

本書は英語原文の理解を補助するための日本語参考訳です。

データ品質マネジメントガイドブック

データと AI の価値を最大化する

(第 1.02 版)

2026-05-14

AISI Japan
AI Safety
Institute

IPA Information-technology
Promotion
Agency, Japan
Digital Infrastructure Center

目次

| | |
|------------------------------|----|
| 1 本ガイドブックについて..... | 6 |
| 1.1 はじめに | 6 |
| 1.2 本ガイドブックの使い方 | 6 |
| 1.3 対象読者 | 7 |
| 2 背景と概要 | 8 |
| 2.1 背景 | 8 |
| 2.2 データ品質と社会的信頼 | 9 |
| 2.3 AI システムとデータ品質 | 10 |
| 2.4 AI セーフティとデータ品質 | 10 |
| 2.5 高品質データの利点 | 11 |
| 2.6 低品質データのリスク | 13 |
| 2.7 低品質データの原因 | 14 |
| 2.8 目的と目標 | 15 |
| 2.9 原則 | 16 |
| 2.10 スコープ | 17 |
| 2.11 データと AI のステークホルダー | 18 |
| 2.12 品質確保の方法 | 19 |
| 2.13 標準とフレームワークの関係 | 20 |
| 2.14 日本における品質関連ガイド | 20 |

| | |
|-------------------------------|----|
| 3 データ品質マネジメントフレームワークの考え方..... | 22 |
| 3.1 マネジメント、ガバナンス、マチュリティ..... | 22 |
| 3.2 フレームワーク概要..... | 23 |
| 3.3 フレームワークの3つのビュー..... | 24 |
| 3.4 データライフサイクル..... | 25 |
| 3.5 品質特性..... | 26 |
| 3.6 管理対象のコンテンツとデータ管理..... | 28 |
| 3.7 AIのためのデータ品質マネジメント..... | 30 |
| 4 評価観点..... | 31 |
| 4.1 プロセスビュー（運用）..... | 31 |
| 4.1.1 概要..... | 31 |
| 4.1.2 データ計画..... | 33 |
| 4.1.3 データ取得..... | 40 |
| 4.1.4 データ準備..... | 44 |
| 4.1.5 データ処理..... | 54 |
| 4.1.6 AI システム..... | 55 |
| 4.1.7 出力評価..... | 60 |
| 4.1.8 結果提供..... | 62 |
| 4.1.9 データ廃棄..... | 65 |
| 4.1.10 ライフサイクル全体にわたるプロセス..... | 67 |

| | |
|--------------------------------|-----|
| 4.2 ガバナンスサイクルビュー | 70 |
| 4.2.1 概要 | 70 |
| 4.2.2 データ品質計画 | 73 |
| 4.2.3 データ品質管理 | 74 |
| 4.2.4 データ品質保証 | 76 |
| 4.2.5 データ品質改善 | 77 |
| 4.2.6 データ関連サポート | 78 |
| 4.2.7 リソース規定 | 80 |
| 4.3 ゲートウェイビュー（品質特性） | 83 |
| 4.3.1 概要 | 83 |
| 4.3.2 データ固有の品質特性 | 84 |
| 4.3.3 固有かつシステム依存のデータ品質特性 | 86 |
| 4.3.4 システム依存のデータ品質 | 90 |
| 4.3.5 AI/ML 向け追加データ品質特性 | 91 |
| 4.3.6 センサーデータ品質特性 | 95 |
| 5 結び | 101 |
| 5.1 メッセージ | 101 |
| 6 文書情報 | 101 |
| 6.1 発行主体・作成体制 | 101 |
| 6.2 参考文献 | 102 |

| | |
|------------------------------------|-----|
| 6.3 改版履歴 | 103 |
| 6.4 日本語参考訳について | 103 |
| 7 付録 1 | 104 |
| 7.1 機械学習品質マネジメントガイドライン 第 4 版 | 104 |
| 8 免責事項 | 107 |

1 本ガイドブックについて

1.1 はじめに

Garbage in, Garbage out.

このガイドブックは、データ品質が低いと AI の性能と信頼が直接損なわれる、というシンプルな前提に基づいている。AI は既存のデータ上の問題をしばしば増幅する。データの品質は、AI の出力に直接影響する。

AI システムは、幅広い社会・産業分野に深く組み込まれつつあり、私たちの日常生活に影響を与え、重要な意思決定を形づくっている。その信頼性と安全性の基盤にあるのは、何よりもまず、それらが依拠するデータの品質である。どれほど高度な AI システムであっても、不正確または信頼できないデータは、誤った結果や有害な結果につながり得る。したがって、**信頼できる AI を確保するための最初かつ最も重要な一歩は、データ品質を確保することである。**

しかし、データ品質を評価し管理する方法は、組織や分野によって異なる。加えて、AI システムで用いられるデータの性質はライフサイクル全体を通じて変化するため、単一で画一的な基準で十分であることはまれである。

本ガイドブックは、**データ品質マネジメントを実務で実装するための共通のフレームワークと考え方を提供することを目的とする。**国際標準を踏まえつつ、実務的なコラムや AI 特有の観点も取り入れることで、現実の場面で適用しやすく、有用なものとすることを目指している。

1.2 本ガイドブックの使い方

- 方針策定の指針として
 - 組織内でデータ品質を確保するための方針を策定する際に、本ガイドブックを基本的な参照資料またはチェックリストとして利用する。
 - 例として、データガバナンス方針、AI 倫理原則、品質マネジメント計画などが挙げられる。

- プロジェクト運営の実務ツールとして
 - データ収集、アノテーション、モデル学習など、AI 開発プロセスの各段階で、何を確認すべきか、誰が責任を持つかを明確にするために利用する。
- 教育・共通理解のための資料として
 - データの取扱いに関わるすべてのステークホルダー、すなわち提供者、保有者、開発者、利用者が共通理解を築くための研修コンテンツとして利用する。

本書は法的拘束力を持つ標準ではなく、データ品質に関する実践知を共有するためのガイドブックである。読者には、それぞれの役割、ニーズ、目的に応じて柔軟に参照し、自らの組織やプロジェクトに最も適した形に調整して活用することを期待する。

また本書は、ISO/IEC 等の関連標準や国内外のガイドを参照し、AI システムにおけるデータ品質マネジメントを実務上検討するための考え方および確認観点を整理したものである。本ガイドブックの利用は、個別の規格適合性、法令適合性、または認証取得を直接保証するものではない。

1.3 対象読者

データ品質マネジメントは、単一の役割や部門だけで完結するものではなく、組織の内外にまたがる多様な機能の協働を必要とする。

| 役割 | 想定される利用方法（例） |
|----------------------|---|
| 経営層 / 方針決定者 | 組織全体および社会全体でデータ品質を確保するための方針や組織体制を整備する際の指針として利用する。 |
| データ管理者 / データガバナンス担当者 | データ品質を維持・評価するための基準やプロセス設計の参照資料として利用する。 |

| 役割 | 想定される利用方法（例） |
|---------------------|---|
| データ・AI エンジニア | AI 開発ライフサイクル全体を通じたデータ品質マネジメントおよび検証の実務指針として利用する。 |
| 事業企画担当者 / サービス設計者 | サービス設計やリスク管理において必要なデータ品質要件を定義し、整合させる際の参考とする。 |
| 監査担当者 / 研究者 / 教育関係者 | データ品質マネジメントの評価、教育、研究のための基礎資料として利用する。 |

2 背景と概要

2.1 背景

AI は、より有用で効率的なサービスを可能にすることで、急速に社会を変革している。しかし、多くの人々はなお AI の利用に懸念を抱いており、安全に導入するためには AI 出力の正確性を確保することが不可欠である。そのためには、AI の構築や運用に用いるデータの品質を高める必要がある。

- AI 社会
 - AI システムは、学習や、予測・推奨・説明などの出力を生成するために大量のデータに依存している。データの品質は、AI システムの性能と有効性に直接影響する。
- データ駆動型社会
 - データ駆動型社会では、正確で信頼できる情報を確保するために高品質なデータが不可欠である。これにより、適切な情報に基づく意思決定が支えられ、誤りが減り、リスクが最小化される。また、高品質なデータは顧客満足度を高め、組織が規制に適合することを助け、法的・財務的問題から組織を守る。

データ品質を確保することは、AI サービスの信頼性を確保するうえで不可欠である。

2.2 データ品質と社会的信頼

AI およびデータの利用には、信頼の確保が極めて重要である。データは AI の基盤である。データが信頼できなければ、全体のプロセスの信頼性も損なわれる。

以下の 2 つのリストは、社会における信頼や信頼できる AI をめぐる幅広い議論の中で、データ品質がどのように位置づけられるかを示したものである。

- データ駆動型社会における信頼の要素
 - サービス品質
 - 倫理
 - 透明性
 - アカウンタビリティ
 - プライバシー保護
 - セキュリティ
 - **データ品質**
 - パートナーシップと連携
- 信頼できる AI の要素
 - **正確性と信頼性**
 - 倫理
 - 透明性
 - アカウンタビリティ
 - ユーザー中心設計
 - 安全性
 - 継続的改善

データ品質は AI の卓越性の基盤であり、信頼できる AI を実現し、それが利用者の受容と関与を促進する。

2.3 AI システムとデータ品質

AI システムにおいて、データは利用時だけでなく、準備段階においても重要な役割を果たす。

下図は、OECD による AI システムの定義に基づき、データ品質に関する領域を示したものである。データは、外部環境と AI モデルをつなぐ AI システムの不可欠な要素である。これらのデータは、モデル構築時に用いる入出力データと、利活用時に用いる入出力データに大別できる。

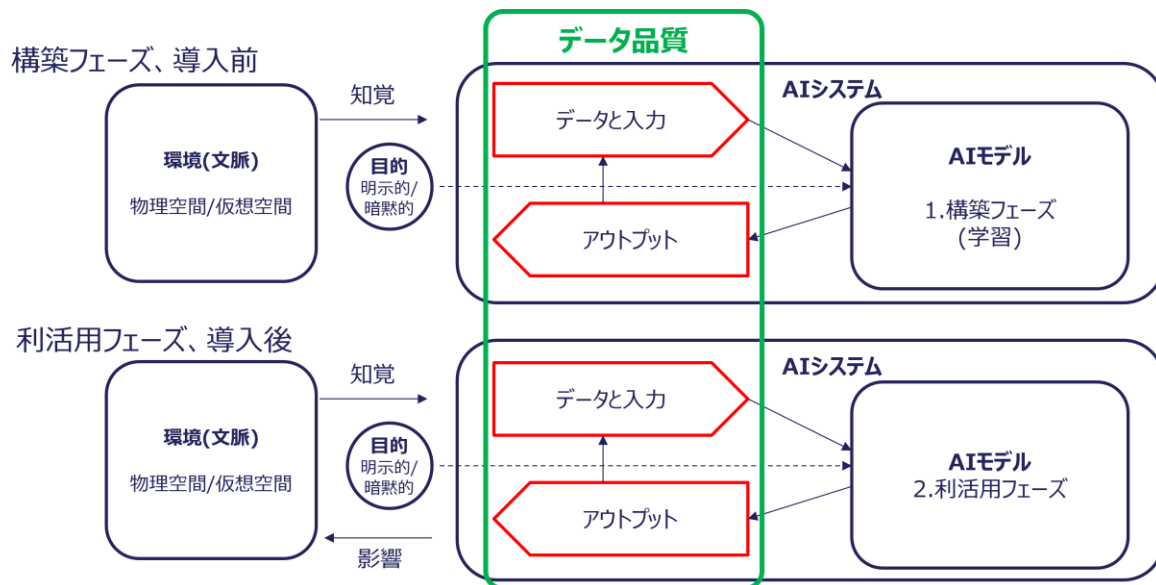


図 1. AI システムとデータ品質の関係

出典：OECD.AI, 2023, “Updates to the OECD’s definition of an AI system explained”

2.4 AI セーフティとデータ品質

低品質な入力データは、AI モデルの性能、出力の信憑性・一貫性・正確性など、多くの重要な側面に影響する。その結果、公衆の信頼が損なわれる可能性がある。データ品質マネジメントは、AI の適切かつ安全な利用を確保するために不可欠である。次のリストは、AI セーフティに関する評価観点を整理したものであり、その一つとしてデータ品質が位置づけられている。また、データ品質は、公平性、

ロバスト性、検証可能性など、AI セーフティに関わる複数の特性とも密接に関係している。

- AI セーフティに関する評価観点
 - 有害情報の出力制御
 - 偽誤情報の出力・誘導の防止
 - 公平性と包摂性
 - ハイリスク利用・目的外利用への対処
 - プライバシー保護
 - セキュリティ確保
 - 説明可能性
 - ロバスト性
 - **データ品質**
 - 検証可能性

出典：AI Safety Institute(AISI), 2025, “[Guide to Evaluation Perspectives on AI Safety](#)”

2.5 高品質データの利点

高品質なデータを利用することには、次のような幅広い利点がある。

1. 正確性の向上
 - データ品質が高いほど、AI の予測や意思決定はより正確かつ一貫したものとなり、不適切な結果のリスクが低減する。
2. インサイトの向上
 - 高品質なデータにより、実行可能なインサイトや信頼できる傾向を特定でき、より良い意思決定が可能になる。
3. 効率の向上
 - クリーンで高品質なデータは、クレンジング、修正、正規化といった前処理作業に費やす時間を減らし、チームがより高い価値のある活動に集中できるようにする。

4. ユーザー体験の改善
 - 正確で関連性の高いデータにより、精緻で個別化され、適時な結果を提供できるため、エンドユーザー全体の体験が向上する。
5. エラーの削減
 - データ内の不整合や誤りを最小化することで、AI 出力の不正確さが減り、信頼性と性能が向上する。
6. コスト削減
 - 高品質データへの投資は、誤りの修正、再処理、低品質データに起因する下流工程の問題への対応にかかるコストを低減する。
7. コンプライアンスとセキュリティ
 - 高品質なデータを維持することで、データガバナンス、プライバシー規制、セキュリティ標準への適合が確保され、組織の健全性が守られる。
8. 信頼の強化
 - 一貫して高品質なデータは、ステークホルダーの AI システムに対する信頼を高め、自動化されたプロセスや意思決定への安心感を醸成する。
9. 拡張性
 - 高品質なデータは、AI システムを拡張するための堅牢な基盤となり、システム拡大時にも一貫した性能を確保する。
10. モデル学習と更新の促進
 - クリーンで整備されたデータは、AI モデルの学習およびチューニングを容易にし、開発サイクルを短縮し、モデルの適応力を高める。
11. 競争優位性
 - 高品質なデータを活用する組織は、優れた AI 駆動型製品やサービスを提供することで、大きな競争優位性を得る。
12. エコシステム統合

- 高品質なデータは、他のシステムやプラットフォームとの円滑な統合を可能にし、相互運用性と効率的なワークフローを確保する。

2.6 低品質データのリスク

低品質なデータの利用は、次のような幅広いリスクを生み出し得る。

1. 意思決定の誤り
 - 不正確または不完全なデータは分析の誤りにつながり、誤った結論、戦略上の失敗、組織のあらゆるレベルでの不適切な意思決定を招く。
2. 非効率な運用
 - 低品質データのクレンジングや修正に時間と資源を割く必要があり、AI モデルのデプロイが遅れ、運用の有効性が低下する。
3. 顧客満足度の低下
 - 不正確または不完全な顧客データは、不十分なパーソナライズ、期待外れの体験、サービスや製品への信頼低下につながる。
4. コスト増加
 - 低品質データは、誤り修正、再処理、AI モデルの再学習、および誤った予測や推奨に起因する潜在的な財務損失を通じてコストを増大させる。
5. 法的・規制上のリスク
 - データ保護法および規制（例：GDPR、CCPA）への非適合は、多額の罰金や訴訟につながる可能性がある。
6. レピュテーション毀損
 - 低品質データに起因する誤りは、顧客、パートナー、ステークホルダーからの信頼や信用を損ない、長期的なブランド毀損につながり得る。
7. 競争上の不利

- より高品質なデータを有する競合は、顧客理解、業務効率、市場対応力といった重要領域で優位に立ち、組織が後れを取る可能性がある。

8. 機会損失

- 低品質データにより重要なインサイトや傾向を見逃すことで、新たな市場機会を捉えられなかったり、有効なイノベーションを起こせなかったりする。

9. モデル性能の低下

- 低品質データで学習した AI モデルは、バイアス、信頼性の低さ、有害なふるまいを示す可能性があり、倫理上の懸念や実世界での有効性低下につながる。

10. セキュリティリスク

- 低品質データには、悪意ある者の悪用可能な脆弱性や誤りが意図せず含まれることがあり、セキュリティ侵害や機微情報の不正利用につながる可能性がある。

11. ステークホルダーの不信

- データ品質が継続的に信頼性と結果を損なう場合、社内チームや社外ステークホルダーは AI システムへの信頼を失う可能性がある。

2.7 低品質データの原因

低品質データの一般的な原因は次のとおりである。

- **ヒューマンエラー**：データ入力やデータ処理時に生じる操作上の誤り。
- **標準化の不足**：データソース間における形式や標準の不一致。
- **不十分なデータガバナンス**：データ管理に関する方針や手順が不足。
- **旧式システム**：現代のデータ要件に対応できないレガシーシステムの利用。

