiment_score) 136: 137: # Calculate confidence
138: total_words = len(words) 139: sentiment_words = len(positive_words) +
len(negative_words) 140: confidence = min(1.0, sentiment_words /
max(total_words, 1)) 141: 142: return { 143: 'sentiment': sentiment, 144:
'score': sentiment_score, 145: 'positive_words': positive_words, 146:
'negative_words': negative_words, 147: 'confidence': confidence, 148:
'total_words': total_words, 149: 'sentiment_words': sentiment_words 150: }
```

```python
151: 152: def analyze_column(self, survey_data: List[Dict[str, Any]], column:
str) -> Dict[str, Any]: 153: """ 154: Analyze sentiment for all responses in a
specific column. 155: 156: Args: 157: survey_data: List of survey responses 158:
column: Column name to analyze 159: 160: Returns: 161: Dictionary containing
aggregated sentiment analysis results 162: """ 163: if not survey_data or column
not in survey_data[0]: 164: return { 165: 'positive': 0, 166: 'negative': 0, 167:
'neutral': 0, 168: 'positive_pct': 0.0, 169: 'negative_pct': 0.0, 170:
'neutral_pct': 0.0, 171: 'avg_score': 0.0, 172: 'total_responses': 0 173: } 174:
175: sentiment_results = [] 176: total_responses = 0 177: 178: for row in
survey_data: 179: text = row.get(column) 180: if text is not None and
str(text).strip(): 181: result = self.analyze_text(str(text)) 182:
sentiment_results.append(result) 183: total_responses += 1 184: 185: if not
sentiment_results: 186: return { 187: 'positive': 0, 188: 'negative': 0, 189:
'neutral': 0, 190: 'positive_pct': 0.0, 191: 'negative_pct': 0.0, 192:
'neutral_pct': 0.0, 193: 'avg_score': 0.0, 194: 'total_responses': 0 195: } 196:
197: # Count sentiments 198: positive_count = sum(1 for r in sentiment_results if
r['sentiment'] == 'positive') 199: negative_count = sum(1 for r in
sentiment_results if r['sentiment'] == 'negative') 200: neutral_count = sum(1 for
r in sentiment_results if r['sentiment'] == 'neutral')
```

```python
201:
202:         # Calculate percentages
203:         total = len(sentiment_results)
204:         positive_pct = (positive_count / total) * 100
205:         negative_pct = (negative_count / total) * 100
206:         neutral_pct = (neutral_count / total) * 100
207:
208:         # Calculate average score
209:         avg_score = sum(r['score'] for r in sentiment_results) / total
210:
211:         return {
212:             'positive': positive_count,
213:             'negative': negative_count,
214:             'neutral': neutral_count,
215:             'positive_pct': positive_pct,
216:             'negative_pct': negative_pct,
217:             'neutral_pct': neutral_pct,
218:             'avg_score': avg_score,
219:             'total_responses': total_responses,
220:             'detailed_results': sentiment_results
221:         }
222:
223:     def analyze_all_text_columns(self, survey_data: List[Dict[str, Any]]) -> Dict[str, Any]:
224:         """
225:         Analyze sentiment for all text columns in the survey data.
226:
227:         Args:
228:             survey_data: List of survey responses
229:
230:         Returns:
231:             Dictionary containing sentiment analysis for all text columns
232:         """
233:         if not survey_data:
234:             return {}
235:
236:         text_columns = []
237:         for col in survey_data[0].keys():
238:             # Check if column contains text data
239:             text_count = 0
240:             total_count = 0
241:
242:             for row in survey_data[:10]:  # Sample first 10 rows
243:                 value = row.get(col)
244:                 if value is not None and str(value).strip():
245:                     total_count += 1
246:                     if len(str(value)) > 20:  # Consider it text if longer than 20 chars
247:                         text_count += 1
248:
249:             if total_count > 0 and (text_count / total_count) > 0.3:
250:                 text_columns.append(col)
```

```python
251: 252:    results = {}
253:        for column in text_columns:
254:            results[column] = self.analyze_column(survey_data, column)
255: 256:        return results
257: 258:    def _clean_text(self, text: str) -> str:
259:        """Clean and normalize text for analysis."""
260:        # Convert to lowercase
261:        text = text.lower()
262: 263:        # Remove extra whitespace
264:        text = re.sub(r'\s+', ' ', text)
265: 266:        # Remove punctuation (keep apostrophes for contractions)
267:        text = re.sub(r'[^\w\s\']', ' ', text)
268: 269:        return text.strip()
270: 271:    def _extract_words(self, text: str) -> List[str]:
272:        """Extract words from text."""
273:        return text.split()
274: 275:    def _is_negated(self, words: List[str], current_index: int) -> bool:
276:        """Check if current word is negated by previous words."""
277:        # Look back up to 3 words for negation
278:        start_index = max(0, current_index - 3)
279: 280:        for i in range(start_index, current_index):
281:            if i < len(words) and words[i].lower() in self.negation_words:
282:                return True
283: 284:        return False
285: 286:    def _categorize_sentiment(self, score: float) -> str:
287:        """Categorize sentiment based on score."""
288:        if score > 1.0:
289:            return 'positive'
290:        elif score < -1.0:
291:            return 'negative'
292:        else:
293:            return 'neutral'
294: 295:    def get_sentiment_summary(self, sentiment_results: List[Dict[str, Any]]) -> Dict[str, Any]:
296:        """Generate a summary of sentiment analysis results."""
297:        if not sentiment_results:
298:            return {}
299: 300:        total_responses = len(sentiment_results)
```

```python
301: positive_count = sum(1 for r in sentiment_results if r['sentiment'] ==
'positive') 302: negative_count = sum(1 for r in sentiment_results if
r['sentiment'] == 'negative') 303: neutral_count = sum(1 for r in
sentiment_results if r['sentiment'] == 'neutral') 304: 305: # Most common
positive and negative words 306: all_positive_words = [] 307: all_negative_words
= [] 308: 309: for result in sentiment_results: 310:
all_positive_words.extend(result['positive_words']) 311:
all_negative_words.extend(result['negative_words']) 312: 313:
positive_word_counts = Counter(all_positive_words) 314: negative_word_counts =
Counter(all_negative_words) 315: 316: return { 317: 'total_responses':
total_responses, 318: 'sentiment_distribution': { 319: 'positive': {'count':
positive_count, 'percentage': (positive_count / total_responses) * 100}, 320:
'negative': {'count': negative_count, 'percentage': (negative_count /
total_responses) * 100}, 321: 'neutral': {'count': neutral_count, 'percentage':
(neutral_count / total_responses) * 100} 322: }, 323: 'top_positive_words':
positive_word_counts.most_common(5), 324: 'top_negative_words':
negative_word_counts.most_common(5), 325: 'average_confidence':
sum(r['confidence'] for r in sentiment_results) / total_responses, 326:
'average_score': sum(r['score'] for r in sentiment_results) / total_responses
327: } 328: 329: def export_sentiment_report(self, sentiment_results: Dict[str,
Any], filename: str = "sentiment_report.txt") -> bool: 330: """Export sentiment
analysis results to a text file.""" 331: try: 332: with open(filename, 'w',
encoding='utf-8') as file: 333: file.write("SENTIMENT ANALYSIS REPORT\n") 334:
file.write("=" * 50 + "\n\n") 335: 336: for column, results in
sentiment_results.items(): 337: file.write(f"Column: {column}\n") 338:
file.write("-" * 30 + "\n") 339: file.write(f"Total Responses:
{results['total_responses']}\n") 340: file.write(f"Positive:
{results['positive']} ({results['positive_pct']:.1f}%)\n") 341:
file.write(f"Negative: {results['negative']}
({results['negative_pct']:.1f}%)\n") 342: file.write(f"Neutral:
{results['neutral']} ({results['neutral_pct']:.1f}%)\n") 343:
file.write(f"Average Score: {results['avg_score']:.2f}\n\n") 344: 345: return
True 346: 347: except Exception as e: 348: print(f"Error exporting sentiment
report: {str(e)}") 349: return False
```

# 7. Pattern Detection Module (pattern_detector.py)

The PatternDetector class identifies correlations, trends, and patterns in survey responses. It analyzes demographic patterns, response combinations, and outlier detection to uncover meaningful insights in the data.

## Pattern Detection Module - pattern_detector.py

```
1: #!/usr/bin/env python3 2: """ 3: Pattern Detector Module 4: Author: Student
Developer 5: Description: Identifies correlations, patterns, and trends in survey
responses 6: using statistical analysis and pattern recognition techniques. 7:
""" 8: 9: import math 10: from typing import List, Dict, Any, Optional, Tuple 11:
from collections import defaultdict, Counter 12: import itertools 13: 14: 15:
class PatternDetector: 16: """Handles pattern detection and correlation analysis
in survey data.""" 17: 18: def __init__(self, survey_data: List[Dict[str, Any]]):
19: """ 20: Initialize the pattern detector. 21: 22: Args: 23: survey_data: List
of dictionaries containing survey responses 24: """ 25: self.survey_data =
survey_data 26: self.total_responses = len(survey_data) 27: self.columns =
list(survey_data[0].keys()) if survey_data else [] 28: 29: def
find_patterns(self) -> List[Dict[str, Any]]: 30: """ 31: Find patterns and
correlations in survey responses. 32: 33: Returns: 34: List of detected patterns
with descriptions and confidence levels 35: """ 36: if not self.survey_data: 37:
return [] 38: 39: patterns = [] 40: 41: # Find demographic patterns 42:
patterns.extend(self._find_demographic_patterns()) 43: 44: # Find response
correlation patterns 45: patterns.extend(self._find_correlation_patterns()) 46:
47: # Find response combination patterns 48:
patterns.extend(self._find_combination_patterns()) 49: 50: # Find outlier
patterns
```

```python
51:         patterns.extend(self._find_outlier_patterns())
52:
53:         # Sort patterns by confidence
54:         patterns.sort(key=lambda x: x['confidence'], reverse=True)
55:
56:         return patterns
57:
58:     def _find_demographic_patterns(self) -> List[Dict[str, Any]]:
59:         """Find patterns related to demographics."""
60:         patterns = []
61:
62:         # Age-based patterns
63:         if 'age' in self.columns:
64:             age_patterns = self._analyze_age_patterns()
65:             patterns.extend(age_patterns)
66:
67:         # Gender-based patterns
68:         if 'gender' in self.columns:
69:             gender_patterns = self._analyze_gender_patterns()
70:             patterns.extend(gender_patterns)
71:
72:         # Regional patterns
73:         if 'region' in self.columns:
74:             regional_patterns = self._analyze_regional_patterns()
75:             patterns.extend(regional_patterns)
76:
77:         # Education-based patterns
78:         if 'education' in self.columns:
79:             education_patterns = self._analyze_education_patterns()
80:             patterns.extend(education_patterns)
81:
82:         return patterns
83:
84:     def _analyze_age_patterns(self) -> List[Dict[str, Any]]:
85:         """Analyze patterns based on age groups."""
86:         patterns = []
87:
88:         # Group responses by age
89:         age_groups = {
90:             '18-25': [],
91:             '26-35': [],
92:             '36-45': [],
93:             '46-55': [],
94:             '56-65': [],
95:             '65+': []
96:         }
97:
98:         for row in self.survey_data:
99:             age_str = row.get('age')
100:             if age_str and str(age_str).isdigit():
```

```python
101: try:
102:     age = int(age_str)
103:     if 18 <= age <= 25:
104:         age_groups['18-25'].append(row)
105:     elif 26 <= age <= 35:
106:         age_groups['26-35'].append(row)
107:     elif 36 <= age <= 45:
108:         age_groups['36-45'].append(row)
109:     elif 46 <= age <= 55:
110:         age_groups['46-55'].append(row)
111:     elif 56 <= age <= 65:
112:         age_groups['56-65'].append(row)
113:     elif age > 65:
114:         age_groups['65+'].append(row)
115: except ValueError:
116:     continue
117:
118: # Analyze patterns for each question
119: for column in self.columns:
120:     if column not in ['age', 'gender', 'region', 'education']:
121:         column_patterns = self._analyze_column_by_age_groups(column, age_groups)
122:         patterns.extend(column_patterns)
123:
124: return patterns
125:
126: def _analyze_column_by_age_groups(self, column: str, age_groups: Dict[str, List[Dict[str, Any]]]) -> List[Dict[str, Any]]:
127:     """Analyze a specific column's responses by age groups."""
128:     patterns = []
129:
130:     # Get most common response for each age group
131:     age_group_responses = {}
132:     for age_group, responses in age_groups.items():
133:         if responses:
134:             values = [r.get(column) for r in responses if r.get(column)]
135:             if values:
136:                 value_counts = Counter(values)
137:                 most_common = value_counts.most_common(1)[0]
138:                 age_group_responses[age_group] = {
139:                     'response': most_common[0],
140:                     'count': most_common[1],
141:                     'percentage': (most_common[1] / len(values)) * 100
142:                 }
143:
144:     # Find patterns
145:     if len(age_group_responses) > 1:
146:         # Find age groups with similar responses
147:         response_groups = defaultdict(list)
148:         for age_group, data in age_group_responses.items():
149:             response_groups[data['response']].append(age_group)
150:
```

