

Functional Practice Statements - Data Profiling
<div><div>Level 1: Initial</div><div><div>1.1 Basic profiling is performed for a data store(s).</div><div>Basic profiling includes such things as analyzing the types or number of distinct values in a column, number or percent of zero, blank or null values, string length, date ranges, patterns, as opposed to the more advanced analysis such as cardinality, frequency distributions, or key integrity</div><div>Example Work Products<ul style="list-style-type: none">Data profiling reportsList of data profiling checks</div></div></div>
<div><div>Level 2: Managed</div><div><div>2.1A data profiling methodology is established and followed.</div><div>The methodology adopted or created by the organization describes the approach to data profiling. The methodology will typically address planning and scoping the effort, profiling techniques, report templates, and presentation formats for summary results. In addition, profiling processes should be reusable and leveraged across multiple data stores and shared data repositories.</div><div>2.2 Data profiling plans are established for projects</div><div>Components typically included in the plan:<ul style="list-style-type: none">Selection of the data store(s) to examineIdentification of the data set(s) to profileList of stakeholders and definition of their involvementObjectives of the profiling activityData quality criteria based on the objectives, which includes referential integrity (parent and child) of the data, consistency of the data with respect to its documented metadata, consistency with established rules and patterns, and standard data quality dimensionsRules to be applied during the profiling activityMethod(s) and tool(s) for data profilingTemplate(s) for documenting resultsSchedule of activities, including resources</div><div>2.3 Plans for profiling a data store are shared with relevant stakeholders and data governance.</div><div>Data profiling activities should not be planned or executed in a vacuum.</div><div>Stakeholders and data governance authorities may have specific needs that should be taken into consideration. In addition, recognizing that the expense and effort to conduct profiling activities is not trivial, it is important for the profiling activities to be aligned with business needs. Sharing plans for profiling can help to ensure agreements and continued alignment.</div><div>2.4 Data profiling activities are conducted according to the plan, and efforts are adjusted when significant deviations from plan are detected.</div><div>Because a profiling effort often produces some unexpected results, the organization needs to be flexible enough to determine, during the profiling initiative, if initial results justify additional time and effort to expand the scope.</div><div>2.5 Data profiling results and recommendations are reported to the stakeholders.</div><div>Results should be used as input to data quality assessment and data cleansing efforts, as well as to inform the data quality strategy. For example, an organization may determine that its product data has an unacceptable percentage of errors. The stakeholders need to weigh in on the impact of these errors and have a role in determining the remediation alternatives and approach.</div><div>Example Work Products<ul style="list-style-type: none">Data profiling methodology documentationApproved data profiling plan and scheduleData profiling findings reports and metricsDefined skill set and training plan for staff with data quality responsibilities</div></div></div>
<div><div>Level 3: Defined</div><div><div>3.1 Data profiling methodologies, processes, practices, tools, and result templates have been defined and standardized.</div><div>Standard profiling tools should be identified and consistently used across the organization to gain efficiencies. An organizational standard method for analyzing and presenting business and technical impacts of data profiling on remediation activities should also be defined and followed. Report templates and metrics are standardized, centrally stored, and published to ensure consistent application across the organization.</div><div>3.2 All techniques identified to meet the profiling objectives are performed.</div><div>While data profiling can be considered primarily a discovery activity, it is typically initiated to meet specified objectives. The profiling techniques and tools employed must support achieving those objectives. Tailoring of techniques and templates may be required. Often a profiling effort has several phases; for example, out of the box profiling checks (e.g., value ranges, ID uniqueness, etc.); standardization analysis (e.g., addresses, syntax); tests for selected business rules; and known issues (which may require complex queries).</div><div>The data profiling team should be fully conversant in all techniques and corresponding tool capabilities selected for the profiling activities.</div><div>3.3 Traceability between data requirements, documented metadata, the physical data, and data quality rules is captured and maintained.</div><div>Data profiling activities should be executed by data profiling experts aware of data requirements, data quality rules, data content, and data structures.</div><div>This is achieved through establishing traceability between data requirements, physical data, and metadata.</div><div>3.4 Data governance is engaged to identify core shared data sets and the corresponding data stores that should be regularly profiled and monitored.</div><div>The organization should have defined rules for when various data sets are profiled (e.g., when data is acquired; or prior to being consolidated, migrated, exported, analyzed, reported for compliance purposes, or structurally transformed).</div><div>3.5 Profiling processes are reusable and deployed across multiple data stores and shared data repositories.