- **不完全なデータ収集**：不十分なデータ収集プロセスに起因する欠損や不完全なデータ。
- **不十分な教育訓練**：データを扱う人員への適切な訓練不足。
- **不十分なデータ統合**：異なるソースからのデータを統合する際に生じる問題。
- **限定的な品質管理**：データ品質に対する確認や検証プロセスが不十分。
- **未検証のソース**：信頼できない、または未検証のソースから得たデータの利用。
- **データ収集におけるバイアス**：内在する偏りを反映したデータの収集。
- **技術的不具合**：ソフトウェアのバグやハードウェア障害に起因する誤り。
- **ドキュメント不足（例：データ辞書、データリネージ図）**：誤解や誤った解釈を招く不十分な文書化。
- **誤解**：要件や指示の誤解によるデータ入力ミス。
- **悪意ある行為**：有害な意図を持つ者による意図的な改ざんやデータポイズニング。
- **不十分なインセンティブ設計**：エンドユーザやデータ入力者にとって、高品質なデータを提供する動機づけが不足。

2.8 目的と目標

本書におけるデータ品質マネジメントの目的と目標は次のとおりである。

- **目的**
 - AI のためのデータ品質を向上させることにより、意思決定、予測、推薦が、正確で一貫性があり、信頼できる情報に基づいて行われるようになる。高品質なデータは、組織が AI を効果的に活用し、リス

クを低減し、その潜在力を最大限に引き出すことを可能にする。社会的な観点では、データ品質の向上は、AI システムへの信頼を醸成し、倫理的なデータ利用を促進し、公平性を支える。これにより、医療、教育、公共サービスにおける革新的な応用に加え、社会的・経済的便益に寄与する。

- 目標
 - データ品質マネジメントの目標は、データ収集、クレンジング、検証、管理に関する標準化された実践を確立することである。それにより、正確で安全かつ利用しやすいデータに基づき、AI システムが最大限の性能を発揮できる環境を整備する。社会的な観点では、データ駆動型イノベーションの便益を分野横断的に広げ、透明性を高め、意思決定プロセスを改善し、技術が人類に責任ある効果的な形で貢献する、持続可能で包摂的なデジタル社会への進展を加速することを目指す。

2.9 原則

本ガイドブックでは、次をデータ品質の原則とみなす。

- 利用者とそのニーズを理解する。
- 目的に適合していることを確保する。
- 品質とコストのトレードオフを考慮する。
- エラーは発生するものと受け止め、完全性を追い求めすぎない。
- 設計を見直し、ライフサイクルを考慮する。
- 適切な能力を備えた成熟したサービスが既に存在する場合は活用する。
- すべてのステークホルダーからフィードバックを得る。
- 確認しやすく、追跡しやすいよう可視化する。

求められるデータ品質の水準は、そのデータが目的に適合しているかどうかによって異なる。AI セーフティの文脈では、AI システムの用途、リスクおよび潜在的な影響に応じて要件を調整するこの考え方は、リスクベースアプローチと呼ばれるこ

とが多い。高リスクな用途では、来歴、代表性、バイアス、個人データ、出力評価、監査可能性などについて、特に注意して確認する必要がある。

品質とコストも重要である。品質を計画する際には、コストを考慮しなければならない。

- 高い品質には、通常、より高いコストがかかる。(a)
 - 目標を上回る品質まで引き上げれば、その分コストは増える。
- 設計のような初期工程で品質を組み込むことが重要である。(b)
 - 設計が適切であれば、運用コストや総コストが低下し、変化にも迅速に対応できる。
- 初期費用だけでなく運用コストも考慮する。(c)
 - システム実装時にデータを最適化すれば、運用コストを削減できる。

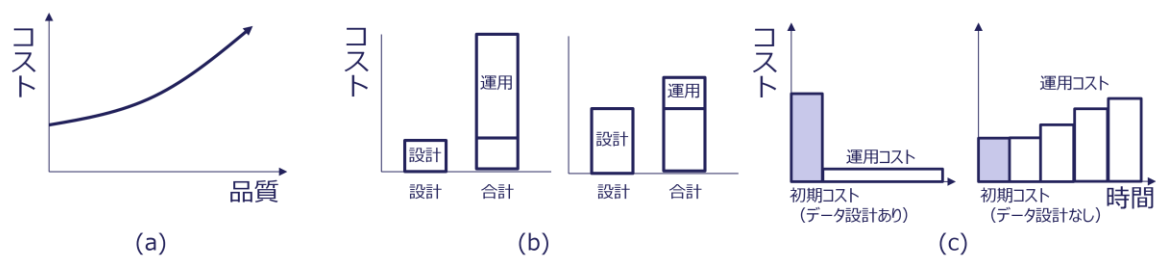


図 2. 品質とコスト

2.10 スコープ

AI とデータは社会において重要な役割を果たすが、そのスコープは広い。したがって、より広い社会的文脈の中で検討する必要がある。

下図は、システムから人とシステムとの関わり、社会的インパクトへと広がる、より大きな階層の中でデータと AI システムがどのように位置づけられるかを示している。本ガイドブックは、この広い文脈の中でも、主としてデータとデータシステムを含む AI システムの領域に焦点を当てる。

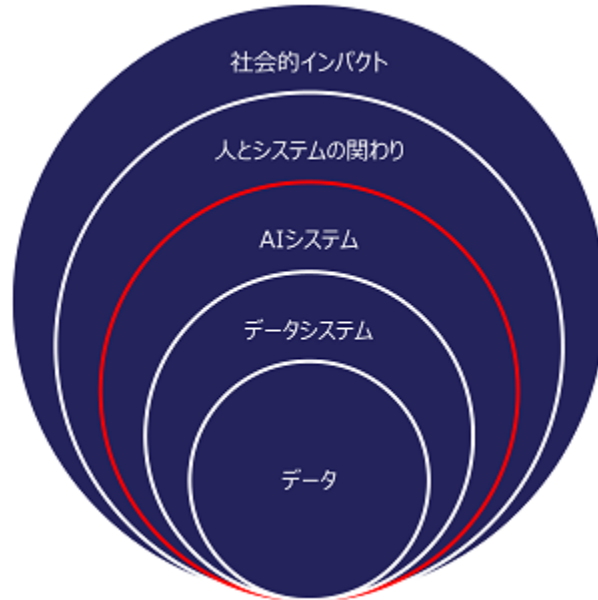


図3. AI とデータのスコープと位置づけ

AI システムは、数値データ、テキスト、画像、音声、動画などの非構造データを含む多種多様なデータを扱う。本ガイドブックは、特定のデータ種別に依存しない一般的なアプローチを示す。AI システムには、従来の特定タスク向け AI システムと生成 AI システムの両方が含まれるが、本ガイドブックでは両者を厳密に区別せずに共通して適用できる考え方に焦点を当てる。ただし、検索拡張生成（Retrieval-Augmented Generation : RAG）のように生成 AI システムに特有の議論も一部含む。

2.11 データと AI のステークホルダー

ステークホルダーは多様であり、その多くは AI の利用者であると同時に、データの提供者でもある。

下図は、データおよび AI の領域に関わる主体と、それらの関係を示している。AI とデータを中心に、広範なステークホルダーが相互につながっている。

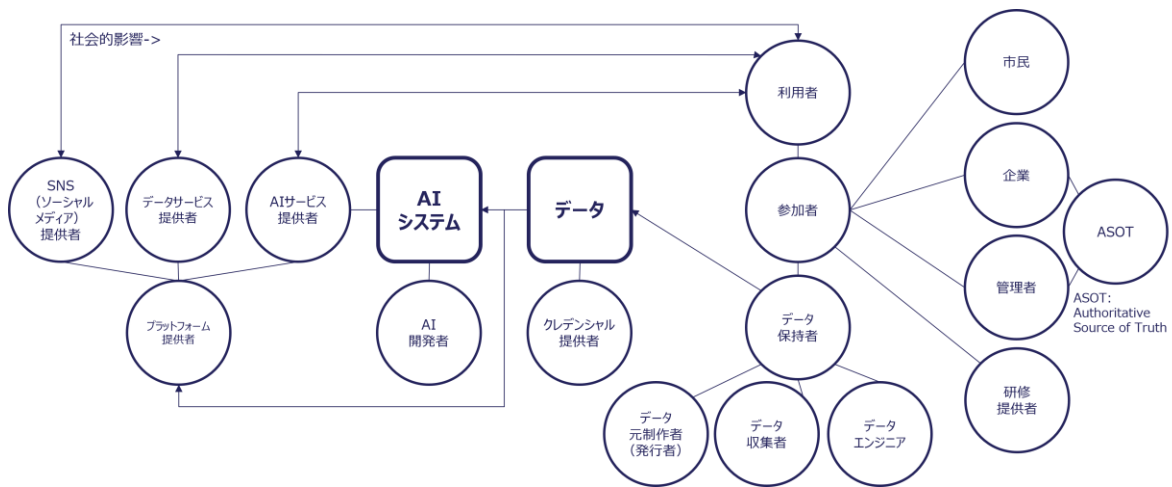


図 4. データと AI のステークホルダー関係図

2.12 品質確保の方法

次の方法が、データ品質を確保するための主なアプローチである。

- **高品質なデータをつくる**
 - 高度なデータ設計を確保するため、適切な参照モデルやデータモデルを活用する。データ作成時の手作業介入を避け、各段階で検証を実施する。
- **低品質なデータを防ぐ**
 - バリデータ、ディテクタ、コンバータ、クレンジングのツールを用いて入力・出力データを検証する。あわせて人による確認も実施する。
- **信頼を構築する**
 - データの来歴を確認する。デジタルコンテンツ透かし技術を活用する。ウォーターマークを埋め込む。
- **正しく使用する**
 - AI 学習におけるデータの適切な利用を確保する。UI/UX を考慮する。さらに、信頼できる唯一の情報源 (Single Source of Truth : SSOT) および権威ある情報源 (Authoritative Source of Truth : ASoT) の考え方

を用いてデータを管理する。あわせて、データに適切なラベルを付与する。

- **ガバナンスを確立する**
 - DataOps の考え方の活用などにより、適切なガバナンスを実装する。

2.13 標準とフレームワークの関係

本書で提供するフレームワークは、ISO、IEC といったデータ品質に関する多くの標準を、実装しやすい形で整理し、導入を導くものである。新たな標準をつくるものではなく、あくまで実務ガイドであり、実装結果を標準化プロセスへフィードバックすることで相乗効果が得られる。

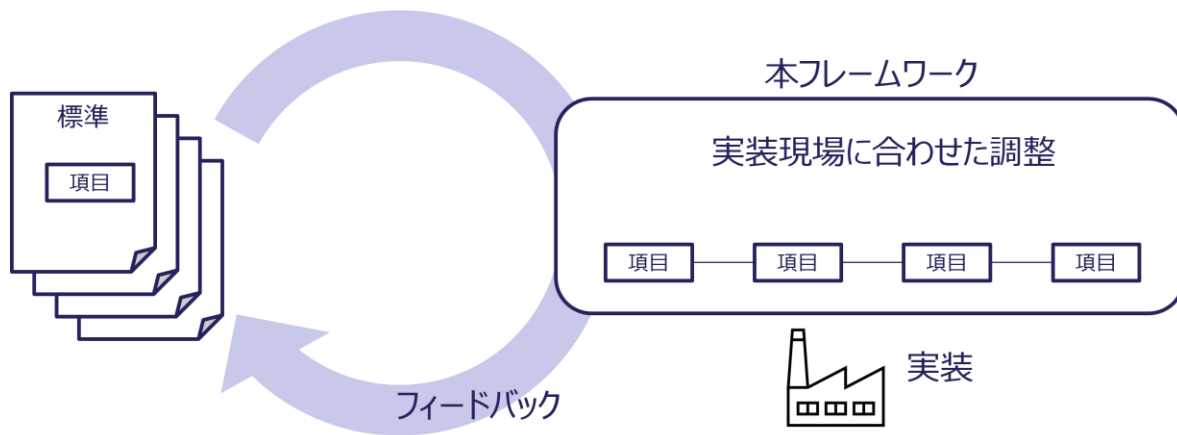


図5. 標準を実務へ適用するためのフレームワーク

2.14 日本における品質関連ガイド

日本では、AI 品質に関する複数のガイドが整備されている。

下図は、日本の主要な AI 品質関連ガイドの階層構造を示している。まず、AI の安全性確保を目的とする「AI 事業者ガイドライン」によって全体像を捉える。次に、AI の利用目的や AI モデルの特性に応じて適切なガイドを選択する。これらの AI システムは、本ガイドブックで示す考え方に基づいて適切に管理されたデータによって支えられる。さらに、データ設計や相互運用性に関するフレームワークが、これらの取組を下支えしてデータ品質向上を促進する。

| AI事業者ガイドライン（総務省、経済産業省） | | | |
|------------------------------|-------|--|------------------------------|
| | | AIシステムの目的 | |
| | | 生成 | 認識・予測 |
| AIモデルのタイプ | 生成 | A. 生成AI品質マネジメントガイドライン (AIST) B. AIセーフティに関する評価観点ガイド (AISI) | 両ガイドライン (A, C) |
| | 認識・予測 | 両ガイドライン (A, C) | C. 機械学習品質マネジメントガイドライン (AIST) |
| データ品質マネジメントガイドブック(AISI, IPA) | | | |
| 政府相互運用性フレームワーク (GIF) (デジタル庁) | | | |

図 6. 日本における AI 品質ガイドラインの構造

下図は、前の図で示した関係をデータの観点から示したものである。



図 7. AI 品質ガイドライン間の関係とデータ品質マネジメントの役割

コラム：日本のサービス品質

日本の高いサービス品質は、製品、人、データの品質に支えられて実現されている。

下図は、製品、人、データがどのように相互作用して高品質なサービスを支えているかを示している。サービス提供者と利用者の双方が品質に対して強い意識を共有しており、その行動がデータや製品の品質に反映され、サービス品質を高める継続的な循環を形成している。

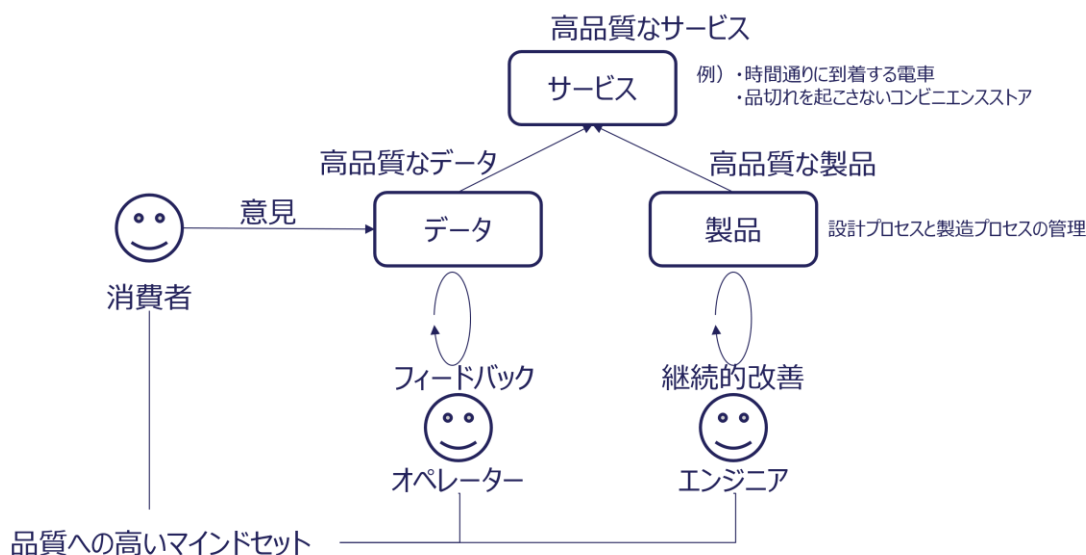


図 8. 日本の高いサービス品質の構造

3 データ品質マネジメントフレームワークの考え方

3.1 マネジメント、ガバナンス、マチュリティ

高品質なデータを継続的に提供することは、AI の普及拡大を促進するうえで不可欠である。効果的なマネジメントとガバナンスは、データ品質を維持し、データおよび AI 活用のマチュリティを高める鍵となる。