```python
151: for response, age_groups_list in response_groups.items(): 152: if
len(age_groups_list) > 1: 153: confidence = min(90, len(age_groups_list) * 20)
154: patterns.append({ 155: 'type': 'age_pattern', 156: 'description': f"Age
groups {', '.join(age_groups_list)} most commonly responded '{response}' to
{column}", 157: 'confidence': confidence, 158: 'sample_size':
sum(len(age_groups[ag]) for ag in age_groups_list), 159: 'response': response,
160: 'affected_groups': age_groups_list 161: }) 162: 163: return patterns 164:
165: def _analyze_gender_patterns(self) -> List[Dict[str, Any]]: 166: """Analyze
patterns based on gender.""" 167: patterns = [] 168: 169: # Group responses by
gender 170: gender_groups = defaultdict(list) 171: for row in self.survey_data:
172: gender = row.get('gender') 173: if gender and str(gender).strip(): 174:
gender_groups[str(gender).strip().title()].append(row) 175: 176: # Analyze
patterns for each question 177: for column in self.columns: 178: if column not in
['age', 'gender', 'region', 'education']: 179: column_patterns =
self._analyze_column_by_gender(column, gender_groups) 180:
patterns.extend(column_patterns) 181: 182: return patterns 183: 184: def
_analyze_column_by_gender(self, column: str, gender_groups: Dict[str,
List[Dict[str, Any]]]) -> List[Dict[str, Any]]: 185: """Analyze a specific
column's responses by gender.""" 186: patterns = [] 187: 188: # Get most common
response for each gender 189: gender_responses = {} 190: for gender, responses in
gender_groups.items(): 191: if responses: 192: values = [r.get(column) for r in
responses if r.get(column)] 193: if values: 194: value_counts = Counter(values)
195: most_common = value_counts.most_common(1)[0] 196: gender_responses[gender]
= { 197: 'response': most_common[0], 198: 'count': most_common[1], 199:
'percentage': (most_common[1] / len(values)) * 100 200: }
```

```python
201: 202: # Find gender differences 203: if len(gender_responses) > 1: 204:
responses = list(gender_responses.values()) 205: if len(set(r['response'] for r
in responses)) > 1: 206: # Different genders have different most common responses
207: confidence = 75 208: patterns.append({ 209: 'type': 'gender_pattern', 210:
'description': f"Gender differences detected in {column} responses", 211:
'confidence': confidence, 212: 'sample_size': sum(len(gender_groups[g]) for g in
gender_responses.keys()), 213: 'details': gender_responses 214: }) 215: 216:
return patterns 217: 218: def _analyze_regional_patterns(self) -> List[Dict[str,
Any]]: 219: """Analyze patterns based on region.""" 220: patterns = [] 221: 222:
# Group responses by region 223: regional_groups = defaultdict(list) 224: for row
in self.survey_data: 225: region = row.get('region') 226: if region and
str(region).strip(): 227:
regional_groups[str(region).strip().title()].append(row) 228: 229: # Analyze
patterns for each question 230: for column in self.columns: 231: if column not in
['age', 'gender', 'region', 'education']: 232: column_patterns =
self._analyze_column_by_region(column, regional_groups) 233:
patterns.extend(column_patterns) 234: 235: return patterns 236: 237: def
_analyze_column_by_region(self, column: str, regional_groups: Dict[str,
List[Dict[str, Any]]]) -> List[Dict[str, Any]]: 238: """Analyze a specific
column's responses by region.""" 239: patterns = [] 240: 241: # Get most common
response for each region 242: regional_responses = {} 243: for region, responses
in regional_groups.items(): 244: if responses: 245: values = [r.get(column) for r
in responses if r.get(column)] 246: if values: 247: value_counts =
Counter(values) 248: most_common = value_counts.most_common(1)[0] 249:
regional_responses[region] = { 250: 'response': most_common[0],
```

```python
251:         'count': most_common[1],
252:         'percentage': (most_common[1] / len(values)) * 100
253:     }
254:
255: # Find regional patterns
256: if len(regional_responses) > 1:
257:     # Find regions with similar responses
258:     response_groups = defaultdict(list)
259:     for region, data in regional_responses.items():
260:         response_groups[data['response']].append(region)
261:
262:     for response, regions_list in response_groups.items():
263:         if len(regions_list) > 1:
264:             confidence = min(85, len(regions_list) * 25)
265:             patterns.append({
266:                 'type': 'regional_pattern',
267:                 'description': f"Regions {', '.join(regions_list)} most commonly responded '{response}' to {column}",
268:                 'confidence': confidence,
269:                 'sample_size': sum(len(regional_groups[r]) for r in regions_list),
270:                 'response': response,
271:                 'affected_regions': regions_list
272:             })
273:
274: return patterns
275:
276: def _analyze_education_patterns(self) -> List[Dict[str, Any]]:
277:     """Analyze patterns based on education level."""
278:     patterns = []
279:
280:     # Group responses by education
281:     education_groups = defaultdict(list)
282:     for row in self.survey_data:
283:         education = row.get('education')
284:         if education and str(education).strip():
285:             education_groups[str(education).strip().title()].append(row)
286:
287:     # Analyze patterns for each question
288:     for column in self.columns:
289:         if column not in ['age', 'gender', 'region', 'education']:
290:             column_patterns = self._analyze_column_by_education(column, education_groups)
291:             patterns.extend(column_patterns)
292:
293:     return patterns
294:
295: def _analyze_column_by_education(self, column: str, education_groups: Dict[str, List[Dict[str, Any]]]) -> List[Dict[str, Any]]:
296:     """Analyze a specific column's responses by education level."""
297:     patterns = []
298:
299:     # Get most common response for each education level
300:     education_responses = {}
```

```
301: for education, responses in education_groups.items(): 302: if responses:
303: values = [r.get(column) for r in responses if r.get(column)] 304: if values:
305: value_counts = Counter(values) 306: most_common =
value_counts.most_common(1)[0] 307: education_responses[education] = { 308:
'response': most_common[0], 309: 'count': most_common[1], 310: 'percentage':
(most_common[1] / len(values)) * 100 311: } 312: 313: # Find education-based
patterns 314: if len(education_responses) > 1: 315: # Find education levels with
similar responses 316: response_groups = defaultdict(list) 317: for education,
data in education_responses.items(): 318:
response_groups[data['response']].append(education) 319: 320: for response,
education_levels in response_groups.items(): 321: if len(education_levels) > 1:
322: confidence = min(80, len(education_levels) * 20) 323: patterns.append({ 324:
'type': 'education_pattern', 325: 'description': f"Education levels {',
'.join(education_levels)} most commonly responded '{response}' to {column}",
326: 'confidence': confidence, 327: 'sample_size': sum(len(education_groups[e])
for e in education_levels), 328: 'response': response, 329:
'affected_education_levels': education_levels 330: }) 331: 332: return patterns
333: 334: def _find_correlation_patterns(self) -> List[Dict[str, Any]]: 335:
"""Find correlation patterns between different questions.""" 336: patterns = []
337: 338: # Analyze correlations between categorical variables 339:
categorical_columns = [col for col in self.columns if col not in ['age',
'gender', 'region', 'education']] 340: 341: for i, col1 in
enumerate(categorical_columns): 342: for col2 in categorical_columns[i+1:]: 343:
correlation_pattern = self._analyze_correlation(col1, col2) 344: if
correlation_pattern: 345: patterns.append(correlation_pattern) 346: 347: return
patterns 348: 349: def _analyze_correlation(self, col1: str, col2: str) ->
Optional[Dict[str, Any]]: 350: """Analyze correlation between two columns."""
```

```
351: # Get unique values for both columns 352: values1 = set() 353: values2 =
set() 354: 355: for row in self.survey_data: 356: val1 = row.get(col1) 357: val2
= row.get(col2) 358: 359: if val1 is not None and str(val1).strip(): 360:
values1.add(str(val1).strip()) 361: if val2 is not None and str(val2).strip():
362: values2.add(str(val2).strip()) 363: 364: if len(values1) < 2 or len(values2)
< 2: 365: return None 366: 367: # Calculate correlation strength 368:
total_responses = 0 369: matching_responses = 0 370: 371: for row in
self.survey_data: 372: val1 = row.get(col1) 373: val2 = row.get(col2) 374: 375:
if val1 is not None and str(val1).strip() and val2 is not None and
str(val2).strip(): 376: total_responses += 1 377: # Check if responses are
related (simplified correlation) 378: if
self._are_responses_related(str(val1).strip(), str(val2).strip()): 379:
matching_responses += 1 380: 381: if total_responses == 0: 382: return None 383:
384: correlation_strength = matching_responses / total_responses 385: 386: if
correlation_strength > 0.6: # Strong correlation threshold 387: return { 388:
'type': 'correlation_pattern', 389: 'description': f"Strong correlation detected
between {col1} and {col2}", 390: 'confidence': min(90, correlation_strength *
100), 391: 'sample_size': total_responses, 392: 'correlation_strength':
correlation_strength, 393: 'columns': [col1, col2] 394: } 395: 396: return None
397: 398: def _are_responses_related(self, response1: str, response2: str) ->
bool: 399: """Check if two responses are related (simplified logic).""" 400: #
This is a simplified correlation check
```

```python
401: # In a real implementation, you might use more sophisticated methods
402:
403: # Check for similar sentiment
404: positive_words = ['good', 'great', 'excellent', 'satisfied', 'happy', 'like', 'love']
405: negative_words = ['bad', 'poor', 'terrible', 'dissatisfied', 'unhappy', 'dislike', 'hate']
406:
407: response1_lower = response1.lower()
408: response2_lower = response2.lower()
409:
410: # Check if both responses have similar sentiment
411: response1_positive = any(word in response1_lower for word in positive_words)
412: response1_negative = any(word in response1_lower for word in negative_words)
413: response2_positive = any(word in response2_lower for word in positive_words)
414: response2_negative = any(word in response2_lower for word in negative_words)
415:
416: if response1_positive and response2_positive:
417: return True
418: if response1_negative and response2_negative:
419: return True
420:
421: # Check for exact matches
422: if response1_lower == response2_lower:
423: return True
424:
425: return False
426:
427: def _find_combination_patterns(self) -> List[Dict[str, Any]]:
428: """Find patterns in response combinations."""
429: patterns = []
430:
431: # Find common response combinations
432: response_combinations = []
433: for row in self.survey_data:
434: combination = []
435: for col in self.columns:
436: value = row.get(col)
437: if value is not None and str(value).strip():
438: combination.append(f"{col}:{str(value).strip()}")
439: if combination:
440: response_combinations.append(tuple(sorted(combination)))
441:
442: # Count combinations
443: combination_counts = Counter(response_combinations)
444: common_combinations = combination_counts.most_common(3)
445:
446: for combo, count in common_combinations:
447: if count > 1: # Only report if combination appears more than once
448: percentage = (count / len(response_combinations)) * 100
449: if percentage > 10: # Only report if more than 10% of responses
450: patterns.append({
```

```
451: 'type': 'combination_pattern', 452: 'description': f"Common response
combination: {', '.join(combo)}", 453: 'confidence': min(85, percentage * 2),
454: 'sample_size': count, 455: 'percentage': percentage, 456: 'combination':
list(combo) 457: }) 458: 459: return patterns 460: 461: def
_find_outlier_patterns(self) -> List[Dict[str, Any]]: 462: """Find outlier
patterns in responses.""" 463: patterns = [] 464: 465: # Find responses that are
significantly different from the norm 466: for column in self.columns: 467: if
column not in ['age', 'gender', 'region', 'education']: 468: outlier_pattern =
self._analyze_outliers(column) 469: if outlier_pattern: 470:
patterns.append(outlier_pattern) 471: 472: return patterns 473: 474: def
_analyze_outliers(self, column: str) -> Optional[Dict[str, Any]]: 475:
"""Analyze outliers in a specific column.""" 476: values = [row.get(column) for
row in self.survey_data if row.get(column)] 477: 478: if not values: 479: return
None 480: 481: # Count responses 482: value_counts = Counter(values) 483:
total_responses = len(values) 484: 485: # Find responses that appear very rarely
(outliers) 486: outlier_threshold = total_responses * 0.05 # 5% threshold 487:
488: outliers = [] 489: for value, count in value_counts.items(): 490: if count
<= outlier_threshold and count > 0: 491: outliers.append({ 492: 'value': value,
493: 'count': count, 494: 'percentage': (count / total_responses) * 100 495: })
496: 497: if outliers: 498: return { 499: 'type': 'outlier_pattern', 500:
'description': f"Outlier responses detected in {column}",
```