</div><div>Sharing and mentoring among peers to build data profiling best practices should occur across the organization, because data quality activities are often highly valuable but resource-intensive. Therefore, it is desirable to leverage efficiencies and avoid rework. Some organizations find that the most effective approach is to operate with a single data profiling team with skilled</div><div>3.6 The SDLC includes data profiling tasks with tailoring criteria, guidance, and governance.</div><div>Most data store development efforts (for example, creation of a new data warehouse) should include data profiling activities as a planned part of the project. Institutionalization of data profiling practices requires that the SDLC include reference and guidelines for these activities, and that tailoring criteria are defined and followed.</div><div>Example Work Products<ul style="list-style-type: none">Data profiling standards, including criteria for processes, standards, best practice criteria, tailoring, and reporting formatsData profiling methodologies tailored from the organizational standardReport showing traceability of data requirements with the data content and characteristics revealed through profiling resultsDocumented tailoring of data-related decisions and rationaleDocumentation that practitioners have required profiling skillsData profiling metricsRecommendations reports from data profiling effortsBusiness and technical impact analysis results templateStandard data profiling report requirementsApproved standard data profiling tool(s)Data profile baselines</div></div></div>
<div><div>Level 4: Measured</div><div><div>4.1 Performance of data profiling processes is measured and used to manage activities across the organization.</div><div><ul style="list-style-type: none">Plans and schedules for data profiling should be managed according to the feedback provided by data quality measurements. Measurement should indicate how well the output of this activity addresses and aligns with the business need and priorities. Decisions on what, when, and how to profile data should be driven by indications of quality and criticality, which may vary by business application. Highly shared data and data sets deemed vital to key business processes should be regularly profiled, as data quality is critical and needs to be frequently monitored.Data quality measures should also indicate how well the staff performed the data profiling activities. The evaluation of plans and execution (actual vs. estimates) should consider such things as use of techniques, impact of results and decisions, compliance with methods and standards, quality of output, and level of effort.Data profiling process performance baselines can be created and used to inform the planning and execution of data profiling activities and results.</div><div>4.2 Data profiling efforts include evaluation of the conformity of data content with its approved metadata and standards.</div><div>Approved standard business terms, meanings, values, and ranges are used as a benchmark for profiling the data content in a data store. Additional documentation is typically found in corporate data dictionaries, data models, system requirements documentation, etc. It is a best practice to update this documentation as needed.</div><div>4.3 During a data profiling activity, actual issues are compared to the statistically predicted issues based on historical profiling results.</div><div>Results should be systematically compared to corresponding historical profiles to evaluate impact of profiling activities on corrective actions and quality improvements.</div><div>4.4 Results are centrally stored, systematically monitored, and analyzed with respect to statistics and metrics to provide insight to data quality improvements over time.</div><div>A consistent impact analysis method can be applied to evaluate business, technical, and cost impacts of remediation. Summary profiling results are provided to data governance bodies and senior management. Results are used to inform data governance and data architecture decisions, especially for highly shared data.</div><div>Example Work Products<ul style="list-style-type: none">Documented profiling methodology, best practices, and standardsProject reports showing application of profiling results to data quality governanceDashboards, scorecards, or other decision support tools for data quality, showing the results of data profiling effortsData quality portal displaying data quality models and results to be used for performance baselines</div></div></div>
<div><div>Level 5: Optimized</div><div><div>5.1 The organization addresses root causes of defects and other issues based on an understanding of the meaning, technical characteristics, and behavior of the data over time.</div><div>5.2 Data profiling processes and other activities are analyzed to identify defects and make improvements based on the quantified expected benefits, estimated costs, and business objectives.</div><div>Data profiling results performed on the same data over time can be statistically analyzed periodically to measure the performance of profiling activities.</div><div>5.3 Real-time or near-real-time automated profiling reports are created for all critical data feeds and repositories.</div><div>Automating the performance and scheduling of profiling improves the data quality program's efficiency and responsiveness to planned and unplanned events.</div><div>Example Work Products<ul style="list-style-type: none">Log of stakeholders' usage of profiling resultsControl charts demonstrating that the processes used across data stores have stabilized (data stores have been sufficiently profiled)Data profiling process objectives for improvement included in standard data management strategies, programs, and reportsReal-time data profiling reports generated on scheduleConclusions drawn from data profiling process analyses and recommendations for improvement</div></div></div>