1. データ・AI マネジメント

- 適切なデータ・AI マネジメントは、データ・AI 品質を確保する。

*本ガイドブックは主にこの領域に焦点を当てる。

2. データ・AI ガバナンス

- 適切なデータ・AI ガバナンスは、持続可能なデータ・AI 品質を確保する。

3. データ・AI マチュリティ

- 品質が確保されたデータが供給されれば、データと AI の活用は拡大する。

下図は、マネジメント、ガバナンス、マチュリティの関係を示している。まず何よりも、データが適切にマネジメントされ、ガバナンスとマチュリティが維持されることで、AI によって価値が創出される。本ガイドブックでは、初期段階として、まずマネジメントに焦点を当てる。

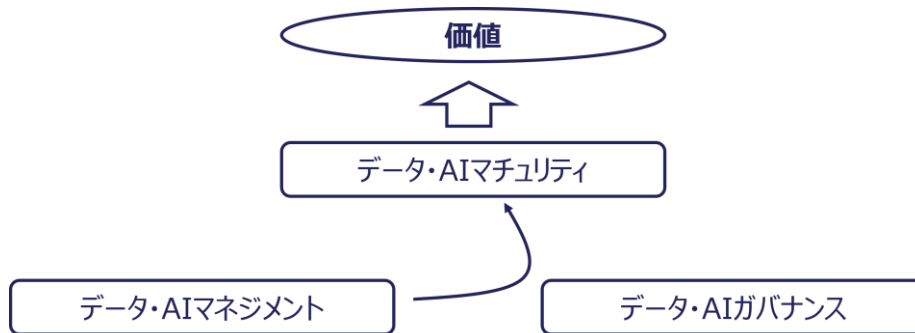


図9. マネジメント、ガバナンス、マチュリティの関係

3.2 フレームワーク概要

データはシステム、組織、国境を越えてやり取りされるため、相互運用性が不可欠である。

データ品質に関する国際標準や業界ガイドラインは数多く存在する。本ガイドブックでは、それらを特性、プロセス、ガバナンスの3つに整理する。

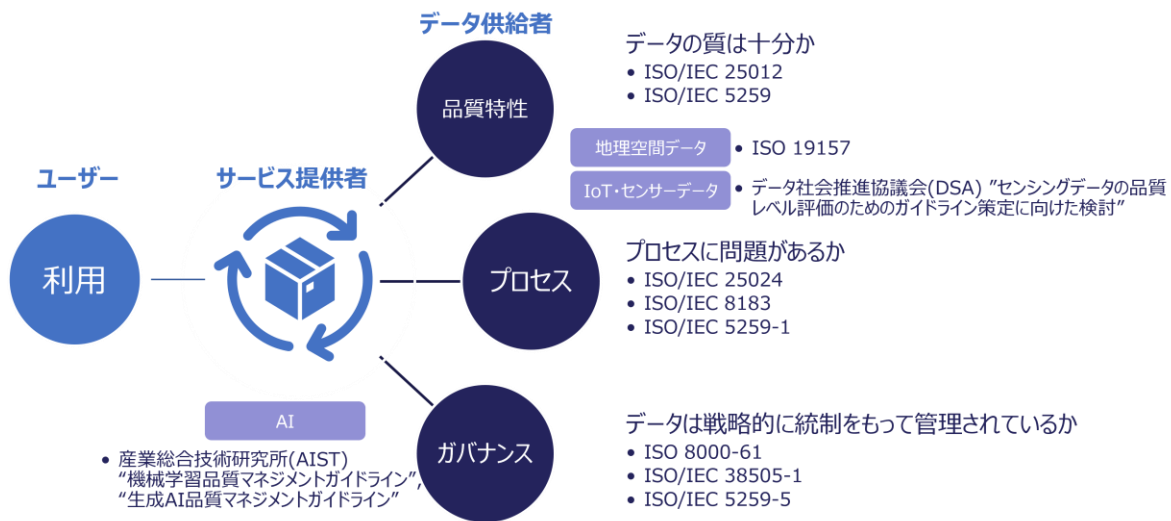


図 10. データ品質に関する国際標準の整理

3.3 フレームワークの 3 つのビュー

本ガイドブックでは、データの品質を多面的に捉えるために、「プロセスビュー」「ガバナンスサイクルビュー」「ゲートウェイ（データ品質特性）ビュー」という 3 つの視点を用いて整理する。

プロセスビューは、データに関わる活動に着目する視点である。データがどのように計画され、収集され、加工され、利用されるかといった一連の活動に加え、設計や要件定義といった上流の活動も含めて捉える。

ガバナンスサイクルビューは、これらの活動を継続的に実行・維持・改善するための体制や仕組みに着目する視点である。方針、役割、責任、評価および改善の仕組みといった観点から、組織としてのデータ品質管理を捉える。

ゲートウェイ（データ品質特性）ビューは、データそのものに着目する視点であり、正確性や一貫性、完全性といった、結果としてのデータ品質を評価する。これには、移植性や追跡可能性のようなデータの説明や形式に関する特性も含まれる。

これら 3 つのビューは、それぞれ以下の役割を持つ：

- プロセスビュー：活動（どのようにその品質が生み出されるか）
- ガバナンスサイクルビュー：仕組み（それらの活動をどのように維持・管理するか）
- ゲートウェイビュー：データの状態（どのような品質が実現されているか）

これら 3 つのビューを組み合わせることで、評価、実行、継続的改善を一体的に管理できるようになる。

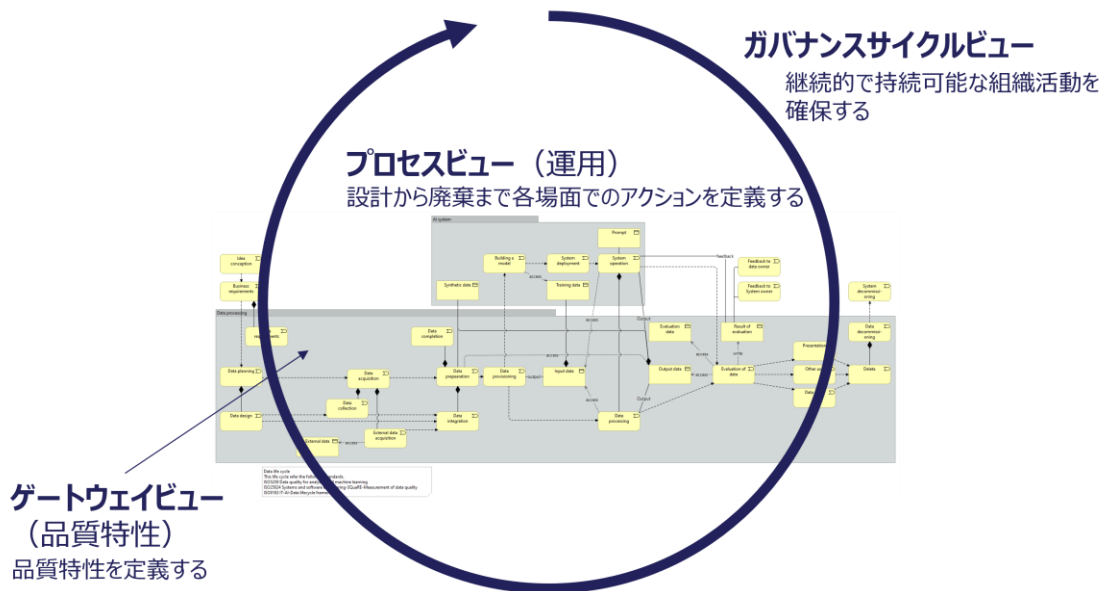


図 11. データ品質マネジメントの 3 ビュー・フレームワーク

3.4 データライフサイクル

データ品質マネジメントにおいては、データの加工や統合といった、データが存在することを前提とした後工程に注目が集まりやすい。しかし、実際のデータ品質は、設計や収集といった上流の段階から大きな影響を受ける。また、データの利用段階のみならず、最終的な廃棄までを視野に入れて管理することは、不要なデータの蓄積を抑制し、適切なデータ管理を実現する上で不可欠である。

したがって、本ガイドではデータを単発の処理対象として捉えるのではなく、設計から廃棄に至るまでの一連の流れとして捉える。この流れはデータライフサイクルと呼ばれる。

データライフサイクルを定義する標準にはさまざまなものがある。ここでは、複数の標準を統合・整理し、下図の通り詳細なライフサイクルを定義する。

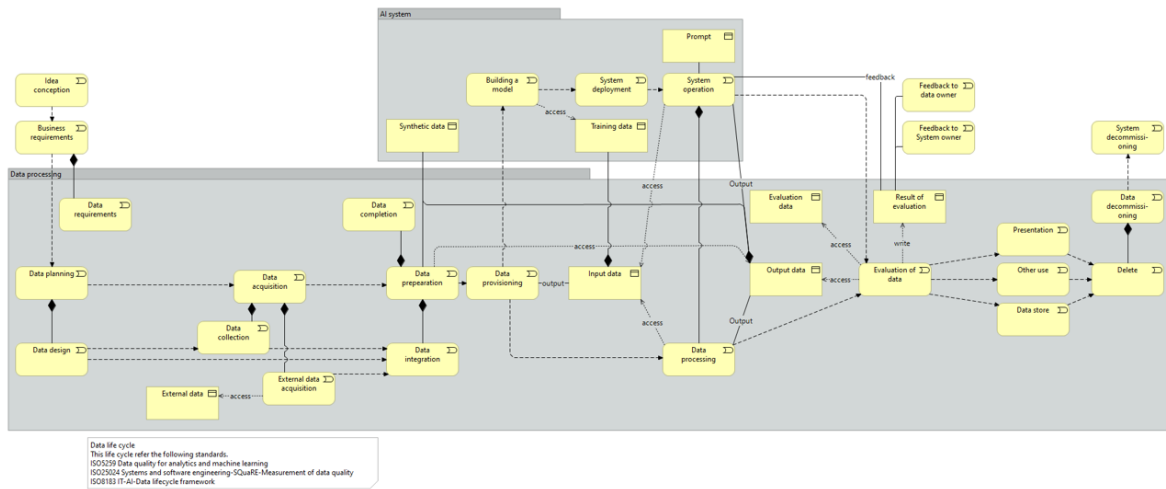


図 12. データライフサイクル

本ガイドでは、データライフサイクルにおける取組の妥当性をプロセスビューにより評価し、ライフサイクルの各段階におけるデータの状態をゲートウェイビューにより評価する管理の考え方を示す。

3.5 品質特性

最初のリストは典型的なデータ上の問題を示し、次のリストはそれら进行评估し対処するための品質特性を示している。

- 問題
 - 不正確なデータ
 - 整形されていないデータ
 - 古いデータ

- 出所不明のデータ
- ポイズニングされたデータ
- スパイクデータ
- 校正されていないデータ
- 偽誤情報・誤情報・悪意ある情報
- バイアス
- 不完全なデータ
- 重複データ
- 不整合なデータ
- 関連性のないデータ
- 破損したデータ
- 不明瞭なデータ
- **データ品質特性 (ISO/IEC 25012 および ISO/IEC 5259-2)**
 - 正確性
 - 完全性
 - 一貫性
 - 信憑性
 - 最新性
 - アクセシビリティ
 - 標準適合性
 - 機密性
 - 効率性
 - 精度
 - 追跡可能性
 - 理解性
 - 可用性
 - 移植性

- 回復性
- 監査可能性
- 均衡性
- 多様性
- 有効性
- 識別可能性
- 関連性
- 代表性
- 類似性
- 適時性

3.6 管理対象のコンテンツとデータ管理

AI は多様なデータを扱い、そのステークホルダーも多様である。下図は、データ品質、透明性、リスク管理との関係における、AI およびデータのエコシステムの多層構造を示している。これは全体理解のための概観図であり、データ品質、透明性、ガバナンスがエコシステム全体でどのように関係するかを示している。

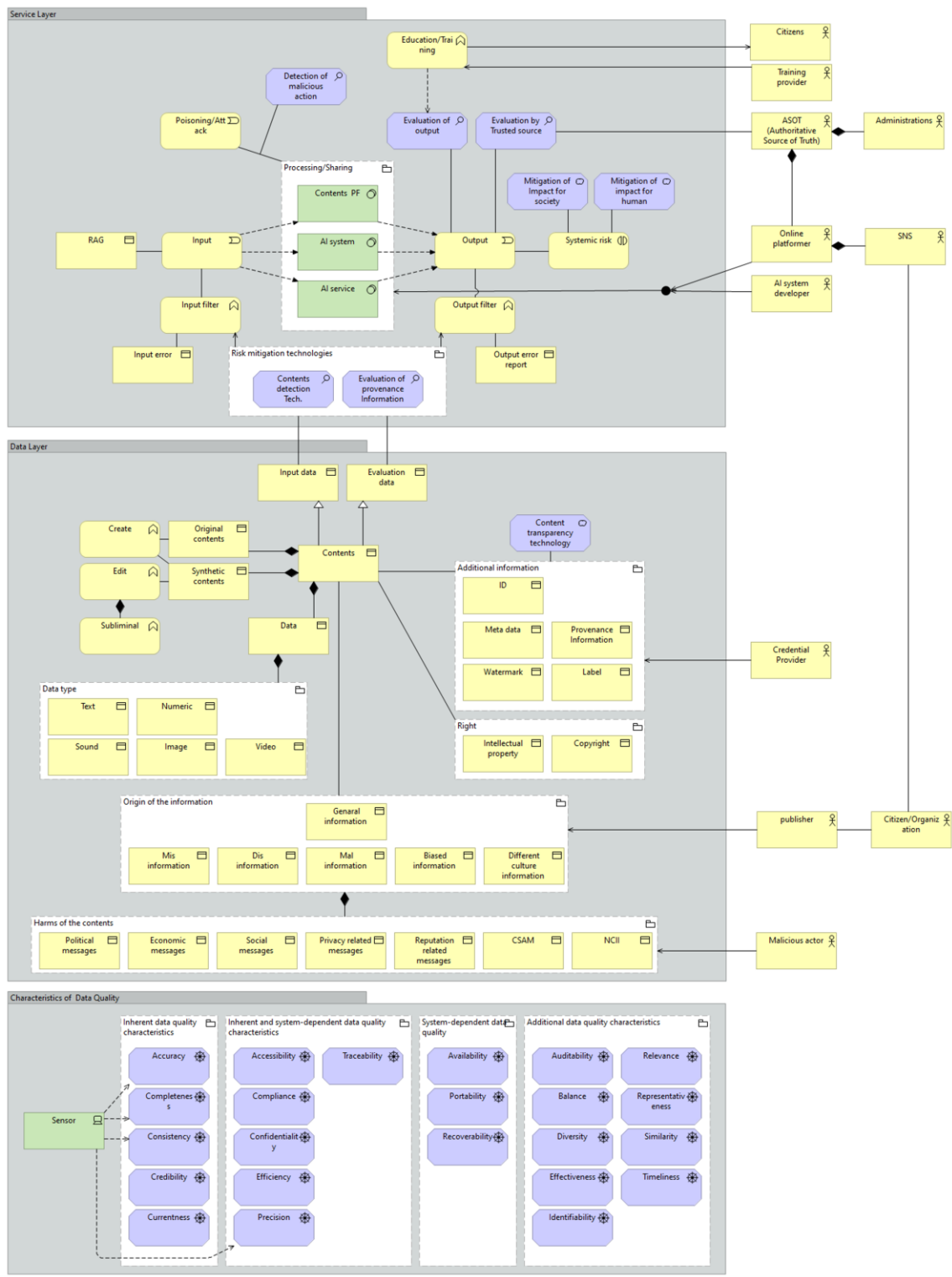


図 13. AI システムのために管理すべきコンテンツの統合構造

3.7 AI のためのデータ品質マネジメント

従来のデータ品質マネジメントでは、正確性、完全性、最新性などが重要である。しかし、データを AI に利用する場合には、次のような追加的な考慮が必要となる。

1. アノテーション（ラベリング）の品質管理
 - 機械学習では、手動または自動によるデータへのラベリングが不可欠である。不正確または一貫しないラベルが多いと、モデル精度の低下や学習結果の偏りにつながる。
2. バイアス確認と公平性の確保
 - データが特定の属性（例：性別、人種、地域）を過度に代表していないことを確認し、そのような偏りが学習過程で増幅されないよう防ぐことが重要である。
3. プライバシーおよびセキュリティ対策
 - 個人情報や機密情報を含むデータを扱う場合、適切な匿名化やマスキングを行い、安全な保管・処理を確保しなければならない。また、学習済みモデルを用いた個人情報の推定を防ぐため、差分プライバシーなどのプライバシー保護技術の検討も必要である。
4. バージョン管理とドリフト管理
 - AI モデルは、学習時点のデータ分布に基づいて訓練される。時間の経過とともに、データ分布や入力と出力の基礎的な関係が変化することがあり（データドリフト、コンセプトドリフト）、モデル性能が急激に低下する場合がある。データのバージョン管理、分布監視、必要に応じた再学習や改良を実施することが重要である。
5. 合成データおよび生成コンテンツへの対応
 - 現実世界のデータが不足している場合、合成データが用いられることがある。しかし、その生成方法や品質を明確にすることが重要である。不適切な合成データの利用は、誤った学習や不正確なモデル出力につながり得る。
6. 説明可能性と透明性の確保