```
501:                'confidence': 70,
502:                'sample_size': total_responses,
503:                'outliers':
     outliers
504:            }
505:
506:        return None
507:
508:    def get_pattern_summary(self) ->
     Dict[str, Any]:
509:        """Generate a summary of detected patterns."""
510:        patterns
     = self.find_patterns()
511:
512:        summary = {
513:            'total_patterns':
     len(patterns),
514:            'pattern_types': Counter(p['type'] for p in patterns),
515:
     'high_confidence_patterns': [p for p in patterns if p['confidence'] >= 80],
516:
     'medium_confidence_patterns': [p for p in patterns if 60 <= p['confidence'] <
     80],
517:            'low_confidence_patterns': [p for p in patterns if p['confidence'] <
     60]
518:        }
519:
520:        return summary
```

# 8. Report Generation Module (report_generator.py)

The ReportGenerator class creates comprehensive, professionally formatted analysis reports. It combines all analysis results into structured reports with executive summaries, key findings, and actionable recommendations.

## Report Generation Module - report_generator.py

```python
1: #!/usr/bin/env python3 2: """ 3: Report Generator Module 4: Description:
Generates comprehensive analysis reports from survey data 5: and saves them to
text files with proper formatting and structure. 6: """ 7: 8: import os 9: from
datetime import datetime 10: from typing import List, Dict, Any, Optional 11: 12:
13: class ReportGenerator: 14: """Handles generation of comprehensive survey
analysis reports.""" 15: 16: def __init__(self): 17: """Initialize the report
generator.""" 18: self.report_sections = [] 19: self.current_datetime =
datetime.now() 20: 21: def generate_report(self, report_data: Dict[str, Any]) ->
bool: 22: """ 23: Generate a comprehensive survey analysis report. 24: 25: Args:
26: report_data: Dictionary containing all analysis results 27: 28: Returns: 29:
True if report was generated successfully, False otherwise 30: """ 31: try: 32: #
Extract data from report_data 33: summary = report_data.get('summary', {}) 34:
sentiment = report_data.get('sentiment', {}) 35: patterns =
report_data.get('patterns', []) 36: file_path = report_data.get('file_path',
'survey_report.txt') 37: 38: # Build report content 39: report_content =
self._build_report_content(summary, sentiment, patterns) 40: 41: # Write report
to file 42: success = self._write_report_to_file(report_content, file_path) 43:
44: return success 45: 46: except Exception as e: 47: print(f"Error generating
report: {str(e)}") 48: return False 49: 50: def _build_report_content(self,
summary: Dict[str, Any], sentiment: Dict[str, Any], patterns: List[Dict[str,
Any]]) -> str:
```

```python
51:     """Build the complete report content."""
52:     report_lines = []
53:
54:     # Header
55:     report_lines.extend(self._generate_header())
56:
57:     # Executive Summary
58:     report_lines.extend(self._generate_executive_summary(summary))
59:
60:     # Survey Overview
61:     report_lines.extend(self._generate_survey_overview(summary))
62:
63:     # Demographic Analysis
64:     report_lines.extend(self._generate_demographic_analysis(summary))
65:
66:     # Sentiment Analysis
67:     report_lines.extend(self._generate_sentiment_analysis(sentiment))
68:
69:     # Pattern Analysis
70:     report_lines.extend(self._generate_pattern_analysis(patterns))
71:
72:     # Statistical Analysis
73:     report_lines.extend(self._generate_statistical_analysis(summary))
74:
75:     # Key Findings
76:     report_lines.extend(self._generate_key_findings(summary, sentiment, patterns))
77:
78:     # Recommendations
79:     report_lines.extend(self._generate_recommendations(summary, sentiment, patterns))
80:
81:     # Footer
82:     report_lines.extend(self._generate_footer())
83:
84:     return "\n".join(report_lines)
85:
86: def _generate_header(self) -> List[str]:
87:     """Generate the report header."""
88:     header = [
89:         "=" * 80,
90:         "SURVEY DATA ANALYSIS REPORT",
91:         "=" * 80,
92:         f"Generated on: {self.current_datetime.strftime('%B %d, %Y at %I:%M %p')}",
93:         f"Report ID: SUR-{self.current_datetime.strftime('%Y%m%d-%H%M%S')}",
94:         "",
95:         "This report provides a comprehensive analysis of survey responses including",
96:         "demographic breakdowns, sentiment analysis, pattern detection, and",
97:         "statistical insights to support data-driven decision making.",
98:         "",
99:         "=" * 80,
100:        ""
```

```python
101:     ]
102:     return header
103:
104:     def _generate_executive_summary(self, summary: Dict[str, Any]) -> List[str]:
105:         """Generate the executive summary section."""
106:         lines = [
107:             "EXECUTIVE SUMMARY",
108:             "-" * 50,
109:             ""
110:         ]
111:
112:         if summary:
113:             total_responses = summary.get('total_responses', 0)
114:             response_rate = summary.get('response_rate', 0)
115:
116:             lines.extend([
117:                 f"■ Total Survey Responses: {total_responses:,}",
118:                 f"■ Response Rate: {response_rate:.1f}%",
119:                 "",
120:                 "Key Highlights:",
121:                 "• Comprehensive analysis of survey data across multiple dimensions",
122:                 "• Demographic breakdowns reveal respondent characteristics",
123:                 "• Sentiment analysis provides insights into respondent attitudes",
124:                 "• Pattern detection identifies correlations and trends",
125:                 "• Statistical analysis supports evidence-based conclusions",
126:                 ""
127:             ])
128:
129:         return lines
130:
131:     def _generate_survey_overview(self, summary: Dict[str, Any]) -> List[str]:
132:         """Generate the survey overview section."""
133:         lines = [
134:             "SURVEY OVERVIEW",
135:             "-" * 50,
136:             ""
137:         ]
138:
139:         if summary:
140:             total_responses = summary.get('total_responses', 0)
141:             data_quality = summary.get('data_quality', {})
142:
143:             lines.extend([
144:                 f"■ Survey Details:",
145:                 f" • Total Responses: {total_responses:,}",
146:                 f" • Data Quality: {self._assess_data_quality(data_quality)}",
147:                 ""
148:             ])
149:
150:         # Data quality metrics
```

```python
151: if data_quality: 152: completeness = data_quality.get('completeness', {})
153: if completeness: 154: lines.append("■ Data Completeness by Column:") 155:
for column, completeness_pct in completeness.items(): 156: lines.append(f" •
{column}: {completeness_pct:.1f}%") 157: lines.append("") 158: 159: return lines
160: 161: def _generate_demographic_analysis(self, summary: Dict[str, Any]) ->
List[str]: 162: """Generate the demographic analysis section.""" 163: lines = [
164: "DEMOGRAPHIC ANALYSIS", 165: "-" * 50, 166: "" 167: ] 168: 169: demographics
= summary.get('demographics', {}) 170: 171: if demographics: 172:
lines.append("■ Respondent Demographics:") 173: lines.append("") 174: 175: for
demo_field, breakdown in demographics.items(): 176: if breakdown: 177:
lines.append(f"■ {demo_field.title()}:") 178: total_demo =
sum(breakdown.values()) 179: 180: for category, count in breakdown.items(): 181:
percentage = (count / total_demo) * 100 182: lines.append(f" • {category}:
{count} ({percentage:.1f}%)") 183: 184: lines.append("") 185: else: 186:
lines.append("No demographic data available for analysis.") 187:
lines.append("") 188: 189: return lines 190: 191: def
_generate_sentiment_analysis(self, sentiment: Dict[str, Any]) -> List[str]: 192:
"""Generate the sentiment analysis section.""" 193: lines = [ 194: "SENTIMENT
ANALYSIS", 195: "-" * 50, 196: "" 197: ] 198: 199: if sentiment: 200:
lines.append("■ Text Response Sentiment Analysis:")
```

```
201: lines.append("") 202: 203: for column, results in sentiment.items(): 204: if
isinstance(results, dict) and 'total_responses' in results: 205:
lines.append(f"■ {column}:") 206: lines.append(f" • Total Responses:
{results['total_responses']}") 207: lines.append(f" • Positive:
{results['positive']} ({results['positive_pct']:.1f}%)") 208: lines.append(f" •
Negative: {results['negative']} ({results['negative_pct']:.1f}%)") 209:
lines.append(f" • Neutral: {results['neutral']}
({results['neutral_pct']:.1f}%)") 210: lines.append(f" • Average Sentiment
Score: {results['avg_score']:.2f}") 211: lines.append("") 212: else: 213:
lines.append("No text responses available for sentiment analysis.") 214:
lines.append("") 215: 216: return lines 217: 218: def
_generate_pattern_analysis(self, patterns: List[Dict[str, Any]]) -> List[str]:
219: """Generate the pattern analysis section.""" 220: lines = [ 221: "PATTERN
ANALYSIS", 222: "-" * 50, 223: "" 224: ] 225: 226: if patterns: 227:
lines.append("■ Detected Patterns and Correlations:") 228: lines.append("") 229:
230: # Group patterns by type 231: pattern_types = {} 232: for pattern in
patterns: 233: pattern_type = pattern.get('type', 'unknown') 234: if pattern_type
not in pattern_types: 235: pattern_types[pattern_type] = [] 236:
pattern_types[pattern_type].append(pattern) 237: 238: for pattern_type,
type_patterns in pattern_types.items(): 239: lines.append(f"■
{pattern_type.replace('_', ' ').title()} Patterns:") 240: for pattern in
type_patterns: 241: lines.append(f" • {pattern['description']}") 242:
lines.append(f" Confidence: {pattern['confidence']:.1f}%") 243: lines.append(f"
Sample Size: {pattern['sample_size']}") 244: lines.append("") 245: else: 246:
lines.append("No significant patterns detected in the survey data.") 247:
lines.append("") 248: 249: return lines 250:
```

```python
251: def _generate_statistical_analysis(self, summary: Dict[str, Any]) ->
List[str]: 252:     """Generate the statistical analysis section.""" 253:     lines = [
254:         "STATISTICAL ANALYSIS", 255:         "-" * 50, 256:         "" 257:     ] 258: 259:
    question_summaries = summary.get('question_summaries', {}) 260: 261:     if
question_summaries: 262:         lines.append("■ Response Distribution Analysis:") 263:
        lines.append("") 264: 265:         for question, q_summary in
question_summaries.items(): 266:             if isinstance(q_summary, dict): 267:
                lines.append(f"■ {question}:") 268:                 lines.append(f"  • Total Responses:
{q_summary.get('total_responses', 0)}") 269:                 lines.append(f"  • Response Rate:
{q_summary.get('response_rate', 0):.1f}%") 270: 271:                 top_responses =
q_summary.get('top_responses', []) 272:                 if top_responses: 273:                     lines.append("  •
Top Responses:") 274:                     for response, count in top_responses[:3]: 275:                         percentage =
(count / q_summary['total_responses']) * 100 276:                         lines.append(f"    - {response}:
{count} ({percentage:.1f}%)") 277: 278:                 lines.append("") 279:     else: 280:
        lines.append("No question response data available for statistical analysis.")
281:     lines.append("") 282: 283:     return lines 284: 285:     def
_generate_key_findings(self, summary: Dict[str, Any], sentiment: Dict[str, Any],
patterns: List[Dict[str, Any]]) -> List[str]: 286:     """Generate the key findings
section.""" 287:     lines = [ 288:         "KEY FINDINGS", 289:         "-" * 50, 290:         "" 291:     ]
292: 293:     findings = [] 294: 295:     # Demographic findings 296:     demographics =
summary.get('demographics', {}) 297:     if demographics: 298:         for demo_field,
breakdown in demographics.items(): 299:             if breakdown: 300:                 most_common =
max(breakdown.items(), key=lambda x: x[1])
```

