

AegisLife Insurance Analytics Project

Title: Data Quality Assessment Report

Presented By: Ayush Kumar Siddharth

Data Quality Challenge

- ▶ Project Scope
 - ▶ 6,082 insurance records across 5 datasets
 - ▶ Data period: 2020-2025
 - ▶ 50 data quality issues identified
 - ▶ 538 records corrected
- ▶ Tools Used
 - ▶ Microsoft Excel - Data validation
 - ▶ Python (Pandas) - Automated cleaning
 - ▶ SQL - Database validation
 - ▶ Power BI - Quality monitoring
- ▶ Goal: Achieve 99%+ data quality for analytics

Executive Summary

KEY METRICS

50	538	89.7% → 99.8%
Total Issues Identified	Records Fixed	Quality Score
(8.85%)		Improvement

ISSUE SEVERITY

- Critical: 5 issues (10%)
- Moderate: 7 issues (14%)
- Minor: 38 issues (76%)

STATUS: Data ready for analytics

Most Common Problems Identified

1. FORMAT INCONSISTENCY (30%)

- Mixed date formats (MM/DD/YYYY, DD-MM-YYYY, YYYY-MM-DD)
- Boolean values (Yes/No, Y/N, True/False, 1/0)
- Text capitalization inconsistent

2. INVALID VALUES (20%)

- Typos: "Mal" instead of "Male"
- Wrong categories: "Activ" vs "Active"
- Misspelled status values

3. MISSING DATA (18%)

- Blank age fields (10 records)
- Null risk scores (5 records)
- Missing premium amounts (6 records)

4. OUTLIERS (10%)

- Extreme claim amounts flagged for review
- Processing times >180 days

5. LOGICAL ERRORS (6%)

- Claim dates before policy start
- Policy end dates before start dates

Data Quality Transformation

METRIC	BEFORE	→	AFTER	CHANGE
Total Records	6,082		6,067	-15 removed
Completeness	92.5%		99.8%	+7.3% ↑
Accuracy	91.2%		99.5%	+8.3% ↑
Format Consistency	85.0%		100.0%	+15.0% ↑
Referential Integrity	100%		100.0%	Maintained
Duplicate Records	0		0	Clean

OVERALL QUALITY SCORE

89.7% → 99.8% (+10.1% improvement)

Only 0.25% of records removed during cleaning

How We Fixed the Data

PHASE 1: EXCEL-BASED CLEANING

- ✓ Standardized date formats 113 records
 - ✓ Fixed gender/category typos 95 records
 - ✓ Imputed missing values 16 records
 - ✓ Validated foreign keys All tables
 - ✓ Removed invalid records 15 records
-

Total Phase 1: 409 records

PHASE 2: PYTHON AUTOMATION

- ✓ Encoding fixes (Latin-1) All files
 - ✓ Automated data type conversion All columns
 - ✓ Merge validation All tables
 - ✓ Calculated derived fields New columns
-

Total Phase 2: 6,082 records

Deep Dive: Standardizing Date Formats

THE PROBLEM

3 different date formats across datasets:

BEFORE:

- 12/31/2023 (MM/DD/YYYY) - 40% of records
- 31-12-2023 (DD-MM-YYYY) - 35% of records
- 2023-12-31 (YYYY-MM-DD) - 25% of records

- ✗ Prevented chronological sorting
- ✗ Caused analysis errors
- ✗ Time-series charts failed

THE SOLUTION

Python Code:

```
df['date'] = pd.to_datetime(df['date'], errors='coerce')  
df['date'] = df['date'].dt.strftime('%Y-%m-%d')
```

AFTER:

- 2023-12-31 (YYYY-MM-DD) - 100%

RESULT: 113 dates standardized to ISO 8601 format

Automated Quality Controls

DATA QUALITY VALIDATION RULES

RULE 1: Age Range

Must be between 18-100 years

Formula: =IF(AND(age>=18, age<=100), "Valid", "Invalid")

RULE 2: Risk Score Range

Must be between 0.00-1.00

RULE 3: Date Format

Must be YYYY-MM-DD

RULE 4: Categorical Values

Must be from predefined list only

RULE 5: Foreign Key Integrity

All foreign keys must exist in parent table

RULE 6: Logical Consistency

End date must be after start date

IMPACT: Prevents 80% of future data quality issues

Action Plan for Future Data Quality

1. IMPLEMENT DROP-DOWN LISTS

- Categorical fields use predefined options
- Reduces typos by 90%

2. USE DATE PICKERS

- Enforces YYYY-MM-DD format automatically
- Eliminates format inconsistencies

3. REAL-TIME VALIDATION

- Validate data at point of entry
- Immediate error feedback to users

4. STANDARDIZE BOOLEAN FIELDS

- Use only "Yes" or "No"
- Remove variations (Y/N, True/False, 1/0)

5. DATABASE CONSTRAINTS

- Enable foreign key constraints
- Prevent invalid relationships

6. MONTHLY DATA AUDITS

- Regular quality checks
- Catch issues early

EXPECTED IMPACT: 80% reduction in data quality issues