- AI モデルは「ブラックボックス」として運用されるべきではない。説明可能な AI (XAI) の手法は、重要な結果の背後にある意思決定過程を明らかにし、信頼とアカウントビリティを高める。

7. 運用・監視体制の確立

- データやモデルの開発時だけでなく、デプロイ後も継続的な監視と保守が不可欠である。予期しないバイアスや性能低下が生じた場合に迅速に対応できる体制を整えることが重要である。

4 評価観点

4.1 プロセスビュー（運用）

4.1.1 概要

プロセスビューでは、データ品質を「データライフサイクルに沿って必要な活動が適切に実施されているか」という観点から捉える。したがって、本ビューではデータそのものでなく、それを生み出すプロセスを対象として評価する。これには、データライフサイクルにおけるステージに対応するプロセスと、複数のステージにまたがって実施される横断プロセスの両方が含まれる。

ここでは、ライフサイクルを次の 8 つの大きなステージに区分し、それぞれのステージに対応するプロセスとして整理する。

1. データ計画
2. データ取得
3. データ準備
4. データ処理
5. AI システム
6. 出力評価
7. 結果提供

8. データ廃棄

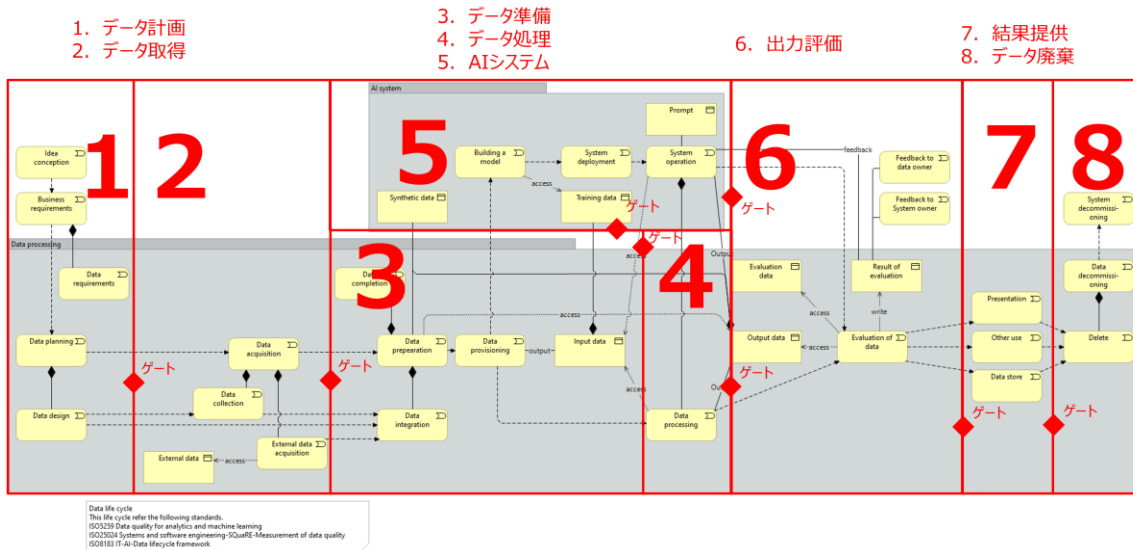


図 14. データライフサイクルプロセス

4.1.1.1 ステークホルダーの役割

下に示すステークホルダーマップは、プロセス全体にどの主体が関与するかを示している。組織横断の役割分担と協働の明確化に役立つ。

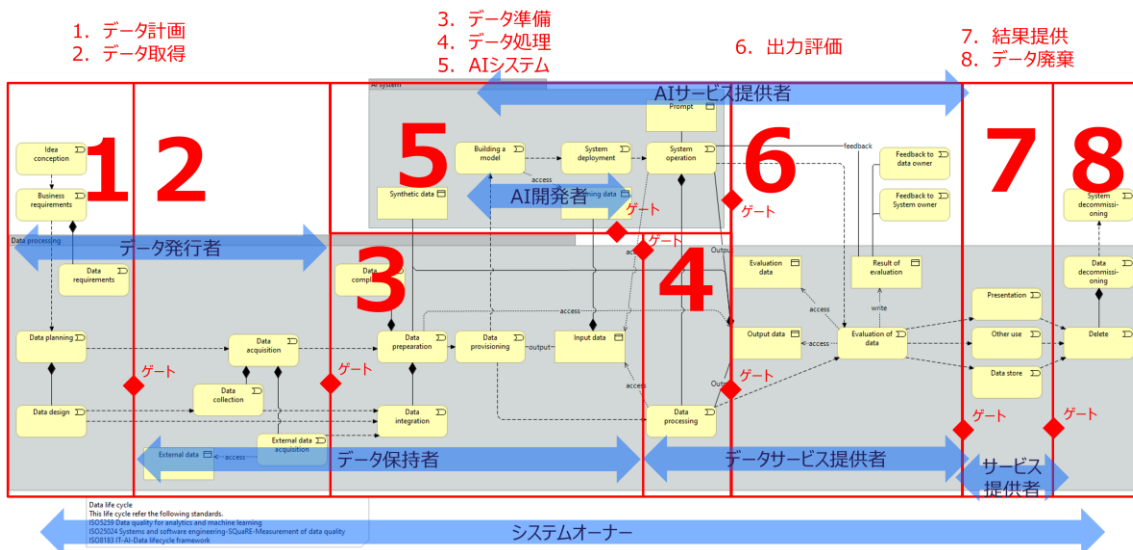


図 15. データライフサイクルにおける役割

4.1.1.2 本項の構成

以降の項は、ステージに対応する 8 つのプロセスとライフサイクル横断のプロセスについて、次の構成で整理する。

- 概要
 - プロセスの目的および役割を示すとともに、主な課題、要件、および得られる価値を示す。なお、プロセスの説明には、そのプロセスが対応するステージの概要が含まれる場合がある。
- プロセス
 - 当該プロセスを構成する個別のプロセスを示す。各プロセスについて、手順の例および実施状況を評価するためのチェックポイントを提示する。チェックポイントは簡潔な説明のためにデータの状態を問う場合があるが、プロセスビューではそれらが適切に確認・管理されているかを評価する。実際の適用にあたっては対象や環境に応じて適宜調整されることを前提とする。

4.1.2 データ計画

4.1.2.1 概要

- 説明
 - データ計画は、データライフサイクルにおける重要なプロセスである。
 - このプロセスでは、データを必要とする意図や動機の定義を扱う。
 - サービス全体での相互運用性と将来の成長に向けた拡張性に関する計画を含む。
 - 取得方法、評価方法、廃棄プロセスを含むデータライフサイクルを明示する。
- 課題
 - 既存データを効果的に管理・整理することの難しさ
 - 現在のニーズに合わせたルールやプロセスの改訂・更新の不足

- 新しい技術を活用してデータ管理実務を最適化することの難しさ
- 要件
 - 標準化され、よく構造化された形式を含む高品質なデータを定義・管理するプロセス
 - システムやサービス間の相互運用性を確保する仕組み
 - データの一貫性を保つための、信頼できる権威的な唯一の情報源を選定・維持するプロセス
- 価値
 - よく準備されたデータによる下流工程の効率向上
 - サービス全体アーキテクチャの持続可能性および拡張性の向上

品質は、コンセプト立案、要件定義、設計の3段階でつくり込まれる。

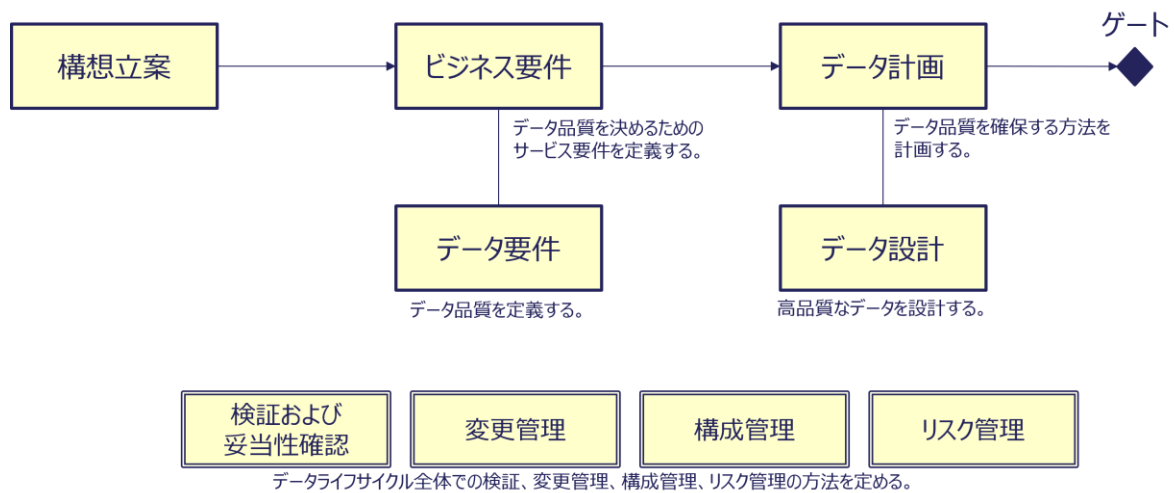


図 16. データ計画プロセス

4.1.2.2 構想立案

- 手順
 1. 利用者ニーズを収集する。
 2. コンセプトを作成する。
 3. 事業方針およびデータ方針を定義する。

- 組織方針と整合させる
- 将来動向を確認する
- 4. ステークホルダーを列挙する。
- 5. 実現可能性を確認する。
- チェックポイント
 - 意思決定者はシステム全体のコンセプトを理解しているか。
 - 意思決定者は高品質データの利点を理解しているか。
 - 意思決定者は低品質データがもたらすリスクを理解しているか。
 - 意思決定者は、必要以上に高いデータ品質を求めるコストと、有効性とコストの均衡の重要性を理解しているか。
 - 組織としての体制と十分なスキルを持つ人材を有しているか。

4.1.2.3 ビジネス要件

- 手順
 - 1. 事業目的と対象範囲を定義する。
 - 2. 価値、要件、制約、リスクを定義する。
- チェックポイント
 - サービス品質は明確に定義されているか。（処理時間、コストなど）

4.1.2.4 データ要件

- 手順
 - 1. 関連するデータを定義する。
 - 入力データ、出力データ、参照データ
 - 制約条件、インターフェース、統計データ
 - 2. データ仕様を明確にする。
 - 説明、目標、要件、考慮事項
 - データ仕様は ISO/IEC 5259-3 を参照する
 - 3. 必要な品質水準を定義する。

- チェックポイント
 - 事業上必要なデータが列挙されているか。
 - 各データについて品質管理項目が定義されているか。
 - 各データについてデータ品質要件レベルが定義されているか。

4.1.2.5 データ計画

- 手順
 1. 既存データを確認する。
 2. データニーズを収集する。
 3. マスターデータを定義する。
 4. データアーキテクチャを定義する。
 5. 設計方針と方法論を定義する。
 - 構造、場所、モデリング、文書化
 6. データソースを定義し、探索する。
 7. データ関連法令の適用有無を確認する。
- チェックポイント
 - 意思決定者は既存データを新モデルへ変換することに合意しているか。
 - 運用や意思決定に関するステークホルダーのデータニーズを理解しているか。
 - データアーキテクチャと設計方針は文書化されているか。
 - 必要なデータは定義され、利用可能か。
 - データ利用に関する法的制約を確認したか。
 - 個人情報、センシティブ情報、プライバシー関連データを含むか。
 - 知的財産に関するデータを含むか。

4.1.2.6 データ設計

- 手順

1. 参照モデルと分類体系を確認する。
 2. データモデルと分類体系を設計する。
 3. メタデータとラベルを設計する。
 4. ルールを設計する。
 5. 法令適合性および標準適合性を確認する。
 - 注：日本では GIF（政府相互運用性フレームワーク）が参照モデルを提供している。
- チェックポイント
 - データ参照モデルや標準化された分類体系を参照しているか。
 - モデリングツールを利用しているか。
 - メタデータはデータカタログ語彙（DCAT）をベースに設計されているか。
 - 利用およびアクセスに関する一般ルールを参照しているか。
 - 法令に適合するよう設計されているか。

コラム：参照モデル

参照モデルとは、特定の業界や分野向けに標準化された構造、プロセス、データフレームワークを提供する概念的枠組みである。これはシステムやプロセスを設計・実装する際の「ひな型」として機能し、用語、定義、方法を統一的に提供する。これにより、組織間・システム間の相互運用性が高まり、開発および運用の効率が向上する。高い抽象度を持つため、参照モデルは幅広いシナリオに適用でき、個別要件に応じてカスタマイズできる。

- 参照モデルがデータ品質を改善する仕組み
 - 参照モデルは、データの一貫性と標準化の基盤を確立する。
 - 共通のデータ定義や命名規則を採用することで、異なるシステム間の不整合を防ぐ。

- さらに、参照モデルを利用することでデータの重複や冗長性を排除し、効率的で簡潔なデータ構造を実現できる。
- 業務プロセスに整合したガバナンスルールや制約を組み込むことで、データの正確性と完全性を高める。

参照モデルは、設計段階でデータ完全性を確保し、運用段階でのエラーや修正コストを低減するうえで重要な役割を果たす。

コラム：データ辞書

データ辞書は、情報システムやデータベースで使用されるデータ要素の名称、定義、形式、関係を体系的に整理・管理するリポジトリである。各データ要素について、属性名、データ型、長さ、制約、業務ルールを定義し、組織全体で一貫した基準を提供する。

この共有されたデータの意味と構造に関する理解は、開発者や運用担当者などのステークホルダーの間で共通認識を形成し、システム開発・運用を効率化し、品質を高める。さらに、データマッピングやシステム統合の際には、異なるシステム間のデータ移行やレポーティングの正確性を高め、一貫性の確保を容易にする。

データガバナンスの観点では、データ辞書はセキュリティ要件、アクセス権、規制適合の基盤として機能する。

これを最新の状態に保つことで、組織は新しいデータ要件にも柔軟に対応でき、変化する事業ニーズに迅速に適応し、堅牢で応答性の高いシステム環境を維持できる。

コラム：データモデリング

データモデリングとは、情報システムやデータベースで使用されるデータを視覚的かつ論理的に整理し、その構造や関係をモデルとして表現するプロセスである。エンティティ（データオブジェクト）、属性（データ要素）、関係（データ同士

のつながり)を定義する。通常、概念データモデル(業務視点)、論理データモデル(技術視点の詳細構造)、物理データモデル(データベース構造を具体化したもの)の3層を含む。

データモデリングは、次のようにデータ品質向上に寄与する。

1. 一貫性と標準化を確保する
2. 重複と冗長性を排除する
3. データ完全性と整合性を確保する
4. データ再利用性を高める
5. データ駆動型意思決定を改善する

データモデリングは単なる技術プロセスではなく、データ品質マネジメントおよびデータガバナンスを支える重要な活動である。

コラム : DCAT

DCAT (Data Catalog Vocabulary) は、インターネット上で公開されるデータカタログ間の相互運用性を高めるために設計された W3C 標準語彙である。これにより、公開者はデータセットやデータカタログを標準化された方法で記述でき、プラットフォームをまたいで発見、共有、再利用しやすくなる。DCAT は、タイトル、説明、URL、ライセンスなどのメタデータをサポートする。一貫したメタデータを促進することで、DCAT はデータセットの発見容易性を高め、データ統合やオープンデータの取組を支援する。現行版の DCAT v3 は、データサービス、バージョン、来歴への対応を拡張し、FAIR データ原則との整合も図っている。

- FAIR : 発見可能 (Findable)、アクセス可能 (Accessible)、相互運用可能 (Interoperable)、再利用可能 (Reusable)

出典 : World Wide Web Consortium(W3C), 2024, [“Data Catalog Vocabulary \(DCAT\) - Version 3”](#)