```
301: findings.append(f"• {demo_field.title()}: {most_common[0]} is the most
common category ({most_common[1]} responses)") 302: 303: # Sentiment findings
304: if sentiment: 305: for column, results in sentiment.items(): 306: if
isinstance(results, dict): 307: dominant_sentiment =
self._get_dominant_sentiment(results) 308: findings.append(f"• {column}:
{dominant_sentiment} sentiment dominates the responses") 309: 310: # Pattern
findings 311: if patterns: 312: high_confidence_patterns = [p for p in patterns
if p.get('confidence', 0) >= 80] 313: if high_confidence_patterns: 314:
findings.append(f"• {len(high_confidence_patterns)} high-confidence patterns
detected in the data") 315: 316: if findings: 317: lines.extend(findings) 318:
else: 319: lines.append("No significant findings to report at this time.") 320:
321: lines.append("") 322: return lines 323: 324: def
_generate_recommendations(self, summary: Dict[str, Any], sentiment: Dict[str,
Any], patterns: List[Dict[str, Any]]) -> List[str]: 325: """Generate the
recommendations section.""" 326: lines = [ 327: "RECOMMENDATIONS", 328: "-" * 50,
329: "" 330: ] 331: 332: recommendations = [] 333: 334: # Data quality
recommendations 335: data_quality = summary.get('data_quality', {}) 336: if
data_quality: 337: completeness = data_quality.get('completeness', {}) 338:
low_completeness = [col for col, comp in completeness.items() if comp < 80] 339:
if low_completeness: 340: recommendations.append(f"• Improve data collection for
columns with low completeness: {', '.join(low_completeness)}") 341: 342: #
Sentiment-based recommendations 343: if sentiment: 344:
negative_sentiment_columns = [] 345: for column, results in sentiment.items():
346: if isinstance(results, dict) and results.get('negative_pct', 0) > 30: 347:
negative_sentiment_columns.append(column) 348: 349: if
negative_sentiment_columns: 350: recommendations.append(f"• Address concerns in
columns with high negative sentiment: {', '.join(negative_sentiment_columns)}")
```

```
351: 352: # Pattern-based recommendations 353: if patterns: 354:
high_confidence_patterns = [p for p in patterns if p.get('confidence', 0) >= 80]
355: if high_confidence_patterns: 356: recommendations.append(f"• Investigate
{len(high_confidence_patterns)} high-confidence patterns for actionable
insights") 357: 358: # General recommendations 359: recommendations.extend([ 360:
"• Consider conducting follow-up surveys to validate findings", 361: "• Implement
targeted improvements based on demographic insights", 362: "• Monitor sentiment
trends over time for continuous improvement", 363: "• Use statistical insights to
inform decision-making processes" 364: ]) 365: 366:
lines.extend(recommendations) 367: lines.append("") 368: 369: return lines 370:
371: def _generate_footer(self) -> List[str]: 372: """Generate the report
footer.""" 373: footer = [ 374: "=" * 80, 375: "REPORT END", 376: "=" * 80, 377:
"", 378: "This report was generated automatically by the Survey Data Analyzer.",
379: "For questions or additional analysis, please contact the development
team.", 380: "", 381: f"Report generated on: {self.current_datetime.strftime('%B
%d, %Y at %I:%M %p')}", 382: "=" * 80 383: ] 384: return footer 385: 386: def
_assess_data_quality(self, data_quality: Dict[str, Any]) -> str: 387: """Assess
overall data quality.""" 388: if not data_quality: 389: return "Unknown" 390:
391: completeness = data_quality.get('completeness', {}) 392: if not
completeness: 393: return "Unknown" 394: 395: avg_completeness =
sum(completeness.values()) / len(completeness) 396: 397: if avg_completeness >=
90: 398: return "Excellent" 399: elif avg_completeness >= 80: 400: return "Good"
```

```python
401: elif avg_completeness >= 70: 402: return "Fair" 403: else: 404: return
"Poor" 405: 406: def _get_dominant_sentiment(self, sentiment_results: Dict[str,
Any]) -> str: 407: """Get the dominant sentiment from results.""" 408:
positive_pct = sentiment_results.get('positive_pct', 0) 409: negative_pct =
sentiment_results.get('negative_pct', 0) 410: neutral_pct =
sentiment_results.get('neutral_pct', 0) 411: 412: if positive_pct > negative_pct
and positive_pct > neutral_pct: 413: return "Positive" 414: elif negative_pct >
positive_pct and negative_pct > neutral_pct: 415: return "Negative" 416: else:
417: return "Neutral" 418: 419: def _write_report_to_file(self, content: str,
file_path: str) -> bool: 420: """Write the report content to a file.""" 421: try:
422: # Ensure directory exists 423: directory = os.path.dirname(file_path) 424:
if directory and not os.path.exists(directory): 425: os.makedirs(directory) 426:
427: # Write content to file 428: with open(file_path, 'w', encoding='utf-8') as
file: 429: file.write(content) 430: 431: return True 432: 433: except Exception
as e: 434: print(f"Error writing report to file: {str(e)}") 435: return False
436: 437: def generate_summary_report(self, summary_data: Dict[str, Any],
file_path: str = "summary_report.txt") -> bool: 438: """Generate a simplified
summary report.""" 439: try: 440: lines = [ 441: "SURVEY SUMMARY REPORT", 442:
"=" * 50, 443: f"Generated: {self.current_datetime.strftime('%B %d, %Y')}", 444:
"", 445: f"Total Responses: {summary_data.get('total_responses', 0)}", 446:
f"Response Rate: {summary_data.get('response_rate', 0):.1f}%", 447: "", 448:
"Key Metrics:", 449: "• Data quality assessment", 450: "• Demographic
breakdowns",
```

```
451: "• Response distributions", 452: "• Pattern detection results", 453: "",
454: "=" * 50 455: ] 456: 457: content = "\n".join(lines) 458: return
self._write_report_to_file(content, file_path) 459: 460: except Exception as e:
461: print(f"Error generating summary report: {str(e)}") 462: return False
```

# 9. Utilities Module (utils.py)

The utils.py file provides utility functions for formatting, validation, and common operations used throughout the application. It includes display formatting, input validation, and helper functions.

## Utilities Module - utils.py

```
1: #!/usr/bin/env python3 2: """ 3: Utilities Module 4: Description: Provides
utility functions for the Survey Data Analyzer CLI, 5: including display
formatting, input validation, and common operations. 6: """ 7: 8: import os 9:
import sys 10: from typing import List, Dict, Any, Optional 11: 12: 13: def
clear_screen(): 14: """Clear the terminal screen.""" 15: os.system('cls' if
os.name == 'nt' else 'clear') 16: 17: 18: def print_header(title: str): 19:
"""Print a formatted header.""" 20: print("\n" + "=" * 60) 21: print(f"TARGET:
{title}") 22: print("=" * 60) 23: 24: 25: def print_menu(options: List[str]): 26:
"""Print a numbered menu.""" 27: print("\nAvailable Options:") 28: for i, option
in enumerate(options, 1): 29: print(f" {i}. {option}") 30: 31: 32: def
print_success(message: str): 33: """Print a success message.""" 34:
print(f"SUCCESS: {message}") 35: 36: 37: def print_error(message: str): 38:
"""Print an error message.""" 39: print(f"ERROR: {message}") 40: 41: 42: def
print_warning(message: str): 43: """Print a warning message.""" 44:
print(f"WARNING: {message}") 45: 46: 47: def print_info(message: str): 48:
"""Print an info message.""" 49: print(f"INFO: {message}") 50:
```

```
51:
52: def validate_file_path(file_path: str) -> bool:
53:     """
54:     Validate if a file path exists and is accessible.
55:
56:     Args:
57:         file_path: Path to the file to validate
58:
59:     Returns:
60:         True if file is valid, False otherwise
61:     """
62:     if not file_path or not isinstance(file_path, str):
63:         return False
64:
65:     file_path = file_path.strip()
66:     if not file_path:
67:         return False
68:
69:     return os.path.exists(file_path) and os.path.isfile(file_path)
70:
71:
72: def validate_csv_file(file_path: str) -> bool:
73:     """
74:     Validate if a file is a valid CSV file.
75:
76:     Args:
77:         file_path: Path to the CSV file
78:
79:     Returns:
80:         True if file is a valid CSV, False otherwise
81:     """
82:     if not validate_file_path(file_path):
83:         return False
84:
85:     # Check file extension
86:     if not file_path.lower().endswith('.csv'):
87:         return False
88:
89:     # Check file size (not empty and not too large)
90:     try:
91:         file_size = os.path.getsize(file_path)
92:         if file_size == 0:
93:             return False
94:         if file_size > 50 * 1024 * 1024:  # 50MB limit
95:             return False
96:     except OSError:
97:         return False
98:
99:     return True
100:
```

```python
101: 102: def format_number(number: float, decimal_places: int = 2) -> str: 103:
""" 104: Format a number with specified decimal places. 105: 106: Args: 107:
number: Number to format 108: decimal_places: Number of decimal places to show
109: 110: Returns: 111: Formatted number string 112: """ 113: return
f"{number:.{decimal_places}f}" 114: 115: 116: def format_percentage(value:
float, total: float) -> str: 117: """ 118: Format a percentage value. 119: 120:
Args: 121: value: The value to calculate percentage for 122: total: The total
value 123: 124: Returns: 125: Formatted percentage string 126: """ 127: if total
== 0: 128: return "0.0%" 129: percentage = (value / total) * 100 130: return
f"{percentage:.1f}%" 131: 132: 133: def format_table(data: List[List[str]],
headers: List[str] = None) -> str: 134: """ 135: Format data as a table. 136:
137: Args: 138: data: List of rows, each row is a list of values 139: headers:
Optional list of column headers 140: 141: Returns: 142: Formatted table string
143: """ 144: if not data: 145: return "No data to display" 146: 147: # Determine
column widths 148: if headers: 149: all_rows = [headers] + data 150: else:
```

```python
151: all_rows = data 152: 153: col_widths = [] 154: for col in
range(len(all_rows[0])): 155: max_width = max(len(str(row[col])) for row in
all_rows) 156: col_widths.append(max_width) 157: 158: # Build table 159: lines =
[] 160: 161: # Header 162: if headers: 163: header_line = " |
".join(f"{headers[i]:<{col_widths[i]}}" for i in range(len(headers))) 164:
lines.append(header_line) 165: lines.append("-" * len(header_line)) 166: 167: #
Data rows 168: for row in data: 169: row_line = " |
".join(f"{row[i]:<{col_widths[i]}}" for i in range(len(row))) 170:
lines.append(row_line) 171: 172: return "\n".join(lines) 173: 174: 175: def
get_user_input(prompt: str, default: str = None) -> str: 176: """ 177: Get user
input with optional default value. 178: 179: Args: 180: prompt: Input prompt to
display 181: default: Optional default value 182: 183: Returns: 184: User input
string 185: """ 186: if default: 187: user_input = input(f"{prompt} (default:
{default}): ").strip() 188: return user_input if user_input else default 189:
else: 190: return input(f"{prompt}: ").strip() 191: 192: 193: def
confirm_action(prompt: str = "Are you sure?") -> bool: 194: """ 195: Get user
confirmation for an action. 196: 197: Args: 198: prompt: Confirmation prompt 199:
200: Returns:
```

```python
201: True if user confirms, False otherwise 202: """ 203: response =
input(f"{prompt} (y/N): ").strip().lower() 204: return response in ['y', 'yes']
205: 206: 207: def display_progress(current: int, total: int, description: str =
"Processing"): 208: """ 209: Display a progress bar. 210: 211: Args: 212:
current: Current progress value 213: total: Total value 214: description:
Description of the operation 215: """ 216: if total == 0: 217: return 218: 219:
percentage = (current / total) * 100 220: bar_length = 30 221: filled_length =
int(bar_length * current // total) 222: bar = '■' * filled_length + '-' *
(bar_length - filled_length) 223: 224: print(f"\r{description}: |{bar}|
{percentage:.1f}% ({current}/{total})", end='') 225: 226: if current == total:
227: print() # New line when complete 228: 229: 230: def safe_divide(numerator:
float, denominator: float, default: float = 0.0) -> float: 231: """ 232: Safely
divide two numbers, returning default if denominator is zero. 233: 234: Args:
235: numerator: The numerator 236: denominator: The denominator 237: default:
Default value if division by zero 238: 239: Returns: 240: Result of division or
default value 241: """ 242: try: 243: return numerator / denominator if
denominator != 0 else default 244: except (TypeError, ValueError): 245: return
default 246: 247: 248: def truncate_text(text: str, max_length: int = 50) -> str:
249: """ 250: Truncate text to a maximum length.
```

```
251:
252: Args:
253:     text: Text to truncate
254:     max_length: Maximum length
255:
256: Returns:
257:     Truncated text
258: """
259: if len(text) <= max_length:
260:     return text
261: return text[:max_length-3] + "..."
262:
263:
264: def format_file_size(size_bytes: int) -> str:
265:     """
266:     Format file size in human-readable format.
267:
268:     Args:
269:         size_bytes: Size in bytes
270:
271:     Returns:
272:         Formatted size string
273:     """
274:     if size_bytes == 0:
275:         return "0 B"
276:
277:     size_names = ["B", "KB", "MB", "GB"]
278:     i = 0
279:     while size_bytes >= 1024 and i < len(size_names) - 1:
280:         size_bytes /= 1024.0
281:         i += 1
282:
283:     return f"{size_bytes:.1f} {size_names[i]}"
284:
285:
286: def validate_age(age_str: str) -> bool:
287:     """
288:     Validate if a string represents a valid age.
289:
290:     Args:
291:         age_str: String to validate as age
292:
293:     Returns:
294:         True if valid age, False otherwise
295:     """
296:     try:
297:         age = int(age_str)
298:         return 0 <= age <= 120
299:     except (ValueError, TypeError):
300:         return False
```

```
301: 302: 303: def validate_gender(gender_str: str) -> bool: 304: """ 305:
Validate if a string represents a valid gender. 306: 307: Args: 308: gender_str:
String to validate as gender 309: 310: Returns: 311: True if valid gender, False
otherwise 312: """ 313: if not gender_str: 314: return False 315: 316:
valid_genders = [ 317: 'male', 'female', 'm', 'f', 'other', 'prefer not to say',
318: 'Male', 'Female', 'M', 'F', 'Other', 'Prefer not to say' 319: ] 320: 321:
return gender_str.strip() in valid_genders 322: 323: 324: def
normalize_text(text: str) -> str: 325: """ 326: Normalize text by removing extra
whitespace and converting to lowercase. 327: 328: Args: 329: text: Text to
normalize 330: 331: Returns: 332: Normalized text 333: """ 334: if not text: 335:
return "" 336: 337: # Remove extra whitespace and convert to lowercase 338:
normalized = " ".join(text.strip().split()).lower() 339: return normalized 340:
341: 342: def count_words(text: str) -> int: 343: """ 344: Count words in text.
345: 346: Args: 347: text: Text to count words in 348: 349: Returns: 350: Number
of words
```

```python
351:     """
352:     if not text:
353:         return 0
354:
355:     return len(text.split())
356:
357:
358: def get_file_extension(file_path: str) -> str:
359:     """
360:     Get file extension from file path.
361:
362:     Args:
363:         file_path: Path to file
364:
365:     Returns:
366:         File extension (including dot)
367:     """
368:     return os.path.splitext(file_path)[1].lower()
369:
370:
371: def ensure_directory_exists(directory_path: str) -> bool:
372:     """
373:     Ensure a directory exists, creating it if necessary.
374:
375:     Args:
376:         directory_path: Path to directory
377:
378:     Returns:
379:         True if directory exists or was created, False otherwise
380:     """
381:     try:
382:         if not os.path.exists(directory_path):
383:             os.makedirs(directory_path)
384:         return True
385:     except Exception:
386:         return False
387:
388:
389: def is_valid_filename(filename: str) -> bool:
390:     """
391:     Check if a filename is valid.
392:
393:     Args:
394:         filename: Filename to validate
395:
396:     Returns:
397:         True if valid filename, False otherwise
398:     """
399:     if not filename:
400:         return False
```

```
401: 402:     # Check for invalid characters 403:         invalid_chars = '<>:"/\\|?*' 404:
return not any(char in filename for char in invalid_chars) 405: 406: 407: def
format_duration(seconds: float) -> str: 408: """ 409:     Format duration in seconds
to human-readable format. 410: 411:     Args: 412:         seconds: Duration in seconds 413:
414:     Returns: 415:         Formatted duration string 416:     """ 417:     if seconds < 60: 418:
return f"{seconds:.1f} seconds" 419:     elif seconds < 3600: 420:         minutes = seconds
/ 60 421:         return f"{minutes:.1f} minutes" 422:     else: 423:         hours = seconds / 3600
424:         return f"{hours:.1f} hours" 425: 426: 427: def print_separator(char: str =
"-", length: int = 60): 428:     """Print a separator line.""" 429:     print(char *
length) 430: 431: 432: def print_bullet_list(items: List[str], indent: int = 2):
433:     """ 434:     Print a bullet list. 435: 436:     Args: 437:         items: List of items to
print 438:         indent: Number of spaces to indent 439:     """ 440:     indent_str = " " *
indent 441:     for item in items: 442:         print(f"{indent_str}• {item}") 443: 444: 445:
def print_key_value_pairs(data: Dict[str, Any], indent: int = 2): 446:     """ 447:
Print key-value pairs in a formatted way. 448: 449:     Args: 450:         data: Dictionary
of key-value pairs
```

```
451:     indent: Number of spaces to indent
452:     """
453:     indent_str = " " * indent
454:     for key, value in data.items():
455:         print(f"{indent_str}{key}: {value}")
```

# 10. Test Suite

The test suite provides comprehensive unit testing for all major components of the application. Each module has corresponding test files that validate functionality, edge cases, and error handling.

## Data Loader Tests - test_data_loader.py

```
1: #!/usr/bin/env python3 2: """ 3: Unit Tests for Data Loader Module 4: Author:
Student Developer 5: Description: Comprehensive unit tests for CSV data loading
and validation functionality. 6: """ 7: 8: import unittest 9: import tempfile 10:
import os 11: import csv 12: import sys 13: import os 14:
sys.path.append(os.path.dirname(os.path.dirname(os.path.abspath(__file__)))) 15:
from data_loader import DataLoader 16: 17: 18: class
TestDataLoader(unittest.TestCase): 19: """Test cases for the DataLoader
class.""" 20: 21: def setUp(self): 22: """Set up test fixtures.""" 23:
self.data_loader = DataLoader() 24: self.temp_dir = tempfile.mkdtemp() 25: 26:
def tearDown(self): 27: """Clean up test fixtures.""" 28: # Clean up temporary
files 29: for file in os.listdir(self.temp_dir): 30:
os.remove(os.path.join(self.temp_dir, file)) 31: os.rmdir(self.temp_dir) 32: 33:
def test_validate_file_existing(self): 34: """Test file validation with existing
file.""" 35: # Create a temporary file 36: temp_file =
os.path.join(self.temp_dir, "test.csv") 37: with open(temp_file, 'w') as f: 38:
f.write("test data") 39: 40: result = self.data_loader._validate_file(temp_file)
41: self.assertTrue(result) 42: 43: def test_validate_file_nonexistent(self):
44: """Test file validation with non-existent file.""" 45: result =
self.data_loader._validate_file("nonexistent_file.csv") 46:
self.assertFalse(result) 47: 48: def test_validate_file_empty(self): 49: """Test
file validation with empty file.""" 50: temp_file = os.path.join(self.temp_dir,
"empty.csv")
```

```python
51: with open(temp_file, 'w') as f:
52:     pass  # Create empty file
53:
54:     result = self.data_loader._validate_file(temp_file)
55:     self.assertFalse(result)
56:
57:     def test_validate_file_too_large(self):
58:         """Test file validation with file too large."""
59:         temp_file = os.path.join(self.temp_dir, "large.csv")
60:         # Create a file larger than 50MB
61:         with open(temp_file, 'w') as f:
62:             f.write("x" * (51 * 1024 * 1024))  # 51MB
63:
64:         result = self.data_loader._validate_file(temp_file)
65:         self.assertFalse(result)
66:
67:     def test_read_csv_with_encoding_utf8(self):
68:         """Test CSV reading with UTF-8 encoding."""
69:         temp_file = os.path.join(self.temp_dir, "test_utf8.csv")
70:
71:         # Create test CSV data
72:         test_data = [
73:             {'name': 'John', 'age': '25', 'city': 'New York'},
74:             {'name': 'Jane', 'age': '30', 'city': 'Los Angeles'},
75:             {'name': 'Bob', 'age': '35', 'city': 'Chicago'}
76:         ]
77:
78:         with open(temp_file, 'w', newline='', encoding='utf-8') as f:
79:             writer = csv.DictWriter(f, fieldnames=['name', 'age', 'city'])
80:             writer.writeheader()
81:             writer.writerows(test_data)
82:
83:         result = self.data_loader._read_csv_with_encoding(temp_file, 'utf-8')
84:         self.assertIsNotNone(result)
85:         self.assertEqual(len(result), 3)
86:         self.assertEqual(result[0]['name'], 'John')
87:
88:     def test_read_csv_with_encoding_latin1(self):
89:         """Test CSV reading with Latin-1 encoding."""
90:         temp_file = os.path.join(self.temp_dir, "test_latin1.csv")
91:
92:         # Create test CSV data
93:         test_data = [
94:             {'name': 'José', 'age': '25', 'city': 'Madrid'},
95:             {'name': 'François', 'age': '30', 'city': 'Paris'}
96:         ]
97:
98:         with open(temp_file, 'w', newline='', encoding='latin-1') as f:
99:             writer = csv.DictWriter(f, fieldnames=['name', 'age', 'city'])
100:            writer.writeheader()
```

```python
101: writer.writerows(test_data) 102: 103: result =
self.data_loader._read_csv_with_encoding(temp_file, 'latin-1') 104:
self.assertIsNotNone(result) 105: self.assertEqual(len(result), 2) 106: 107: def
test_clean_and_validate_data_valid(self): 108: """Test data cleaning and
validation with valid data.""" 109: test_data = [ 110: {'age': '25', 'gender':
'Male', 'region': 'North'}, 111: {'age': '30', 'gender': 'Female', 'region':
'South'}, 112: {'age': '35', 'gender': 'Male', 'region': 'East'} 113: ] 114: 115:
result = self.data_loader._clean_and_validate_data(test_data) 116:
self.assertIsNotNone(result) 117: self.assertEqual(len(result), 3) 118: 119: def
test_clean_and_validate_data_invalid_age(self): 120: """Test data cleaning with
invalid age data.""" 121: test_data = [ 122: {'age': '25', 'gender': 'Male',
'region': 'North'}, 123: {'age': '150', 'gender': 'Female', 'region': 'South'}, #
Invalid age 124: {'age': '35', 'gender': 'Male', 'region': 'East'} 125: ] 126:
127: result = self.data_loader._clean_and_validate_data(test_data) 128:
self.assertIsNotNone(result) 129: # Should filter out invalid age 130:
self.assertEqual(len(result), 2) 131: 132: def
test_clean_and_validate_data_invalid_gender(self): 133: """Test data cleaning
with invalid gender data.""" 134: test_data = [ 135: {'age': '25', 'gender':
'Male', 'region': 'North'}, 136: {'age': '30', 'gender': 'Invalid', 'region':
'South'}, # Invalid gender 137: {'age': '35', 'gender': 'Female', 'region':
'East'} 138: ] 139: 140: result =
self.data_loader._clean_and_validate_data(test_data) 141:
self.assertIsNotNone(result) 142: # Should filter out invalid gender 143:
self.assertEqual(len(result), 2) 144: 145: def
test_clean_and_validate_data_empty(self): 146: """Test data cleaning with empty
data.""" 147: result = self.data_loader._clean_and_validate_data([]) 148:
self.assertIsNone(result) 149: 150: def test_validate_row_valid(self):
```

```python
151:     """Test row validation with valid data."""
152:     valid_row = {'age': '25', 'gender': 'Male', 'region': 'North'}
153:     result = self.data_loader._validate_row(valid_row, 1)
154:     self.assertTrue(result)
155:
156: def test_validate_row_invalid_age(self):
157:     """Test row validation with invalid age."""
158:     invalid_row = {'age': '150', 'gender': 'Male', 'region': 'North'}
159:     result = self.data_loader._validate_row(invalid_row, 1)
160:     self.assertFalse(result)
161:
162: def test_validate_row_invalid_gender(self):
163:     """Test row validation with invalid gender."""
164:     invalid_row = {'age': '25', 'gender': 'Invalid', 'region': 'North'}
165:     result = self.data_loader._validate_row(invalid_row, 1)
166:     self.assertFalse(result)
167:
168: def test_validate_row_insufficient_fields(self):
169:     """Test row validation with insufficient fields."""
170:     invalid_row = {'age': '25'}  # Only one field
171:     result = self.data_loader._validate_row(invalid_row, 1)
172:     self.assertFalse(result)
173:
174: def test_get_data_summary(self):
175:     """Test data summary generation."""
176:     test_data = [
177:         {'age': '25', 'gender': 'Male', 'region': 'North'},
178:         {'age': '30', 'gender': 'Female', 'region': 'South'},
179:         {'age': '35', 'gender': 'Male', 'region': 'East'}
180:     ]
181:
182:     summary = self.data_loader.get_data_summary(test_data)
183:
184:     self.assertIn('total_rows', summary)
185:     self.assertIn('columns', summary)
186:     self.assertIn('missing_values', summary)
187:     self.assertIn('unique_values', summary)
188:
189:     self.assertEqual(summary['total_rows'], 3)
190:     self.assertEqual(len(summary['columns']), 3)
191:
192: def test_export_sample_data(self):
193:     """Test sample data export functionality."""
194:     temp_file = os.path.join(self.temp_dir, "sample_export.csv")
195:
196:     result = self.data_loader.export_sample_data(temp_file)
197:     self.assertTrue(result)
198:     self.assertTrue(os.path.exists(temp_file))
199:
200:     # Verify the exported file has content
```