4.1.3 データ取得

4.1.3.1 概要

- 説明
 - データ取得では、内部および外部ソースからのデータの生成、収集、取得を扱う。
 - このプロセスには、センサー由来であれ手動入力であれ、データ取得時の誤りを防ぐための措置が含まれる。
 - 外部ソースからデータを取得する際に、低品質データを除外するための検証活動を含む。
 - 現実世界のデータセットを合成データで補完する場合、このプロセスには、不適切な合成コンテンツの混入を防ぐための文書化と管理が含まれる。
 - 追跡可能性とデータ品質管理を支えるため、データプロファイルや来歴情報などのメタデータを付与することを含む。
- 課題
 - データ収集時に低品質または不適切なデータが混入するリスク
- 要件
 - 高品質なデータを効率的に収集し、不正確さを最小化するためのプロセス
 - 取得時点で品質を維持し、データクレンジングなど大規模な前処理の必要性を減らすための措置
- 価値
 - データ準備およびクレンジングに伴う時間とコストの削減
 - データへの信頼向上による、より良い意思決定と AI アプリケーションの信頼性向上

目的達成に必要なデータを生成・収集し、組織内に存在しない場合は外部から取得する。

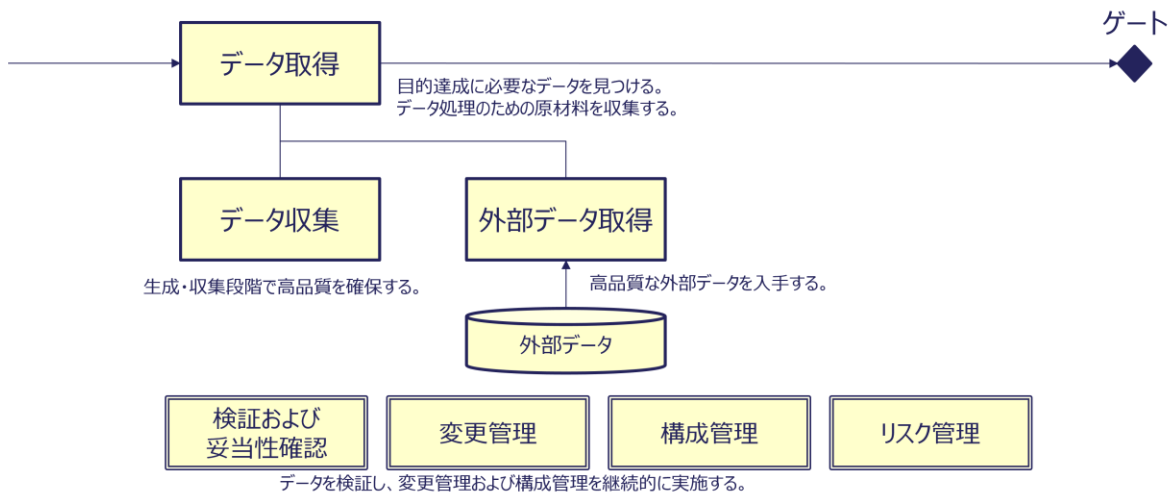


図 17. データ取得プロセス

4.1.3.2 データ取得

- 手順
 1. 必要なデータを特定し、所在を確認する。
 2. 来歴情報を確認する。
 3. 利用条件を確認する。
- チェックポイント
 - データは信頼できるソースから取得しているか。
 - データの来歴情報に問題はないか。
 - データの利用方法に制約はあるか。

4.1.3.3 データ収集

- 手順
 1. デバイス（センサー）を確認する。
 2. データを収集または入力する。
 - Web フォームや API を用いてエラーを防止する
 3. データを検証する。
 - 範囲外データや不適切データを除去する

- 不整合データを修正する
- 4. 匿名化およびマスキングを行う。
- 5. メタデータを作成する。
 - 一般情報
 - 品質情報
 - データの測定または収集方法
- チェックポイント
 - 不適切なデータを防止する対策を講じているか。
 - 範囲外のデータが入力されないようにしているか。
 - データの一貫性を確認しているか。
 - メタデータに DCAT を参照しているか。
 - データの測定方法および収集方法を記述しているか。

4.1.3.4 外部データ取得

- 手順
 - 1. メタデータおよび来歴情報を確認する。
 - 2. データを検証する。
 - 品質特性を確認する
 - 不適切なコンテンツ（合成データを含む）を検出する
 - 3. メタデータおよび来歴情報を追加する。
- チェックポイント
 - データソースは信頼できるか。
 - データの来歴情報は明確か。
 - データ品質は仕様を満たしているか。
 - 対象システムとの間で内容や意味を損なわずに効果的に交換・移転できるポータブルデータの要件を満たしているか。
 - システムの目的に照らして不適切なデータを含んでいないか。

- 取り込むデータ項目以外に重要情報（例：個人識別情報（Personally Identifiable Information：PII））を含んでいないか。
- 合成データであることが明確に識別されていない合成データを含んでいないか。

コラム：入力方法の改善

システム間をインタフェースで接続してデータを連携することで、利用者による入力ミスを防ぐことができる。これにより、次のような効果が期待される。

- 正確で迅速
- 伝達ミスがない

人がデータを入力する場合には、入力フォームを用いることで、入力中に利用者自身がデータを修正・調整できる。具体的には、次のことを促進できる。

- データ形式の確認
- データの正確性確認
- 入力ゆれの防止
- コントロールド・ボキャブラリの利用

コラム：来歴データ

来歴データとは、データの生成、変換、移動、利用を含む履歴と起源を記録する情報を指す。誰が、いつ、どのようにデータを作成、変更、利用したかを追跡する。来歴データは、データ処理の透明性を高めるメタデータとして機能し、データの信頼性と完全性が重要な医療、金融、科学研究などの分野で特に有用である。通常は、データのソース、処理履歴、適用された業務ルールやアルゴリズム、関連するシステムや利用者の情報を含む。

- 来歴データがデータ品質を改善する仕組み

- 来歴データは、その履歴を透明化することでデータの信頼性を確保する。
- データの起源とライフサイクルが明確になることで、正確性や適切性の検証、エラーや改ざんの特特定が容易になる。
- また、データ作成や更新の過程に関するインサイトを与え、不整合や不完全な記録の修正を可能にする。
- さらに、監査や標準適合性確認の過程では、来歴データはデータが適切に取り扱われたことの証拠となり、その信頼性への確信を高める。

これにより、データ分析や意思決定の正確性と有効性が高まり、潜在的リスクの低減とデータガバナンスの強化にもつながる。

4.1.4 データ準備

4.1.4.1 概要

- 説明
 - データ準備は、利用に適した品質水準でデータを利用可能な状態にするための最終的なプロセスである。
 - このプロセスには、準備したデータから誤りや範囲外データを除去し、分類体系、意味、粒度を調整するデータクレンジングが含まれる。
 - 必要に応じて不足データを補うことを含む。
 - クレンジングしたデータを統合し、統一されたデータセットを作成することを含む。
 - 必要に応じて、真正性の確認支援や不正利用抑止のためにデータへウォーターマークを付与することを含む。
 - 機械学習技術を用いたデータ拡張用の追加データ作成を含む場合がある。

- 課題
 - データ定義および品質の精緻化および標準化の必要性
 - 不明確または不完全なデータプロファイル
 - 不十分なデータ品質管理または重要なデータ要素の欠落
- 要件
 - 想定用途に適したデータ品質を確保するためのプロセス
 - データに関する包括的かつ正確なメタデータを作成するためのプロセス
 - データの出所および変換履歴を追跡するための来歴情報の管理
 - データの真正性および信頼性を保証するための措置
- 価値
 - 高品質な入力データの準備による、より正確で信頼できる処理
 - データライフサイクル全体を通じたデータ検証、確認、追跡可能性の向上

さまざまなデータセットを統合して入力データを作成し、それを AI モデルの学習や後続処理に利用する。

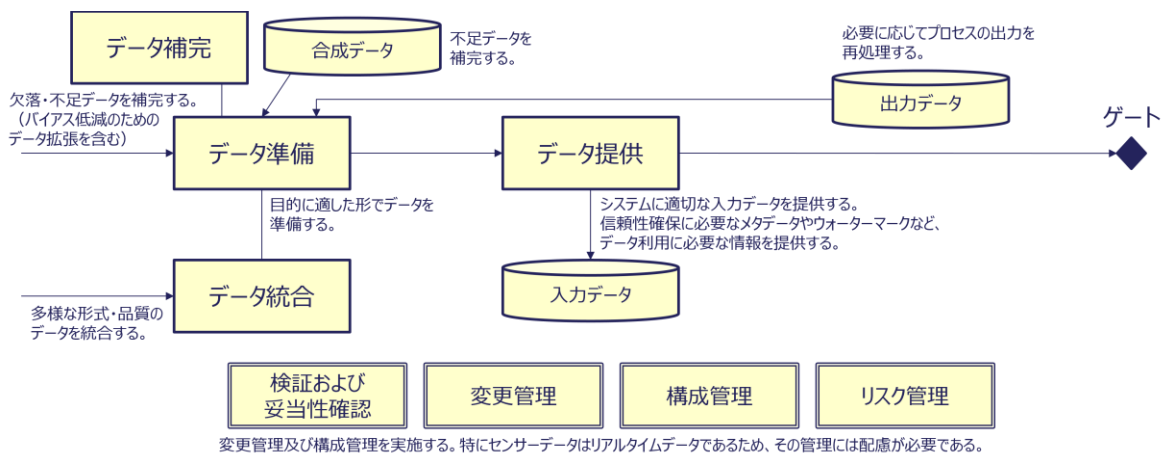


図 18. データ準備プロセス

4.1.4.2 データ準備

- 手順
 1. 収集したデータをドキュメント化する。
 2. データ統合方針を定義する。
 3. データ補完方針を定義する。
 4. データをクレンジングする。
 5. ラベルを追加する。
- チェックポイント
 - 必要なデータがすべてドキュメント化されているか。
 - データ統合方針が整備されているか。
 - データ補完に関する方針が定義されているか。
 - データセットは AI 学習および検証のために論理的に分割されているか。

4.1.4.3 データ統合

- 手順
 1. 統合後データモデルを定義する。
 2. 分類体系および管理語彙の変換表を作成する。
 3. データ項目の意味を確認し、マッチング方法を決定する。
 4. 精度と単位を調整する。
 5. ファイル形式を変換する（例：XML、JSON、CSV）。
 6. 意味的な差をドキュメント化する。
 7. 統合関連メタデータを作成する。
- チェックポイント
 - 変換表は、データ項目の意味を考慮して作成されているか。
 - 重複データに対する優先順位ルールはあるか。
 - 空欄データはデータなしとして扱っているか。

- データ変換はツールにより自動で行われているか。
- 統合データの一貫性を確認しているか。
- メタデータは DCAT 準拠で作成されているか。

4.1.4.4 データ補完

- 手順
 1. データの状態を分析し、補完が必要かどうかを判断する。
 2. 補完データを追加する。
- チェックポイント
 - データ補完に適したデータを使用しているか。

4.1.4.5 データ拡張 (AI システム向け)

- 手順
 1. 元データを準備する。
 2. 元データに基づいて学習データを生成する。
- チェックポイント
 - 元データにバイアスはないか。
 - 少数の元データへの依存による偏りはないか。

4.1.4.6 合成データ取得

- 手順
 1. メタデータおよび来歴情報を確認する。
 2. データの妥当性を確認する。
- チェックポイント
 - コンテンツが合成物であることが明確に示されているか。
 - 不適切なコンテンツまたは同意なく作成されたコンテンツを含んでいないか。

4.1.4.7 データ提供

- 手順
 1. カタログ情報を登録する。
 2. データインタフェースを提供する。
 3. バージョンを管理する。
 4. データサンプルを提供する。
 5. ウォーターマークなど、コンテンツ保護情報を追加する。
 6. アクセス制御を管理する。
 7. 必要に応じて利用状況を追跡する。
- チェックポイント
 - データは機械可読なインタフェースで提供されているか。
 - ウォーターマークなど、データの不正利用を困難にする仕組みを含んでいるか。
 - 不正利用を防ぐアクセス制御が整備されているか。

コラム：ID

データに一意的識別子（ID）を付与することは、データ品質向上に不可欠である。各レコードに ID を付与することで、複数ソースからの情報統合や、システム横断での一貫したデータ管理が可能になる。これにより、類似または関連するデータを統合・比較する際の混乱や重複を防げる。たとえば、2つの部門が別々のデータベースで顧客を管理していても、同じ顧客 ID を使っていれば2つのデータセットを正確に結合できる。

ID 付与は、データ統合と一貫性の支援に加え、次のような利点ももたらす。

1. 追跡可能性と監査可能性：一意的 ID により、データレコードの起源や履歴を追いやすくなり、監査やコンプライアンスの確認に役立つ。

2. 効率的な更新：特定のデータを更新または削除する必要があるとき、ID があれば対象レコードを迅速かつ正確に特定できる。
3. 拡張性：データ量が増えても、一意 ID があれば新規レコードを競合や混乱なく追加でき、大規模データベースでも明確性を保てる。
4. データガバナンスの強化：ID は明確な所有関係やガバナンスポリシーを支援、適切なチームやシステムに責任を割り当てやすくする。

総じて、ID の利用は、信頼できるデータ統合と一貫性を実現するだけでなく、データ全体の管理しやすさと価値を高める基本的な実践である。

コラム：データクレンジング

データクレンジングとは、データ品質と信頼性を高めるために、データセット内の誤り、不整合、不正確さを特定し、修正または除去するプロセスである。これにより、データが完全で、正確で、一貫した状態となり、分析、報告、意思決定に、より適したものになる。データクレンジングでは、欠損値、重複レコード、誤った形式、外れ値などの問題に対処する。とくに複数データセットを統合する場合や、高度な分析のためにデータを準備する場合には、データの完全性を維持するうえで重要である。

- データクレンジングの手法
 - 重複の除去：重複レコードを特定して除去し、冗長または誤解を招くデータを防ぐ。
 - 欠損データへの対応：必要に応じて補間、補完、削除などの手法で欠損値を埋める。
 - データ標準化：日付形式や大文字・小文字など、データが一貫した形式に従うようにする。
 - データ検証：事前定義ルールや参照データセットに照らしてデータを確認し、正確性と一貫性を確保する。

- エラー修正：タイプミス、不正な入力、値の不一致を特定して是正する。

効果的なデータクレンジングは高品質なデータを確保し、より良いインサイトと、より正確な意思決定につながる。

コラム：データマッチング

データマッチングとは、異なるデータセット間、または同一データセット内のレコードを比較・特定し、それらが同一の実体を表しているかを判断するプロセスである。重複排除、システム横断での顧客情報連携、不正検知、データ統合などの場面で広く利用される。名称、住所、電話番号、その他の識別子といった属性を評価し、多少の揺れや不整合があっても一致を検出する。マッチングの度合いはスコアで表されることもある。

重複または関連するレコードを統合することで、データマッチングはデータ品質を向上させ、正確な分析と報告を可能にする。

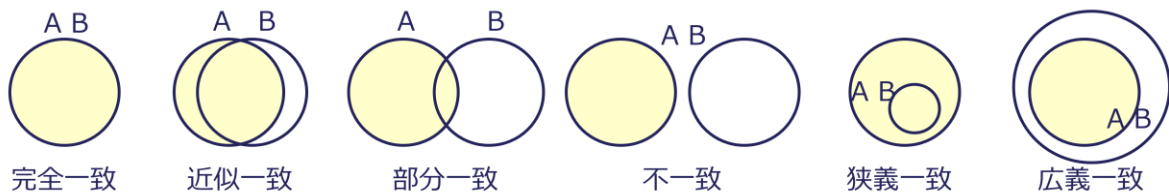


図 19. データマッチングのパターン

コラム：データ補完

データ補完には、データ拡張やデータエンリッチメントなどという考え方がある。

データ拡張とは、既存データにさまざまな変換を加えることで、データセットの規模と多様性を高める技法である。変換の例として、画像データであれば回転、反転、切り抜き、明るさ・コントラスト調整、ノイズ付加、テキストデータであ

れば言い換えや逆翻訳などがある。目的は学習データの多様性を高め、機械学習モデルの性能とロバスト性を向上させることである。

データエンリッチメントとは、外部または内部ソースから追加情報を補うことで既存データを強化するプロセスである。これにより、生データへ価値が付加され、より包括的で、正確で、分析や意思決定に有用なものになる。エンリッチメントには、欠損データポイントの追加、データ精度の向上、人口統計、地理、行動に関する文脈情報の統合などが含まれることが多い。

コラム：合成データ

合成データには、利用によって期待される効果がある一方、予想されるリスクもある。

効果

AI が生成する合成データの利用は、大きな効果と有益な機会をもたらす。

1. プライバシー保護の強化
 - 合成データは、実在の個人情報をさらすことなく分析やモデル学習を可能にし、個人のプライバシー保護に寄与する。
2. データ利用可能性の拡大と安全な検証
 - 合成データは、実データが希少、機微、収集困難な場合の不足を補い、管理された条件下で AI システムを安全に試験・検証することを可能にする。
3. AI 開発の加速
 - 多様で高品質な学習データを生成することで、合成データは AI の性能、ロバスト性、公平性の向上に寄与する。
4. コストとリスクの低減

- 合成データの利用は、データ収集やアノテーションのコストを下げるとともに、実データを扱うことに伴う規制リスクやセキュリティリスクを最小化する。

5. イノベーションと連携の促進

- オープンに共有される合成データセットは、機密性を維持しつつ、業界横断の連携や研究の加速を促進し得る。

リスク

AI が生成する合成コンテンツは、大きな影響を持ち得る。

1. フェイクニュースと偽誤情報の拡散

- 合成コンテンツは、非常に現実的な画像、音声、テキストを生成でき、虚偽情報の拡散に意図的に悪用される可能性がある。

2. プライバシー侵害

- ディープフェイクのような技術は、個人の顔や声を模倣するコンテンツを作成でき、そのプライバシーを侵害する可能性がある。

3. ブランドまたは企業評判への損害

- 悪意を持って作成された合成コンテンツは、企業やブランドの評判を損なう可能性がある。

4. 法的リスクと倫理的課題

- 合成コンテンツの作成と配布は、著作権侵害、パブリシティ権侵害、差別的コンテンツなど、法的・倫理的問題を引き起こし得る。

5. 信頼の低下

- 合成コンテンツがあふれることで、消費者や企業がデジタルコンテンツを信頼しにくくなる。

6. 国家安全保障および政治的リスク

- 合成コンテンツが政治的に悪用される可能性がある。

コラム：コンテンツ透明化技術

コンテンツ透明化技術とは、デジタルコンテンツに明確性、真正性、アカウントビリティをもたらすことを目的としたツール、フレームワーク、システムを指す。利用者が接する情報の起源、文脈、信頼性を理解できるようにすることを目指すものである。メディア、広告、ソーシャルプラットフォーム、データ活用産業などで、偽誤情報、偏ったコンテンツ、不透明なアルゴリズムといった課題への対応に用いられることが多い。

コンテンツ透明化技術の主な特徴は次のとおりである。

- ソース検証：コンテンツの起源や作成者を確認することで真正性を確保する。
- 追跡可能性：編集履歴、所有権、流通経路の明確な履歴を提供する。
- コンテンツラベリング：スポンサー提供、利用者生成、機械生成などを示すためにメタデータやマーカを付与する。

コラム：データ品質の表示

データ品質指標の水準は、内部管理向けと外部公開向けとで自然に異なるべきである。データ品質は、しばしば組織内部でデータ提供者や利用者によって詳細に管理される。しかし、データが社会でより広く流通する場合には、より単純で分かりやすい指標が必要になる。これは家電製品の扱いに似ている。製造や物流では詳細な仕様や品質確認が用いられるが、消費者が店頭で目にするのは省エネルギーや安全マークのような簡潔な指標である。同様に、データ品質マネジメントにおいても、目的や関与するステークホルダーに応じて、内部管理向けの詳細な指標と、外部共有や取引向けの簡潔な指標を設計することが重要である。

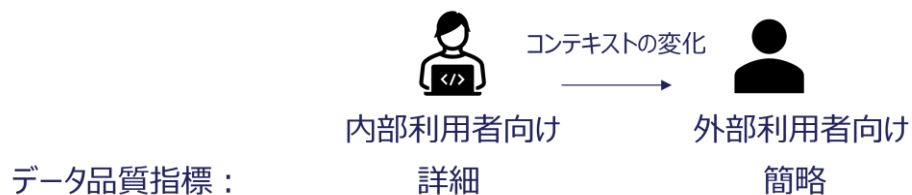


図 20. 内部管理用と外部利用用におけるデータ品質指標のコンテキストの変化

4.1.5 データ処理

4.1.5.1 概要

- 説明
 - データ処理は、データを効果的かつ効率的に処理することでサービス価値を創出する。
 - このプロセスには、エラー発生時にデータ経路の中で根本原因を特定し、実行可能なフィードバックを提供することが含まれる。
 - このワークフローでは、処理パイプライン全体を通じてデータ完全性とシームレスな運用を維持することを重視する。
- 課題
 - 処理パイプラインに支障を来すデータの不整合、不一致、不完全な記録
 - 処理済みデータの信頼性と正確性を維持することの難しさ
 - 堅牢な診断ツールや明確なプロセスモデルがない場合の、時間を要する根本原因特定
- 要件
 - エラーの根本原因を特定し対処するための包括的なプロセスと堅牢なデータモデル
 - 必要に応じたリアルタイム監視を含む、エラー解消の迅速化とサービス停止の最小化のための監視・診断の仕組み
- 価値
 - 効果的なエラー対応と高品質なデータ処理による、サービスの正確性と信頼性の向上
 - 顧客の信頼向上、運用リスクの低減、およびデータ駆動型イノベーションの支援

入力されたデータを処理する。処理中にデータ不整合が発生する可能性があるため、フィードバック機構が必要である。

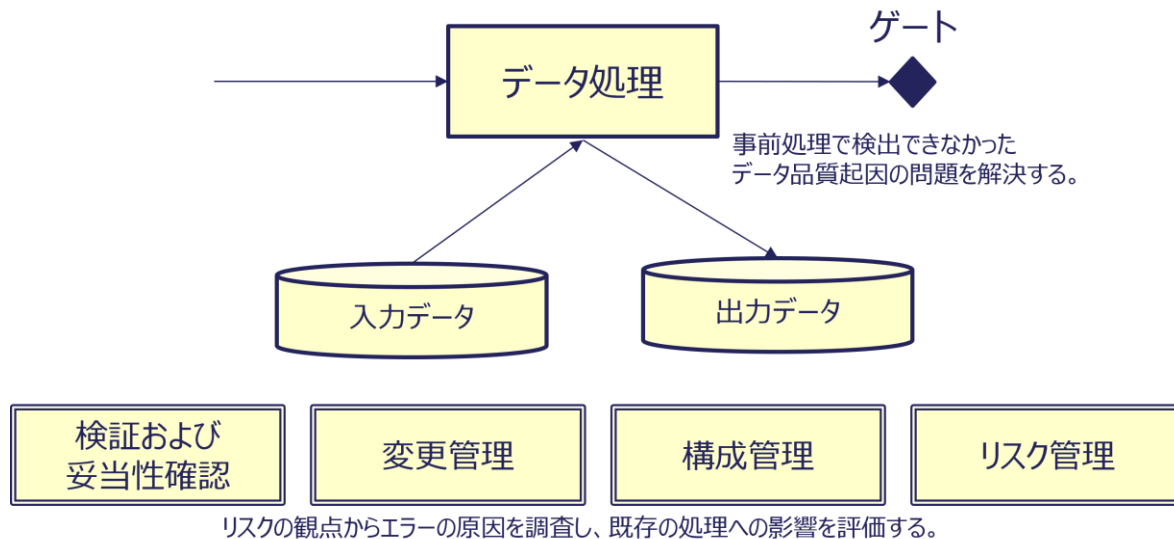


図 21. データ処理プロセス

4.1.5.2 データ処理

- 手順
 1. 一貫性を確認する。
 2. データを処理する。
 3. 処理エラーを報告する。
- チェックポイント
 - データ異常がある場合にエラー通知を提供しているか。
 - データ処理は検証のために可視化されているか。

4.1.6 AI システム

4.1.6.1 概要

- 説明

- AI システムプロセスでは、キュレーションされたデータセットを用いた AI モデルの学習を扱う。
- 学習した AI モデルの実装、利用、運用を含む。
- 出力評価に基づく継続的な再学習と改善を含む。
- たとえば、外部の情報源を参照して正確性と信頼性を高めるために、RAG を利用できる。
- 課題
 - データが限定的または不明確な場合に、虚偽または不正確な出力を含むハルシネーションが発生するリスク
 - 出力に社会的またはデータ由来のバイアスが表れ、公平性、包摂性、倫理上の課題を生じるリスク
 - 学習および再学習に必要な時間、資源、監督の負荷の高さによる拡張性の低下
- 要件
 - 利用者の問い合わせに対する回答の正確性と信頼性を向上させるためのプロセス
 - 信頼構築と意思決定プロセスの透明性確保のため、検証可能なソースや推論根拠などの参照情報を提示する仕組み
- 価値
 - 正確で迅速かつ文脈に即した応答による、より高いサービス水準
 - 信頼でき、偏りがなく、透明性が高い出力による、利用者やステークホルダーの AI システムへの信頼向上

目的適合性のある高品質データを入力することで、運用時の精度が向上する。

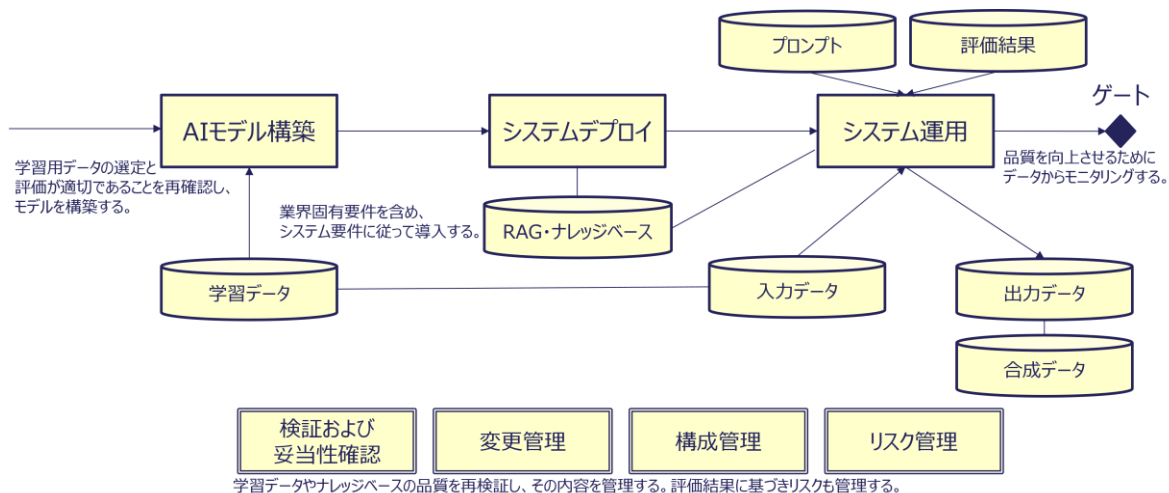


図 22. AI システムプロセス

4.1.6.2 AI モデル構築

- 手順

1. AI に求める水準を確認する。
2. 学習データを決定する。
 - バイアスのあるデータや不適切なデータを防ぐ
3. RAG を使うかどうかを決定する。
4. 個人情報または知的財産がないか確認する。
5. モデルを訓練する。
6. モデルをテストする。

- チェックポイント

- 学習データに信頼できるデータを使用しているか。
- 適切な場面で精度向上のために RAG を実装しているか。
- 個人データまたは知的財産につながり得る情報がある場合、不要な情報を学習データから除外するなど、表示を防ぐ措置を講じているか。

4.1.6.3 システムデプロイ

- 手順
 1. システムをデプロイする。
 2. 利用者教育を行う。
- チェックポイント
 - デプロイ方針を有し、それに従っているか。

4.1.6.4 システム運用

- 手順
 1. システムを運用する。
 2. 継続的に監視する。
- チェックポイント
 - 出力の再利用が AI システムへバイアスを持ち込んでいないか。

コラム：バイアス

AI を実装するうえで、バイアス低減は不可欠である。バイアスに対処することで、AI システムはより広い社会的受容を得やすくなり、倫理的な利用を確保し、長期的な信頼性と持続可能性を維持できる。

1. 不公平および差別の防止
 - バイアスを含むデータで学習した AI は、性別、年齢、人種、地域などに基づいて特定集団に不利な判断を行う可能性があり、倫理上の懸念を生む。
2. 信頼と公平性の確保
 - バイアスは AI システムへの信頼を損なう。信頼できる AI を構築するには、公平性と中立性を維持しなければならない。
3. 事業リスクの回避

- バイアスに関する問題は企業の評判を損ない、法的リスクにさらす可能性がある。
4. 社会的影響の最小化
- AI の普及は、その社会的影響力を増幅する。もし偏っていれば、不平等を悪化させたり、社会の分断を深めたりする可能性がある。
5. 法的・倫理的基準への適合
- AI に関する規制は、公平性やバイアス低減をますます重視している。非適合は罰則や特定市場からの排除につながる可能性がある。

コラム：モデル崩壊

モデル崩壊とは何か

生成 AI モデルが、現実の人間データではなく合成データで繰り返し学習されることで、時間の経過とともに多様性と正確性を徐々に失っていく現象である。

次の図は、合成データの反復利用によってモデル崩壊がどのように進行するかを示している。

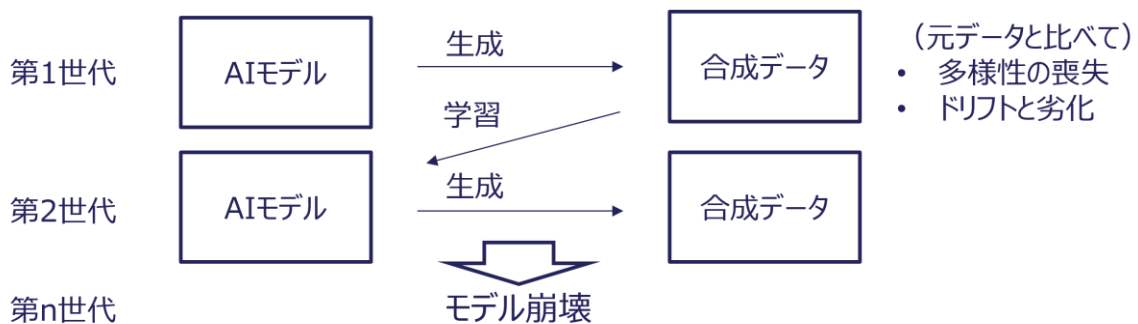


図 23. モデル崩壊のメカニズム

予防と緩和

- 実データと合成データを明確に区別し、適切に管理する。
- 検証済みの現実世界データセットを用いて定期的に再学習する。

- データセットの完全性のための継続的な品質監視とフィードバックループを実装する。

出典：Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, Yarin Gal, 2024, “[AI models collapse when trained on recursively generated data](#)”

4.1.7 出力評価

4.1.7.1 概要

- 説明
 - 出力評価は、AI システムの出力をエンドユーザーへ提供する前の最終確認として機能する。
 - このプロセスには、倫理面、情報の正確性、潜在的バイアス、意図した目的との整合について出力を評価することが含まれる。
 - 自動システムと人のレビュアーにより、信頼できる情報源と事前に設定した評価用データセットに照らして出力を評価する。
 - 特定されたエラーを、継続的改善のためにデータ保有者またはシステムオーナーへ報告することを含む。
- 課題
 - 利用者の信頼とシステム信頼性を低下させる不正確または不適切な応答のリスク
 - 不整合を引き起こす、権威ある定義や標準からの逸脱
 - エラーを検知、分類、修正する仕組みの不足
- 要件
 - 不正確または不適切な出力を防ぐため、正確性、倫理、バイアスに対処する堅牢な措置
 - 利用者確認を要する出力への明確なフラグ付け
 - 同じエラーの再発防止とシステム性能改善のためのフィードバックループ

- 価値
 - 信頼性、正確性、使いやすさに関する利用者の期待を満たす、または上回る出力による、より高いサービス水準
 - AI システムの倫理的かつ透明性が高い運用による、利用者、ステークホルダー、社会からの信頼と信用の向上