```python
201: with open(temp_file, 'r', encoding='utf-8') as f:
202:     content = f.read()
203: self.assertIn('age', content)
204: self.assertIn('gender', content)
205: self.assertIn('region', content)
206:
207: def test_load_csv_integration(self):
208:     """Test complete CSV loading integration."""
209:     temp_file = os.path.join(self.temp_dir, "integration_test.csv")
210:
211:     # Create test CSV data
212:     test_data = [
213:         {'age': '25', 'gender': 'Male', 'region': 'North', 'satisfaction': 'High'},
214:         {'age': '30', 'gender': 'Female', 'region': 'South', 'satisfaction': 'Medium'},
215:         {'age': '35', 'gender': 'Male', 'region': 'East', 'satisfaction': 'Low'}
216:     ]
217:
218:     with open(temp_file, 'w', newline='', encoding='utf-8') as f:
219:         writer = csv.DictWriter(f, fieldnames=['age', 'gender', 'region', 'satisfaction'])
220:         writer.writeheader()
221:         writer.writerows(test_data)
222:
223:     result = self.data_loader.load_csv(temp_file)
224:     self.assertIsNotNone(result)
225:     self.assertEqual(len(result), 3)
226:     self.assertEqual(result[0]['age'], '25')
227:     self.assertEqual(result[0]['gender'], 'Male')
228:
229:
230: if __name__ == '__main__':
231:     unittest.main()
```

# Sentiment Analyzer Tests - test_sentiment_analyzer.py

```python
1: #!/usr/bin/env python3
2: """
3: Unit Tests for Sentiment Analyzer Module
4: Author: Student Developer
5: Description: Comprehensive unit tests for sentiment analysis functionality.
6: """
7:
8: import unittest
9: import sys
10: import os
11: sys.path.append(os.path.dirname(os.path.dirname(os.path.abspath(__file__))))
12: from sentiment_analyzer import SentimentAnalyzer
13:
14:
15: class TestSentimentAnalyzer(unittest.TestCase):
16:     """Test cases for the SentimentAnalyzer class."""
17:
18:     def setUp(self):
19:         """Set up test fixtures."""
20:         self.sentiment_analyzer = SentimentAnalyzer()
21:         self.test_data = [
22:             {'feedback': 'This is excellent! I love it.', 'satisfaction': 'Very Satisfied'},
23:             {'feedback': 'Terrible experience, very bad service.', 'satisfaction': 'Dissatisfied'},
24:             {'feedback': 'It was okay, nothing special.', 'satisfaction': 'Neutral'},
25:             {'feedback': 'Great product and amazing support!', 'satisfaction': 'Very Satisfied'},
26:             {'feedback': 'Poor quality and bad service.', 'satisfaction': 'Dissatisfied'}
27:         ]
28:
29:     def test_initialization(self):
30:         """Test SentimentAnalyzer initialization."""
31:         self.assertIsInstance(self.sentiment_analyzer.positive_keywords, dict)
32:         self.assertIsInstance(self.sentiment_analyzer.negative_keywords, dict)
33:         self.assertIsInstance(self.sentiment_analyzer.neutral_keywords, dict)
34:         self.assertIsInstance(self.sentiment_analyzer.negation_words, set)
35:         self.assertIsInstance(self.sentiment_analyzer.intensifier_words, dict)
36:
37:         # Check that keywords are loaded
38:         self.assertGreater(len(self.sentiment_analyzer.positive_keywords), 0)
39:         self.assertGreater(len(self.sentiment_analyzer.negative_keywords), 0)
40:
41:     def test_analyze_text_positive(self):
42:         """Test sentiment analysis with positive text."""
43:         text = "This is excellent! I love the product."
44:         result = self.sentiment_analyzer.analyze_text(text)
45:
46:         self.assertIsInstance(result, dict)
47:         self.assertIn('sentiment', result)
48:         self.assertIn('score', result)
49:         self.assertIn('positive_words', result)
50:         self.assertIn('negative_words', result)
```

```python
51:     self.assertIn('confidence', result)
52:
53:     self.assertEqual(result['sentiment'], 'positive')
54:     self.assertGreater(result['score'], 0)
55:     self.assertGreater(len(result['positive_words']), 0)
56:
57: def test_analyze_text_negative(self):
58:     """Test sentiment analysis with negative text."""
59:     text = "This is terrible! I hate the service."
60:     result = self.sentiment_analyzer.analyze_text(text)
61:
62:     self.assertEqual(result['sentiment'], 'negative')
63:     self.assertLess(result['score'], 0)
64:     self.assertGreater(len(result['negative_words']), 0)
65:
66: def test_analyze_text_neutral(self):
67:     """Test sentiment analysis with neutral text."""
68:     text = "It was okay, nothing special."
69:     result = self.sentiment_analyzer.analyze_text(text)
70:
71:     self.assertEqual(result['sentiment'], 'neutral')
72:     self.assertAlmostEqual(result['score'], 0, places=1)
73:
74: def test_analyze_text_with_negation(self):
75:     """Test sentiment analysis with negation."""
76:     text = "This is not good at all."
77:     result = self.sentiment_analyzer.analyze_text(text)
78:
79:     # Should be negative due to negation
80:     self.assertEqual(result['sentiment'], 'negative')
81:     self.assertLess(result['score'], 0)
82:
83: def test_analyze_text_with_intensifier(self):
84:     """Test sentiment analysis with intensifier."""
85:     text = "This is very excellent!"
86:     result = self.sentiment_analyzer.analyze_text(text)
87:
88:     self.assertEqual(result['sentiment'], 'positive')
89:     self.assertGreater(result['score'], 0)
90:
91: def test_analyze_text_empty(self):
92:     """Test sentiment analysis with empty text."""
93:     result = self.sentiment_analyzer.analyze_text("")
94:
95:     self.assertEqual(result['sentiment'], 'neutral')
96:     self.assertEqual(result['score'], 0)
97:     self.assertEqual(len(result['positive_words']), 0)
98:     self.assertEqual(len(result['negative_words']), 0)
99:
100: def test_analyze_text_none(self):
```

```python
101:     """Test sentiment analysis with None text."""
102:     result = self.sentiment_analyzer.analyze_text(None)
103:
104:     self.assertEqual(result['sentiment'], 'neutral')
105:     self.assertEqual(result['score'], 0)
106:
107: def test_analyze_column(self):
108:     """Test sentiment analysis for a column."""
109:     result = self.sentiment_analyzer.analyze_column(self.test_data, 'feedback')
110:
111:     self.assertIsInstance(result, dict)
112:     self.assertIn('positive', result)
113:     self.assertIn('negative', result)
114:     self.assertIn('neutral', result)
115:     self.assertIn('positive_pct', result)
116:     self.assertIn('negative_pct', result)
117:     self.assertIn('neutral_pct', result)
118:     self.assertIn('avg_score', result)
119:     self.assertIn('total_responses', result)
120:
121:     # Check that percentages sum to approximately 100
122:     total_pct = result['positive_pct'] + result['negative_pct'] + result['neutral_pct']
123:     self.assertAlmostEqual(total_pct, 100.0, places=1)
124:
125: def test_analyze_column_empty_data(self):
126:     """Test sentiment analysis with empty data."""
127:     result = self.sentiment_analyzer.analyze_column([], 'feedback')
128:
129:     self.assertEqual(result['positive'], 0)
130:     self.assertEqual(result['negative'], 0)
131:     self.assertEqual(result['neutral'], 0)
132:     self.assertEqual(result['total_responses'], 0)
133:
134: def test_analyze_column_missing_column(self):
135:     """Test sentiment analysis with missing column."""
136:     result = self.sentiment_analyzer.analyze_column(self.test_data, 'nonexistent_column')
137:
138:     self.assertEqual(result['positive'], 0)
139:     self.assertEqual(result['negative'], 0)
140:     self.assertEqual(result['neutral'], 0)
141:     self.assertEqual(result['total_responses'], 0)
142:
143: def test_analyze_all_text_columns(self):
144:     """Test sentiment analysis for all text columns."""
145:     result = self.sentiment_analyzer.analyze_all_text_columns(self.test_data)
146:
147:     self.assertIsInstance(result, dict)
148:     self.assertIn('feedback', result)
149:
150:     feedback_result = result['feedback']
```

```python
151: self.assertIn('positive', feedback_result) 152: self.assertIn('negative',
feedback_result) 153: self.assertIn('neutral', feedback_result) 154: 155: def
test_clean_text(self): 156: """Test text cleaning functionality.""" 157:
dirty_text = " This is a TEST!!! " 158: cleaned =
self.sentiment_analyzer._clean_text(dirty_text) 159: 160:
self.assertEqual(cleaned, "this is a test") 161: 162: def
test_extract_words(self): 163: """Test word extraction.""" 164: text = "This is a
test sentence." 165: words = self.sentiment_analyzer._extract_words(text) 166:
167: self.assertIsInstance(words, list) 168: self.assertEqual(len(words), 5)
169: self.assertIn('This', words) 170: self.assertIn('is', words) 171:
self.assertIn('a', words) 172: self.assertIn('test', words) 173:
self.assertIn('sentence.', words) 174: 175: def test_is_negated(self): 176:
"""Test negation detection.""" 177: words = ['This', 'is', 'not', 'good'] 178:
179: # Check if 'good' is negated 180: is_negated =
self.sentiment_analyzer._is_negated(words, 3) # Index of 'good' 181:
self.assertTrue(is_negated) 182: 183: # Check if 'This' is negated (should not
be) 184: is_negated = self.sentiment_analyzer._is_negated(words, 0) # Index of
'This' 185: self.assertFalse(is_negated) 186: 187: def
test_categorize_sentiment(self): 188: """Test sentiment categorization.""" 189:
# Test positive sentiment 190: sentiment =
self.sentiment_analyzer._categorize_sentiment(2.0) 191:
self.assertEqual(sentiment, 'positive') 192: 193: # Test negative sentiment 194:
sentiment = self.sentiment_analyzer._categorize_sentiment(-2.0) 195:
self.assertEqual(sentiment, 'negative') 196: 197: # Test neutral sentiment 198:
sentiment = self.sentiment_analyzer._categorize_sentiment(0.5) 199:
self.assertEqual(sentiment, 'neutral') 200:
```

```python
201:         sentiment = self.sentiment_analyzer._categorize_sentiment(-0.5)
202:         self.assertEqual(sentiment, 'neutral')
203:
204:     def test_get_sentiment_summary(self):
205:         """Test sentiment summary generation."""
206:         # Create sentiment results
207:         sentiment_results = [
208:             {'sentiment': 'positive', 'score': 2.0, 'positive_words': ['excellent'], 'negative_words': [], 'confidence': 0.8},
209:             {'sentiment': 'negative', 'score': -1.5, 'positive_words': [], 'negative_words': ['terrible'], 'confidence': 0.7},
210:             {'sentiment': 'neutral', 'score': 0.0, 'positive_words': [], 'negative_words': [], 'confidence': 0.5}
211:         ]
212:
213:         summary = self.sentiment_analyzer.get_sentiment_summary(sentiment_results)
214:
215:         self.assertIsInstance(summary, dict)
216:         self.assertIn('total_responses', summary)
217:         self.assertIn('sentiment_distribution', summary)
218:         self.assertIn('top_positive_words', summary)
219:         self.assertIn('top_negative_words', summary)
220:         self.assertIn('average_confidence', summary)
221:         self.assertIn('average_score', summary)
222:
223:         self.assertEqual(summary['total_responses'], 3)
224:         self.assertEqual(summary['sentiment_distribution']['positive']['count'], 1)
225:         self.assertEqual(summary['sentiment_distribution']['negative']['count'], 1)
226:         self.assertEqual(summary['sentiment_distribution']['neutral']['count'], 1)
227:
228:     def test_export_sentiment_report(self):
229:         """Test sentiment report export."""
230:         sentiment_results = {
231:             'feedback': {
232:                 'positive': 2,
233:                 'negative': 2,
234:                 'neutral': 1,
235:                 'positive_pct': 40.0,
236:                 'negative_pct': 40.0,
237:                 'neutral_pct': 20.0,
238:                 'avg_score': 0.2,
239:                 'total_responses': 5
240:             }
241:         }
242:
243:         # Test successful export
244:         result = self.sentiment_analyzer.export_sentiment_report(sentiment_results, "test_sentiment_report.txt")
245:         self.assertTrue(result)
246:
247:         # Clean up
248:         import os
249:         if os.path.exists("test_sentiment_report.txt"):
250:             os.remove("test_sentiment_report.txt")
```

```python
251:
252:     def test_analyze_text_complex_sentence(self):
253:         """Test sentiment analysis with complex sentence."""
254:         text = "This product is absolutely fantastic and I would highly recommend it to everyone!"
255:         result = self.sentiment_analyzer.analyze_text(text)
256:
257:         self.assertEqual(result['sentiment'], 'positive')
258:         self.assertGreater(result['score'], 0)
259:         self.assertGreater(len(result['positive_words']), 0)
260:
261:     def test_analyze_text_mixed_sentiment(self):
262:         """Test sentiment analysis with mixed sentiment."""
263:         text = "The product is good but the service is terrible."
264:         result = self.sentiment_analyzer.analyze_text(text)
265:
266:         # Should have both positive and negative words
267:         self.assertGreater(len(result['positive_words']), 0)
268:         self.assertGreater(len(result['negative_words']), 0)
269:
270:     def test_analyze_text_no_sentiment_words(self):
271:         """Test sentiment analysis with no sentiment words."""
272:         text = "The product arrived on time and was delivered to the correct address."
273:         result = self.sentiment_analyzer.analyze_text(text)
274:
275:         self.assertEqual(result['sentiment'], 'neutral')
276:         self.assertEqual(len(result['positive_words']), 0)
277:         self.assertEqual(len(result['negative_words']), 0)
278:
279:
280: if __name__ == '__main__':
281:     unittest.main()
```

## Statistical Analyzer Tests - test_stats_analyzer.py

```python
1: #!/usr/bin/env python3
2: """
3: Unit Tests for Statistical Analyzer Module
4: Author: Student Developer
5: Description: Comprehensive unit tests for statistical analysis functionality.
6: """
7:
8: import unittest
9: import sys
10: import os
11: sys.path.append(os.path.dirname(os.path.dirname(os.path.abspath(__file__))))
12: from stats_analyzer import StatsAnalyzer
13:
14:
15: class TestStatsAnalyzer(unittest.TestCase):
16:     """Test cases for the StatsAnalyzer class."""
17:
18:     def setUp(self):
19:         """Set up test fixtures."""
20:         self.test_data = [
21:             {'age': '25', 'gender': 'Male', 'region': 'North', 'satisfaction': 'High'},
22:             {'age': '30', 'gender': 'Female', 'region': 'South', 'satisfaction': 'Medium'},
23:             {'age': '35', 'gender': 'Male', 'region': 'East', 'satisfaction': 'Low'},
24:             {'age': '28', 'gender': 'Female', 'region': 'North', 'satisfaction': 'High'},
25:             {'age': '42', 'gender': 'Male', 'region': 'South', 'satisfaction': 'Medium'},
26:             {'age': '33', 'gender': 'Female', 'region': 'East', 'satisfaction': 'Low'}
27:         ]
28:         self.stats_analyzer = StatsAnalyzer(self.test_data)
29:
30:     def test_initialization(self):
31:         """Test StatsAnalyzer initialization."""
32:         self.assertEqual(self.stats_analyzer.total_responses, 6)
33:         self.assertEqual(len(self.stats_analyzer.columns), 4)
34:         self.assertIn('age', self.stats_analyzer.columns)
35:         self.assertIn('gender', self.stats_analyzer.columns)
36:         self.assertIn('region', self.stats_analyzer.columns)
37:         self.assertIn('satisfaction', self.stats_analyzer.columns)
38:
39:     def test_cross_tabulate_valid_columns(self):
40:         """Test cross-tabulation with valid columns."""
41:         crosstab = self.stats_analyzer.cross_tabulate('gender', 'satisfaction')
42:
43:         self.assertIsInstance(crosstab, list)
44:         self.assertGreater(len(crosstab), 1)  # Should have header + data rows
45:
46:         # Check header row
47:         header = crosstab[0]
48:         self.assertIn('', header)  # First column should be empty (row labels)
49:
50:         # Check data rows
```

```
51: for row in crosstab[1:]: 52: self.assertIsInstance(row, list) 53:
self.assertEqual(len(row), len(header)) 54: 55: def
test_cross_tabulate_invalid_columns(self): 56: """Test cross-tabulation with
invalid columns.""" 57: with self.assertRaises(ValueError): 58:
self.stats_analyzer.cross_tabulate('invalid_column', 'gender') 59: 60: def
test_chi_square_test_valid_data(self): 61: """Test chi-square test with valid
data.""" 62: result = self.stats_analyzer.chi_square_test('gender',
'satisfaction') 63: 64: self.assertIsInstance(result, dict) 65:
self.assertIn('chi_square', result) 66: self.assertIn('p_value', result) 67:
self.assertIn('df', result) 68: self.assertIn('significant', result) 69: 70: #
Check data types 71: self.assertIsInstance(result['chi_square'], float) 72:
self.assertIsInstance(result['p_value'], float) 73:
self.assertIsInstance(result['df'], int) 74:
self.assertIsInstance(result['significant'], bool) 75: 76: def
test_chi_square_test_insufficient_data(self): 77: """Test chi-square test with
insufficient data.""" 78: # Create data with only one value per category 79:
limited_data = [ 80: {'gender': 'Male', 'satisfaction': 'High'}, 81: {'gender':
'Female', 'satisfaction': 'Low'} 82: ] 83: limited_analyzer =
StatsAnalyzer(limited_data) 84: 85: result =
limited_analyzer.chi_square_test('gender', 'satisfaction') 86: 87:
self.assertIn('error', result) 88: self.assertEqual(result['chi_square'], 0.0)
89: self.assertEqual(result['p_value'], 1.0) 90: 91: def
test_correlation_analysis_numeric_data(self): 92: """Test correlation analysis
with numeric data.""" 93: numeric_data = [ 94: {'age': '25',
'satisfaction_score': '8'}, 95: {'age': '30', 'satisfaction_score': '7'}, 96:
{'age': '35', 'satisfaction_score': '6'}, 97: {'age': '28',
'satisfaction_score': '9'}, 98: {'age': '42', 'satisfaction_score': '5'} 99: ]
100: numeric_analyzer = StatsAnalyzer(numeric_data)
```

```
101: 102: result = numeric_analyzer.correlation_analysis('age',
'satisfaction_score') 103: 104: self.assertIsInstance(result, dict) 105:
self.assertIn('correlation', result) 106: self.assertIn('sample_size', result)
107: self.assertIn('strength', result) 108: 109:
self.assertIsInstance(result['correlation'], float) 110:
self.assertIsInstance(result['sample_size'], int) 111:
self.assertIsInstance(result['strength'], str) 112: 113: def
test_correlation_analysis_insufficient_data(self): 114: """Test correlation
analysis with insufficient data.""" 115: insufficient_data = [ 116: {'age': '25',
'score': '8'} 117: ] 118: insufficient_analyzer =
StatsAnalyzer(insufficient_data) 119: 120: result =
insufficient_analyzer.correlation_analysis('age', 'score') 121: 122:
self.assertIn('error', result) 123: self.assertEqual(result['correlation'], 0.0)
124: self.assertEqual(result['sample_size'], 0) 125: 126: def
test_calculate_correlation(self): 127: """Test correlation calculation.""" 128:
x = [1, 2, 3, 4, 5] 129: y = [2, 4, 6, 8, 10] # Perfect positive correlation 130:
131: correlation = self.stats_analyzer._calculate_correlation(x, y) 132:
self.assertAlmostEqual(correlation, 1.0, places=5) 133: 134: def
test_calculate_correlation_negative(self): 135: """Test negative correlation
calculation.""" 136: x = [1, 2, 3, 4, 5] 137: y = [10, 8, 6, 4, 2] # Perfect
negative correlation 138: 139: correlation =
self.stats_analyzer._calculate_correlation(x, y) 140:
self.assertAlmostEqual(correlation, -1.0, places=5) 141: 142: def
test_calculate_correlation_no_correlation(self): 143: """Test correlation
calculation with no correlation.""" 144: x = [1, 2, 3, 4, 5] 145: y = [1, 1, 1,
1, 1] # No variation in y 146: 147: correlation =
self.stats_analyzer._calculate_correlation(x, y) 148:
self.assertEqual(correlation, 0.0) 149: 150: def
test_interpret_correlation(self):
```

```
151: """Test correlation interpretation.""" 152: # Test very strong positive
correlation 153: strength = self.stats_analyzer._interpret_correlation(0.9) 154:
self.assertEqual(strength, "Very Strong") 155: 156: # Test strong positive
correlation 157: strength = self.stats_analyzer._interpret_correlation(0.7) 158:
self.assertEqual(strength, "Strong") 159: 160: # Test moderate correlation 161:
strength = self.stats_analyzer._interpret_correlation(0.5) 162:
self.assertEqual(strength, "Moderate") 163: 164: # Test weak correlation 165:
strength = self.stats_analyzer._interpret_correlation(0.3) 166:
self.assertEqual(strength, "Weak") 167: 168: # Test very weak correlation 169:
strength = self.stats_analyzer._interpret_correlation(0.1) 170:
self.assertEqual(strength, "Very Weak") 171: 172: def
test_analyze_response_patterns(self): 173: """Test response pattern analysis."""
174: patterns = self.stats_analyzer.analyze_response_patterns() 175: 176:
self.assertIsInstance(patterns, dict) 177: self.assertIn('common_combinations',
patterns) 178: self.assertIn('response_clusters', patterns) 179:
self.assertIn('outliers', patterns) 180: 181:
self.assertIsInstance(patterns['common_combinations'], list) 182:
self.assertIsInstance(patterns['response_clusters'], list) 183:
self.assertIsInstance(patterns['outliers'], list) 184: 185: def
test_get_statistical_summary(self): 186: """Test statistical summary
generation.""" 187: summary = self.stats_analyzer.get_statistical_summary() 188:
189: self.assertIsInstance(summary, dict) 190: self.assertIn('total_responses',
summary) 191: self.assertIn('total_columns', summary) 192:
self.assertIn('numeric_columns', summary) 193:
self.assertIn('categorical_columns', summary) 194:
self.assertIn('statistical_tests', summary) 195: 196:
self.assertEqual(summary['total_responses'], 6) 197:
self.assertEqual(summary['total_columns'], 4) 198:
self.assertIsInstance(summary['numeric_columns'], list) 199:
self.assertIsInstance(summary['categorical_columns'], list) 200:
```

```python
201:    def test_perform_multiple_chi_square_tests(self):
202:        """Test multiple chi-square tests."""
203:        results = self.stats_analyzer.perform_multiple_chi_square_tests('gender')
204:
205:        self.assertIsInstance(results, list)
206:
207:        for result in results:
208:            self.assertIn('column1', result)
209:            self.assertIn('column2', result)
210:            self.assertEqual(result['column1'], 'gender')
211:            self.assertNotEqual(result['column2'], 'gender')
212:
213:    def test_calculate_expected_frequencies(self):
214:        """Test expected frequency calculation."""
215:        observed = [[2, 1], [1, 2]]  # 2x2 contingency table
216:
217:        expected = self.stats_analyzer._calculate_expected_frequencies(observed)
218:
219:        self.assertIsInstance(expected, list)
220:        self.assertEqual(len(expected), 2)
221:        self.assertEqual(len(expected[0]), 2)
222:
223:        # Check that expected frequencies sum to total
224:        total_observed = sum(sum(row) for row in observed)
225:        total_expected = sum(sum(row) for row in expected)
226:        self.assertAlmostEqual(total_expected, total_observed, places=5)
227:
228:    def test_chi_square_p_value(self):
229:        """Test chi-square p-value calculation."""
230:        # Test with small chi-square value
231:        p_value = self.stats_analyzer._chi_square_p_value(1.0, 2)
232:        self.assertGreater(p_value, 0.0)
233:        self.assertLessEqual(p_value, 1.0)
234:
235:        # Test with large chi-square value
236:        p_value = self.stats_analyzer._chi_square_p_value(10.0, 2)
237:        self.assertGreaterEqual(p_value, 0.0)
238:        self.assertLess(p_value, 1.0)
239:
240:        # Test with zero degrees of freedom
241:        p_value = self.stats_analyzer._chi_square_p_value(5.0, 0)
242:        self.assertEqual(p_value, 1.0)
243:
244:
245: if __name__ == '__main__':
246:     unittest.main()
```

# 11. Sample Data

The project includes sample survey data for testing and demonstration purposes. The sample data contains realistic survey responses with various demographic and response patterns.

## Sample Survey Data - sample_survey.csv