出力評価は、システムのセーフガードとして重要である。不適切な応答を防止する。

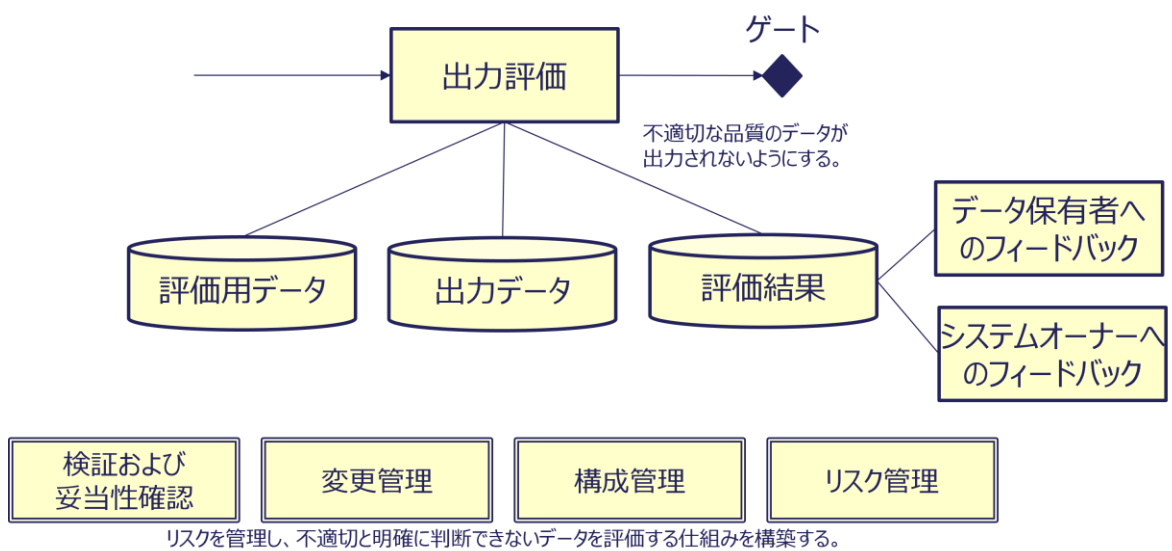


図 24. 出力評価プロセス

4.1.7.2 出力評価

- 手順
 1. 評価用データを準備する。
 2. 信頼できるデータに照らして出力を検証する。
 3. 評価用データに基づき不適切な応答を除去する。
 4. エラー報告を作成する。
 5. データ保有者へフィードバックする。
 6. システムオーナーへフィードバックする。

- チェックポイント
 - 非倫理的な応答を防ぐセーフガードは整備されているか。
 - 応答に対して判断が必要な場合、そのプロセスに人の介入が含まれているか。
 - 不適切な応答の原因を調査し、フィードバックする仕組みを有しているか。
 - 信頼できるデータと一致しない出力はないか。

コラム : ASoT

権威ある情報源 (ASoT) とは、特定のシステム、プロセス、文脈において、決定的かつ信頼される情報源を指す。ここは、正確で、最新で、完全なデータが維持される場所であり、その情報に依拠するシステムやチーム間の一貫性と信頼性を確保する。ASoT を用いてデータを検証することで、不正確なデータを排除し、データ品質を向上できる。政府が提供するオープンデータは、データ品質向上のための重要な社会基盤として機能する。



図 25. ASoT とオープンデータを用いたデータ検証のイメージ

4.1.8 結果提供

4.1.8.1 概要

- 説明
 - 結果提供では、処理済みデータと AI 出力を、明確で利用しやすい形でエンドユーザーへ提供することを扱う。

- このプロセスには、誤解を防ぎ、信頼と利用性を支える提供上の措置が含まれる。
- 文脈と信頼性を提供するためのメタデータ管理を含む。
- データ保護や知的財産保護のためのウォーターマーキングを含む場合がある。
- 課題
 - 不十分な提供による誤解のリスク
 - 非機械可読または非標準化インタフェースに起因する相互運用性の阻害
 - 適切な指針やセーフガードがない場合の出力の不適切利用のリスク
- 要件
 - 多様な利用者にとっての理解しやすさとアクセシビリティを支える提供プロセス
 - データ主権、所有権、管理権、規制適合を維持するためのプロセス
 - システム間の相互運用性を確保するための標準化された形式およびインタフェース
- 価値
 - 明確なデータ共有、十分な情報に基づく意思決定、信頼による AI の社会的効果の拡大
 - データ駆動型エコシステムにおける倫理的利用と持続的成長の支援

出力は、必要に応じて公開し、または他のシステムで利用できるよう提供することが望ましい。誤解およびその社会的影響を防ぐことが重要である。

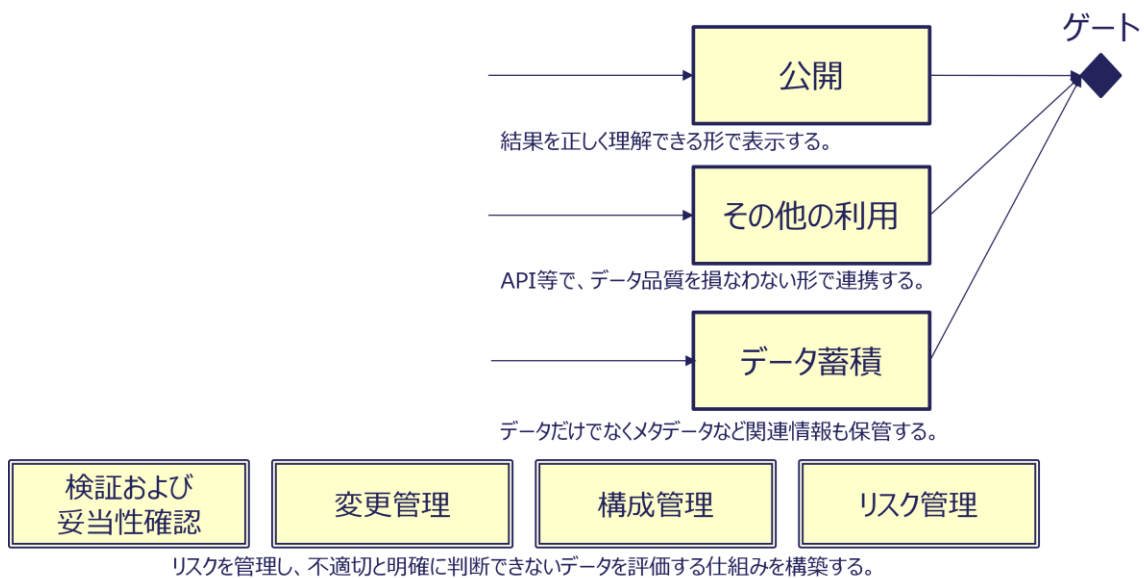


図 26. 結果提供プロセス

4.1.8.2 公開

- 手順

1. アクセス制御を管理する。
2. 処理結果を提示する。
 - 分かりやすさ
 - 可視化
3. 利用者の意見を収集する。

- チェックポイント

- 回答や表現は誤解を招かないか。
- 来歴情報を確認可能にするなど、回答の根拠を確認できる仕組みがあるか。

4.1.8.3 その他の利用

- 手順

1. アクセス制御を管理する。

2. API と関連情報を提供する。

- チェックポイント
 - 機械可読なインタフェースを提供しているか。

4.1.8.4 データ蓄積

- 手順
 1. 処理結果を保存する。
 2. データをバックアップする。
 3. データ保管を最適化する。
- チェックポイント
 - データが失われないよう保護措置を講じているか。

コラム：社会的・人的影響

低品質データは、システムに誤った情報を提供させ、利用者へ害を及ぼしたり、誤解させたり、不利益を与えたりする可能性がある。さらに、そのような誤情報は事故や、より広い経済的影響につながる場合もある。データ提供前に評価することに加え、何らかの悪影響が発生した場合の対策も準備しておくことが重要である。誤情報が公開された場合には、正確なデータを提供し検証できる信頼できる情報源を持つことが重要である。

4.1.9 データ廃棄

4.1.9.1 概要

- 説明
 - データ廃棄では、データ品質管理プロセスの正式な完了とデータの正式な廃棄扱いを扱う。
 - このプロセスには、誤用や意図しない利用を防ぐため、データが廃棄済みであることを明確に示す措置が含まれる。

- 古いデータによる誤りを防ぐことで、AI システムの完全性を支える。
- 課題
 - 利用者が廃棄済みまたは古いデータを知らずに利用し、不適切な結論や結果に至るリスク
- 要件
 - データが廃棄済みであることを利用者に明確に通知するためのシステムおよびプロセス
 - データの状態に関する曖昧さをなくすためのメタデータ更新、自動アラート、明確なドキュメント化
- 価値
 - 古いデータの利用によるエラーと手戻りの削減
 - 廃棄状態の明確な伝達による、AI システムの信頼性と信頼、運用効率、データガバナンスの向上

データ提供を停止する。廃棄を事前告知し、その後は保守や保証の対象外であることを明示する。

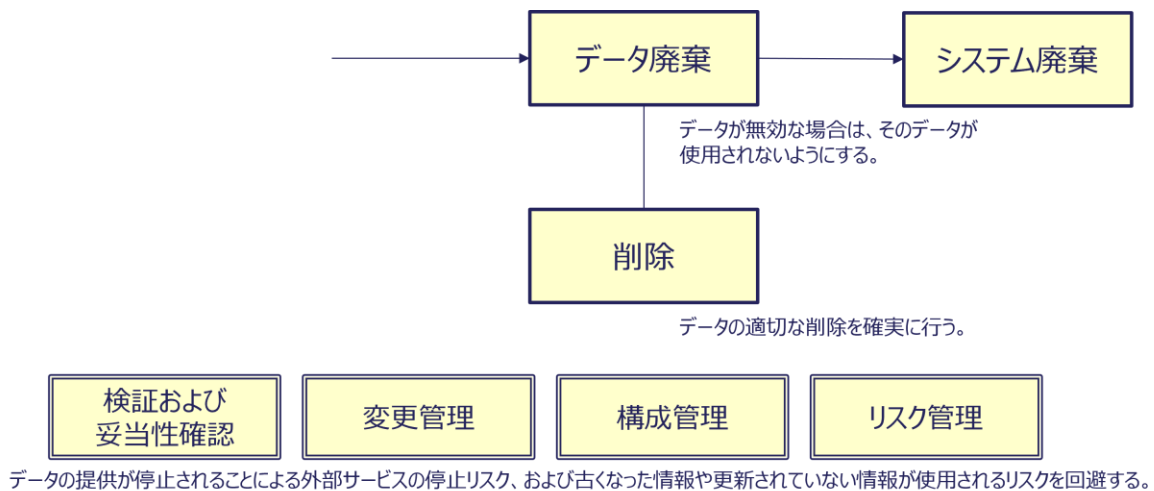


図 27. データ廃棄プロセス

4.1.9.2 データ廃棄

- 手順

1. 利用者にデータ廃棄を通知する。
 2. 廃棄済みデータに関する情報を提供する。
 3. データ移管時には、データ本体、メタデータ、関連文書も移管する。
- チェックポイント
 - データ停止前に十分な周知を行ったか。
 - 他者へ移管する場合、引継ぎに十分な情報を提供したか。

4.1.9.3 削除

- 手順
 1. 必要に応じてデータをアーカイブする。
 2. 復元できないようにデータを消去する。
 3. 一部削除する場合は事前に通知する。
- チェックポイント
 - 削除結果を確認したか。

4.1.9.4 システム廃棄

- 手順
 1. システム構成要素を廃棄する。
- チェックポイント
 - システム構成要素は適切に廃棄されているか。

4.1.10 ライフサイクル全体にわたるプロセス

4.1.10.1 概要

- 説明
 - ライフサイクル全体にわたるプロセスでは、データライフサイクル全体を通じた検証および妥当性確認、変更管理、構成管理、リスク管理を扱う。
 - これらの取組は、プロセス横断でのデータ品質の向上を支える。

- 課題
 - システムまたはサービス全体に対する管理の欠如による、問題発生時の見直しや迅速な対応の難しさ
- 要件
 - ステークホルダーにとって理解しやすい仕組みによるデータ品質関連情報の管理
 - データ品質関連情報の追跡可能性
- 価値
 - 標準適合性と透明性の向上
 - 透明性向上による信頼の向上

4.1.10.2 検証および妥当性確認

- 手順
 1. システムの目的に必要な特性、目標値、許容値を定義する。
 2. ライフサイクル全体を通じてシステムを管理する。
- チェックポイント
 - データ品質管理が現場の負担になっていないか。

4.1.10.3 変更管理

- 手順
 1. データ統合、処理、修正などの変更に関する基本方針を定義する。
 2. データ更新など、時間経過に伴うデータ劣化への方針を定義する。
 3. 変更を記録する。
- チェックポイント
 - 変更履歴を容易に参照できるか。

4.1.10.4 構成管理

- 手順

1. ソフトウェア構成を管理する。
 2. データ構成を管理する。
 3. 関連文書の構成を管理する。
- チェックポイント
 - ソフトウェア、データ、文書の一覧を管理しているか。

4.1.10.5 リスク管理

- 手順
 1. データ品質リスクと対応方針を定義する。
 2. アクセス制御と継続監視によりデータ品質リスクを管理する。
 3. 重要なリスク要因へ対応する。
 4. 事業継続計画（BCP）を策定する。
- チェックポイント
 - 自身のシステムまたはサービスのリスクを理解しているか。
 - リスクが発見されたとき、組織には根本原因から見直す文化があるか。

コラム：アクセス制御

アクセス制御は、データポイズニングやその他の不正な改変を防止し、データ品質を維持する上で不可欠である。どの利用者、デバイス、プロセスにデータへのアクセス、変更、削除を許可するかを定義し制御することで、組織は悪意ある攻撃や偶発的エラーのリスクを大幅に低減できる。効果的なアクセス制御の主要要素には、明確に定義された利用者ロールと権限の実装、多要素認証の強制、アクセスログの定期監査、異常または不正な活動の継続監視が含まれる。この構造化されたアプローチにより、データはライフサイクル全体を通じて正確で、一貫しており、安全な状態に保たれる。

コラム：AI によるデータ品質改善

AI とデータ品質は、相互に強化し合う関係にある。データ品質マネジメントが AI 性能を高める一方で、AI 自体もデータ品質向上のために活用できる。

AI 支援型のデータ品質マネジメントとは、人工知能を用いてデータエラーを検知し、修正し、防止することである。AI アルゴリズムは、データパターンを学習することで、不整合、欠損値、重複を自動で特定できる。また、改善提案を行い、時間経過に伴うデータ品質を監視し、変化するデータ環境に適応する。これにより、意思決定に必要な正確で、信頼でき、最新の情報を確保しやすくなる。

次の図は、AI とデータ品質の関係が一方方向ではなく、双方向で相互に強化し合う関係であることを示している。

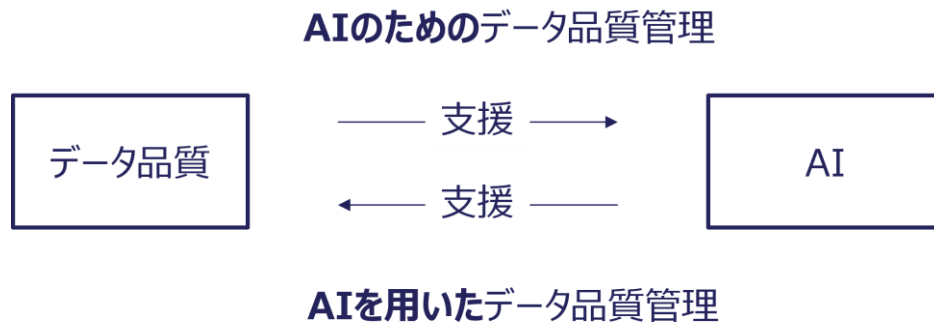


図 28. AI とデータ品質の相互関係

4.2 ガバナンスサイクルビュー

4.2.1 概要

データ品質を管理するには、個々の活動を正しく実施するとともに、それらの活動が持続可能で円滑に運用されるよう、組織的統制を確立する必要がある。

本ビューは、データ品質のガバナンスに着目するものであり、プロセスビューで定義されたデータライフサイクルプロセスが、組織として適切に統制され、監視

され、継続的に改善されていることを確保することを目的とする。本ビューでの「プロセス」とは、データライフサイクルプロセスの実行を統括・支援するマネジメントプロセスを指す。

マネジメントプロセスは、場合によってはライフサイクルプロセスと重複することがある。これは、検証やリスク管理などのプロセスが、実行面とガバナンス面の両方の性質を持つためである。

このビューでは、データ品質ガバナンス構造を ISO 8000-61 を参照している。

下図は、ISO 8000-61 で定義されるデータ品質マネジメントプロセスの基本構造を示している。継続的改善を実現するため、このプロセスは PDCA サイクルに従う。各フェーズ内のサブプロセスも示している。データ関連サポートプロセスは、実装を可能にする情報と技術を提供する。