```
1: age,gender,region,education,satisfaction,feedback,recommend 2:
25,Female,North,Bachelor,Very Satisfied,Great experience with the product!,Yes
3: 32,Male,South,Master,Satisfied,Good but could be better.,Yes 4:
45,Female,East,High School,Neutral,It was okay, nothing special.,Maybe 5:
28,Male,West,Bachelor,Dissatisfied,Poor quality and bad service.,No 6:
35,Female,North,PhD,Very Satisfied,Excellent product and amazing support!,Yes 7:
42,Male,South,Bachelor,Satisfied,Overall good experience.,Yes 8:
29,Female,East,Master,Neutral,Average service, nothing exceptional.,Maybe 9:
38,Male,West,High School,Dissatisfied,Very disappointed with the service.,No 10:
31,Female,North,Bachelor,Very Satisfied,Love the product! Highly recommend.,Yes
11: 27,Male,South,Master,Satisfied,Good value for money.,Yes 12:
50,Female,East,PhD,Neutral,It meets basic expectations.,Maybe 13:
33,Male,West,Bachelor,Dissatisfied,Not worth the money at all.,No 14:
26,Female,North,High School,Very Satisfied,Amazing experience!,Yes 15:
39,Male,South,Bachelor,Satisfied,Decent product, would buy again.,Yes 16:
44,Female,East,Master,Neutral,It's okay, nothing special.,Maybe 17:
36,Male,West,PhD,Dissatisfied,Terrible customer service.,No 18:
23,Female,North,Bachelor,Very Satisfied,Best purchase ever!,Yes 19:
41,Male,South,High School,Satisfied,Good quality product.,Yes 20:
30,Female,East,Bachelor,Neutral,Average experience.,Maybe 21:
47,Male,West,Master,Dissatisfied,Complete waste of money.,No
```

## 12. Requirements and Dependencies

The project uses only Python standard library modules, making it easy to run without external dependencies. This ensures maximum compatibility and ease of deployment.

# 13. Project Documentation

The README.md file provides comprehensive documentation including installation instructions, usage examples, technical specifications, and development guidelines.

## Project Documentation - README.md

```
1: # Survey Data Analyzer and Response Processor 2: 3: **A comprehensive Command
Line Interface (CLI) tool for analyzing survey responses from CSV files** 4: 5:
## Project Overview 6: 7: The Survey Data Analyzer is a modular Python CLI
application designed to perform structured statistical analysis, basic text
sentiment analysis, identify response patterns, and generate readable reports
from survey data stored in CSV format. 8: 9: ### Key Features 10: 11: - **CSV
Input Handling**: Load and validate structured survey data from CSV files 12: -
**Survey Summary**: Calculate response counts, percentages, and demographic
breakdowns 13: - **Statistical Analysis**: Perform cross-tabulation and
Chi-Square tests 14: - **Sentiment Analysis**: Analyze text responses using
keyword-based sentiment scoring 15: - **Pattern Recognition**: Identify
correlations and trends between survey responses 16: - **Report Generation**:
Create comprehensive analysis reports in text format 17: 18: ## Project Structure
19: 20: ``` 21: SurveyAnalyzer/ 22: ■■■ main.py # CLI menu and application entry
point 23: ■■■ data_loader.py # CSV loading and validation 24: ■■■
survey_summary.py # Response counts and demographics 25: ■■■ stats_analyzer.py #
Cross-tabulation and chi-square tests 26: ■■■ sentiment_analyzer.py # Text
sentiment analysis 27: ■■■ pattern_detector.py # Pattern correlation detection
28: ■■■ report_generator.py # Report generation and file output 29: ■■■
utils.py # Utility functions and helpers 30: ■■■ sample_survey.csv # Example
input file 31: ■■■ requirements.txt # Project dependencies (stdlib only) 32:
■■■ README.md # This documentation 33: ■■■ test/ # Unit tests 34: ■■■
test_data_loader.py 35: ■■■ test_stats_analyzer.py 36: ■■■
test_sentiment_analyzer.py 37: ``` 38: 39: ## Quick Start 40: 41: ###
Prerequisites 42: 43: - Python 3.8 or higher 44: - No external dependencies
required (uses only standard library) 45: 46: ### Installation 47: 48: 1. Clone
or download the project files 49: 2. Navigate to the SurveyAnalyzer directory 50:
3. Run the application:
```

```
51: 52: ```bash 53: python main.py 54: ``` 55: 56: ### Usage Example 57: 58: 1.
**Load Survey Data**: Choose option 1 and provide the path to your CSV file 59:
2. **View Summary Statistics**: Choose option 2 to see demographic breakdowns 60:
3. **Analyze Sentiment**: Choose option 3 for text sentiment analysis 61: 4.
**Cross-tabulate Results**: Choose option 4 for statistical analysis 62: 5.
**Detect Patterns**: Choose option 5 to find correlations 63: 6. **Generate
Report**: Choose option 6 to create a comprehensive report 64: 65: ## Sample Data
66: 67: The project includes `sample_survey.csv` with example data containing:
68: - Age, gender, region, education demographics 69: - Satisfaction ratings 70:
- Text feedback responses 71: - Recommendation preferences 72: 73: ## Technical
Requirements 74: 75: ### Python Version 76: - **Minimum**: Python 3.8 77: -
**Recommended**: Python 3.9+ 78: 79: ### Standard Library Modules Used 80: -
`csv`: CSV file handling 81: - `re`: Regular expressions for text processing 82:
- `math`: Mathematical calculations 83: - `statistics`: Statistical functions 84:
- `os`: Operating system interface 85: - `sys`: System-specific parameters 86: -
`datetime`: Date and time handling 87: - `itertools`: Iterator building blocks
88: - `collections`: Specialized container datatypes 89: - `typing`: Type hints
support 90: 91: ### Object-Oriented Design 92: - **Class Hierarchies**: Modular
class structure for each component 93: - **Inheritance**: Base classes for common
functionality 94: - **Polymorphism**: Flexible method implementations 95: 96: ##
User Guide 97: 98: ### 1. Loading Survey Data 99: 100: The application accepts
CSV files with the following requirements:
```

101: - **File Format**: Comma-separated values (CSV) 102: - **Encoding**: UTF-8, Latin-1, or CP1252 103: - **Size Limit**: Maximum 50MB 104: - **Required Columns**: At least 2 columns with valid data 105: 106: ### 2. Survey Summary Analysis 107: 108: View comprehensive statistics including: 109: - Total response count and response rate 110: - Demographic breakdowns (age, gender, region, education) 111: - Response distributions for each question 112: - Data quality metrics 113: 114: ### 3. Sentiment Analysis 115: 116: Analyze text responses using: 117: - **Positive Keywords**: excellent, amazing, great, good, etc. 118: - **Negative Keywords**: terrible, awful, bad, poor, etc. 119: - **Negation Handling**: "not good" → negative sentiment 120: - **Intensifier Recognition**: "very good" → amplified positive 121: 122: ### 4. Statistical Analysis 123: 124: Perform advanced statistical tests: 125: - **Cross-tabulation**: Compare responses across different questions 126: - **Chi-Square Tests**: Determine statistical significance 127: - **Correlation Analysis**: Find relationships between variables 128: 129: ### 5. Pattern Detection 130: 131: Identify meaningful patterns: 132: - **Demographic Patterns**: Age, gender, regional differences 133: - **Response Correlations**: Related answers across questions 134: - **Combination Patterns**: Common response combinations 135: - **Outlier Detection**: Unusual or rare responses 136: 137: ### 6. Report Generation 138: 139: Generate comprehensive reports including: 140: - Executive summary 141: - Demographic analysis 142: - Sentiment analysis results 143: - Pattern detection findings 144: - Statistical insights 145: - Key findings and recommendations 146: 147: ## Testing 148: 149: ### Running Tests 150:

```bash
# Run all tests
python -m unittest discover test/

# Run specific test file
python -m unittest test.test_data_loader

# Run with verbose output
python -m unittest discover test/ -v
```

### Test Coverage

- **Data Loading**: File validation, encoding detection, data cleaning
- **Statistical Analysis**: Cross-tabulation, chi-square calculations
- **Sentiment Analysis**: Keyword matching, negation handling
- **Pattern Detection**: Correlation analysis, demographic patterns
- **Report Generation**: File output, formatting

## Technical Documentation

### Core Classes

#### DataLoader
- **Purpose**: Handle CSV file loading and validation
- **Key Methods**: `load_csv()`, `_validate_file()`, `_clean_and_validate_data()`
- **Features**: Multi-encoding support, delimiter detection, data cleaning

#### SurveySummary
- **Purpose**: Generate summary statistics and demographic breakdowns
- **Key Methods**: `generate_summary()`, `_analyze_demographics()`, `get_age_distribution()`
- **Features**: Response counting, percentage calculations, trend analysis

#### StatsAnalyzer
- **Purpose**: Perform statistical analysis and hypothesis testing
- **Key Methods**: `cross_tabulate()`, `chi_square_test()`, `correlation_analysis()`
- **Features**: Cross-tabulation matrices, chi-square tests, correlation coefficients

#### SentimentAnalyzer
- **Purpose**: Analyze text sentiment using keyword-based approach
- **Key Methods**: `analyze_text()`, `analyze_column()`, `_clean_text()`
- **Features**: Keyword dictionaries, negation handling, intensifier recognition

#### PatternDetector
- **Purpose**: Identify patterns and correlations in survey responses
- **Key Methods**: `find_patterns()`, `_analyze_age_patterns()`, `_find_correlation_patterns()`
- **Features**: Demographic analysis, correlation detection, outlier identification

#### ReportGenerator
- **Purpose**: Generate comprehensive analysis reports

201: - **Key Methods**: `generate_report()`, `_build_report_content()`,
`_write_report_to_file()` 202: - **Features**: Structured reports, multiple
sections, file output 203: 204: ### Error Handling 205: 206: The application
implements robust error handling: 207: - **File Validation**: Check file
existence, size, and format 208: - **Data Validation**: Validate age ranges,
gender values, required fields 209: - **Graceful Degradation**: Continue
processing even with partial data 210: - **User-Friendly Messages**: Clear error
messages with suggestions 211: 212: ### Performance Considerations 213: 214: -
**Memory Efficient**: Process data in chunks for large files 215: - **Fast
Processing**: Optimized algorithms for statistical calculations 216: -
**Responsive UI**: Non-blocking operations with progress indicators 217: 218: ##
Educational Features 219: 220: ### Learning Objectives 221: 222: This project
demonstrates: 223: - **Modular Design**: Clean separation of concerns 224: -
**Object-Oriented Programming**: Classes, inheritance, polymorphism 225: -
**Error Handling**: Robust exception management 226: - **Documentation**:
Comprehensive docstrings and comments 227: - **Testing**: Unit tests for all
major components 228: - **CLI Design**: User-friendly command-line interface 229:
230: ### Code Quality 231: 232: - **Readability**: Clear variable names and
function purposes 233: - **Maintainability**: Modular structure for easy updates
234: - **Extensibility**: Easy to add new analysis features 235: -
**Documentation**: Inline comments and comprehensive docstrings 236: 237: ##
Contributing 238: 239: ### Development Guidelines 240: 241: 1. **Code Style**:
Follow PEP 8 conventions 242: 2. **Documentation**: Add docstrings for all
functions and classes 243: 3. **Testing**: Write unit tests for new features 244:
4. **Error Handling**: Implement proper exception handling 245: 5. **User
Experience**: Ensure intuitive CLI interface 246: 247: ### Adding New Features
248: 249: 1. Create new module in appropriate directory 250: 2. Implement class
with clear interface

251: 3. Add comprehensive unit tests 252: 4. Update main.py to integrate new feature 253: 5. Update documentation 254: 255: ## License 256: 257: This project is designed for educational purposes. Feel free to use, modify, and distribute for learning and teaching. 258: 259: ## Author 260: 261: **Student Developer** - Survey Data Analyzer Project 262: 263: --- 264: 265: **Note**: This application uses only Python standard library modules, making it easy to run without external dependencies. Perfect for educational environments and learning Python programming concepts.