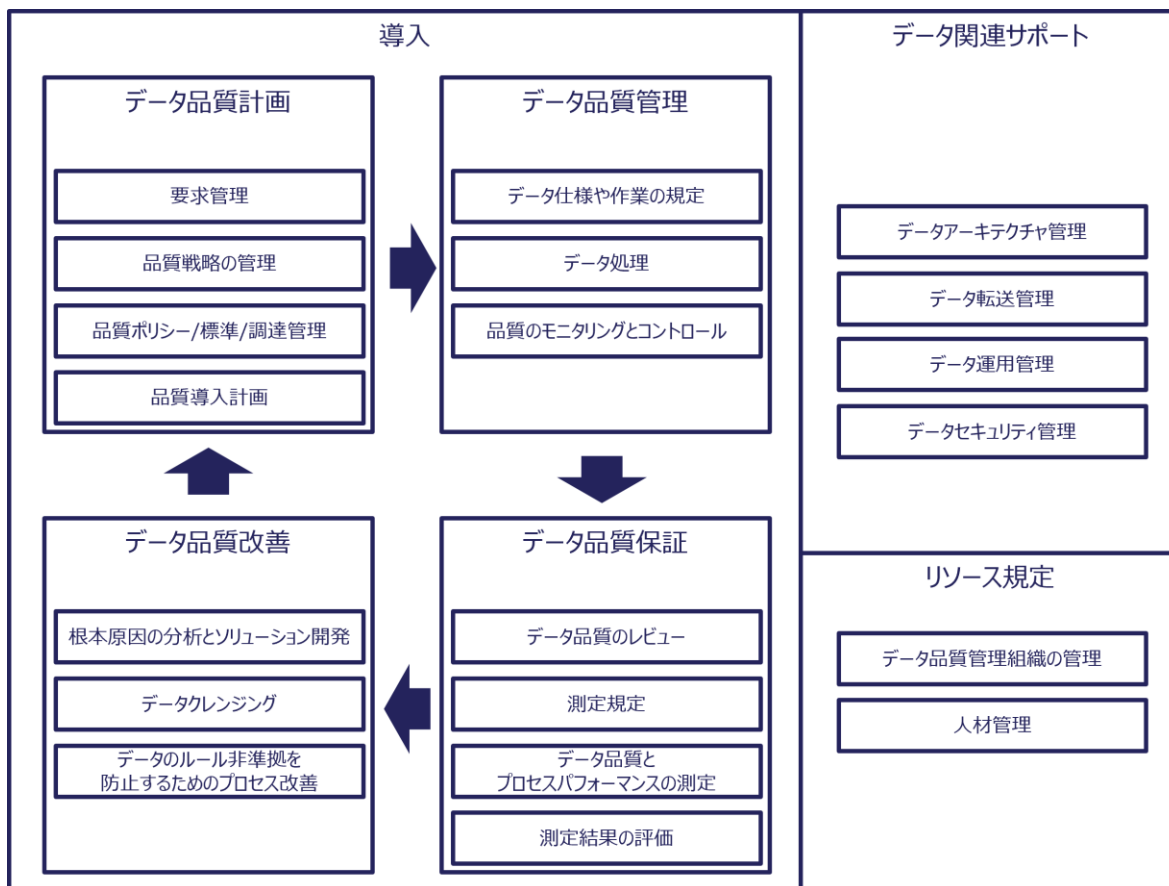


図 29. データ品質マネジメントのガバナンスサイクル

4.2.1.1 本節の構成

本節では、各プロセスについて、以下の2つの構成で整理する。

- 概要
 - 各マネジメントプロセスの概要を示す。これには、目的、役割、データ品質確保における価値に加え、関連するリスク、適用範囲、関係者、および実務上の留意点を含む。
- プロセス
 - 各マネジメントプロセスの中で実施されるサブプロセスを示す。併せて、それらに関連して実施される代表的な活動を示す。

4.2.2 データ品質計画

4.2.2.1 概要

データ品質計画は、ライフサイクル全体でデータ品質を管理するための組織的アプローチを定義する。AI システム向けの正確で、信頼でき、一貫したデータを支え、バイアス、エラー、不適切な意思決定などのリスクを低減する。このプロセスには、データ要件の特定、品質基準の定義、方針の策定、実行可能な計画の作成が含まれる。実務上の重要な留意点は、経営層がデータ品質の重要性を理解していることを確保することである。

4.2.2.2 要求管理

- 事業ニーズの特定：AI システムの具体的な目的と、その達成に必要なデータを定める。
- データ品質特性の定義：プロジェクトに関連する重要特性（例：正確性、完全性、一貫性、適時性）を定める。
- 現行データの評価：改善が必要な領域を特定するため、既存データのギャップ分析を行う。
- 要件の文書化：後続プロセスを導くため、データ品質要件を明確に文書化する。

4.2.2.3 品質戦略の管理

- ビジョンの策定：AI システムの目的と整合する長期的なデータ品質目標を定義する。
- 測定可能な目標の設定：時間経過に伴うデータ品質を評価する重要業績評価指標（KPI）を設定する。
- ステークホルダー連携：データエンジニア、アナリスト、意思決定者を含むすべてのステークホルダーとの整合を確保する。
- リスク管理：AI 活用における低品質データに伴う潜在リスクを特定し、緩和する。

4.2.2.4 品質ポリシー/標準/調達管理

- 方針の定義：データガバナンスおよび品質管理に関する組織方針を確立する。
- 標準の作成：一貫性確保のため、データ収集、保管、処理、検証に関する具体的標準を策定する。
- 手順の文書化：定期監査や検証プロトコルなど、データ品質の維持・改善に関する詳細手順を整理する。
- コンプライアンスの確認：すべての方針と標準が関連法令・規制要件に適合していることを確認する。

4.2.2.5 品質導入計画

- アクションプランの策定：スケジュールと責任を含む、データ品質対策の段階的实施計画を作成する。
- 役割の割当て：データ品質マネジメントに関わるチームメンバーの役割と責任を明確に定義する。
- ツールとプロセスの実装：データクレンジング、検証、監視、報告のためのツールを導入する。
- 監視と精緻化：フィードバックや性能指標に基づき、データ品質を継続的に監視し、実装計画を必要に応じて精緻化する。

4.2.3 データ品質管理

4.2.3.1 概要

データ品質管理は、AI システムで用いられるデータが、正確性、一貫性、信頼性などについて必要な品質水準を満たすことを目指す。AI 性能の最適化、バイアス低減、信頼できる結果の実現を支える。このプロセスには、明確な仕様、堅牢な処理技術、品質課題を特定し対処するための継続監視が含まれる。実務上の重要な留意点は、自動チェックなどの技術的措置を用いることで、現場に過度な負担をかけないことである。

4.2.3.2 データ仕様や作業の規定

- データ要件の定義：形式、構造、許容値範囲を含む明確なデータ要件を定義する。
- 作業指示の提供：統一性を確保するため、データ収集、ラベリング、処理に関する詳細指示を提供する。
- メタデータ指針の整備：データソース、タイムスタンプ、文脈情報を追跡するためのメタデータ作成指針を整備する。
- 検証チェックリストの作成：品質基準遵守を確認するため、提供者向け検証チェックリストを作成する。

4.2.3.3 データ処理

- データクレンジングと前処理：不整合、重複、不要情報を除去するため、データクレンジングと前処理を行う。
- データ正規化：システム間の互換性を確保するため、データを正規化する（例：形式変換、スケーリング）。
- アノテーションとラベリング：教師あり学習モデル向けに正確にアノテーションおよびラベリングを行う。
- 自動検証の実装：処理中の異常やエラーを特定するため、自動データ検証ツールを実装する。
- 処理の文書化：追跡可能性と再現可能性を維持するため、各処理ステップを文書化する。

4.2.3.4 データ品質のモニタリングとコントロール

- 品質指標の設定：主要品質指標（例：正確性、完全性、一貫性、適時性）を設定する。
- 監視の実装：データ品質問題を迅速に検出するため、必要に応じたリアルタイム監視を含む適時の監視システムを導入する。
- 定期監査のスケジューリング：データ完全性を確認・検証するため、定期監査のスケジュールを設定する。

- AI ツールの活用：データセット内のエラーパターンや潜在的バイアスを特定するため、AI ツールを活用する。
- フィードバックループの構築：発見事項に基づき、データ品質基準を継続的に改善するためのフィードバックループを構築する。

4.2.4 データ品質保証

4.2.4.1 概要

データ品質保証（DQA）は、AI 開発および運用に用いるデータが、正確性、一貫性、完全性、信頼性などに関する必要基準を満たしていることについて保証を提供する。AI システムにおけるエラー、バイアス、非効率の防止に役立つ。このプロセスには、潜在的なデータ品質問題の特定、評価基準の定義、データ品質の測定、結果分析による継続的改善の誘導が含まれる。

4.2.4.2 データ品質のレビュー

- 潜在的問題の特定：データセット内の欠損値、不整合、不正確さ、冗長性などの潜在的問題を特定する。
- 根本原因分析：収集、処理、保存におけるエラーを含め、データ品質問題の根本原因を分析する。
- 問題の文書化：既知の問題を文書化し、AI モデル性能への潜在的影響を評価する。

4.2.4.3 測定規定

- 測定基準の定義：データ品質特性（例：正確性、完全性、適時性、一貫性、妥当性）に対する明確で測定可能な基準を定義する。
- ベンチマークと閾値の設定：利用に適するとみなされるためにデータが満たすことを期待されるベンチマークと閾値を設定する。
- 基準の整合：AI モデル固有の要件および事業目標と基準を整合させる。

4.2.4.4 データ品質とプロセスパフォーマンスの測定

- データセット評価：設定した測定基準に照らして、データセットを体系的に評価する。

- データプロファイリング：データプロファイリングツールを用いて異常を検知し、品質属性を評価する。
- ワークフロー監視：品質基準への一貫した準拠を確保するため、データ処理ワークフローを監視する。

4.2.4.5 測定結果の評価

- 測定結果の分析：測定結果を分析し、ギャップと改善領域を特定する。
- 影響の定量化：データ品質問題が AI システム性能および意思決定へ与える影響を定量化する。
- レポートとダッシュボードの作成：発見事項をステークホルダーへ共有するため、レポートやダッシュボードを作成する。

4.2.5 データ品質改善

4.2.5.1 概要

データ品質改善は、データ品質問題に対処することで、AI システムの信頼性と有効性を向上させるマネジメントプロセスである。データ問題の特定と解消、データ一貫性の向上、時間の経過によるデータの劣化に対する予防措置の実装により、持続的な品質を支える。このプロセスには、データのルール非準拠を防止するための根本原因分析、ソリューション開発、データクレンジング、プロセス改善が含まれる。

単にデータを修正するだけでなく、持続的改善のためには、データ入力の難しさやプロセス自体の非効率といった根本原因に対処することが重要である。

4.2.5.2 根本原因の分析とソリューション開発

- 問題源の特定：誤入力、システムエラー、古いデータなど、データ品質問題の根本原因を調査する。
- 対象を絞った解決策の策定：ワークフロー更新、検証ルール改善、データ入力自動化など、特定した根本原因に対する是正策を設計・実装する。
- 結果の監視と検証：実装した解決策の有効性を継続的に評価し、問題が解消され再発しないことを確認する。

4.2.5.3 データクレンジング

- 不正確さの特定：データセット内の重複、不完全、不整合な記録を検出する。
- データ形式の標準化：日付形式、単位、命名規則などの一貫性を確保する。
- エラーの修正または除去：不正確な入力を修正し、欠損値を埋め、不要または古いデータを除去してデータ完全性を維持する。
- クレンジング処理の自動化：ツールやアルゴリズムを活用し、反復的なクレンジング作業を自動化して人的ミスを減らす。

4.2.5.4 データのルール非準拠を防止するためのプロセス改善

- データガバナンスポリシーの確立：品質問題を防ぐため、データ収集、管理、利用に関する明確な方針を定義し徹底する。
- データ検証メカニズムの強化：入力時または取り込み時に、必要に応じたリアルタイム検証を含む検証チェックを実装し、不適合を早期に特定する。
- ステークホルダー教育：データ品質維持のベストプラクティスと基準遵守の重要性について、従業員やデータ取扱者を教育する。
- データプロセスの監視と監査：データワークフローを定期的に評価し、不適合の潜在要因を検出し、既定方針への適合を確認する。

4.2.6 データ関連サポート

4.2.6.1 概要

データ関連サポートは、AI の開発および活用における効果的なデータ品質ガバナンスに必要な情報、技術、運用上の支援を提供する。ライフサイクル全体でデータの完全性、一貫性、信頼性を支える。このプロセスには、データアーキテクチャ管理、データ転送管理、データ運用管理、データセキュリティ管理が含まれる。これらの取組は、標準適合性、セキュリティ、運用上のニーズに対処しつつ、データ管理を最適化する枠組み、プロセス、ツールを整備する。

4.2.6.2 データアーキテクチャ管理

- データアーキテクチャとメタデータリポジトリの管理：一貫性確保のため、データモデル、スキーマ、標準を定義し、データの発見容易性と追跡可能性を高めるため、メタデータリポジトリを整備・運用する。
- データリネージ追跡：データフローと変換を監視するため、データリネージ追跡を実装する。
- 拡張性と柔軟性の確保：進化する AI ニーズに対応できる拡張性と柔軟性を確保する。
- アーキテクチャの整合：組織目標および標準適合要件とアーキテクチャを整合させる。

4.2.6.3 データ転送管理

- 安全なデータ転送プロトコルの整備：安全なデータ転送のためのプロトコル（例：暗号化、VPN）を整備する。
- データフロー監視と監査記録の維持：ボトルネック防止と適時の可用性を確保するため、必要な場合のリアルタイム可用性を含めてデータフローを監視・最適化し、アカウントビリティと追跡可能性のため、データ転送ログと監査記録を維持する。
- 越境データ転送方針の定義：法規制（例：GDPR、CCPA）に適合するよう、越境データ転送方針を定義する。
- データ転送の自動化：適用可能な場合は、手作業エラー削減のためデータ転送プロセスを自動化する。

4.2.6.4 データ運用管理

- データ検証とクレンジング：正確性維持のため、定期的なデータ検証とクレンジングを実施する。
- データ保管管理：アクセシビリティと拡張性を確保するよう、データ保管ソリューションを管理する。

- データワークフロー整備と性能監視：データ取り込み、処理、統合のワークフローを整備し、システム性能を監視し、データ運用に関する問題を解決する。
- データセットのバージョン管理：変更を追跡し一貫性を確保するため、データセットのバージョン管理を実施する。

4.2.6.5 データセキュリティ管理

- アクセス制御の実装：ロールベースアクセス制御（Role-Based Access Control：RBAC）や多要素認証を含むアクセス制御を実装する。
- セキュリティ監査と脆弱性評価：定期的なセキュリティ監査と脆弱性評価を実施する。
- 暗号化とインシデント監視：保管時および転送時データに対する暗号化プロトコルを導入し、不審な活動を監視し、セキュリティインシデントへ迅速に対応する。
- 規制適合の確保：国内外のデータ保護規制への適合を確保する。

4.2.7 リソース規定

4.2.7.1 概要

リソース規定は、効果的なデータ品質ガバナンスに必要な組織構造、人材、運用能力が整っていることを確保する。AI システムへの高品質なデータの提供を支援、リスクを最小化し、AI モデルの正確性、信頼性、公平性を高める。このプロセスには、チーム編成、スキルを持つ人材の配置、明確な役割と責任の定義が含まれる。

4.2.7.2 データ品質管理組織の管理

- データ品質ガバナンスチームの設置：専任のデータ品質ガバナンスチームまたは会議体を設置する。
- 役割と責任の定義：データ品質マネジメントに関する役割と責任（例：データスチュワード、データ品質アナリスト）を定義する。

- レポーティングラインとアカウントビリティの整備：データ品質問題に対する明確なレポーティングラインとアカウントビリティを整備する。
- 部門横断的な協働枠組みの構築：さまざまな部門のステークホルダーを巻き込む横断的協働の枠組みを構築する。
- 組織構造の見直しと更新：変化するデータおよび AI ニーズに適應できるように、組織構造を定期的に見直し更新する。

4.2.7.3 人材管理

- データ品質専門人材の採用：データガバナンス、データアーキテクチャ、AI 統合に精通したデータ品質専門人材を採用する。
- 研修プログラムの提供：従業員のデータ品質評価、クレンジング、検証に関するスキルを高める研修プログラムを提供する。
- キャリアパスと育成計画の定義：データガバナンス専門職の定着向上のため、キャリアパスと育成計画を定義する。
- 人員体制の確保：継続的なデータ品質監視と改善業務に十分な人員体制を確保する。
- データ品質文化の醸成：組織全体でデータ品質意識とアカウントビリティの文化を醸成する。

コラム：人材とチーム

データガバナンスチームは、AI チームと連携してデータ品質を向上させる。

下図は、AI チーム（CAIO および AI オフィス）とデータガバナンスチーム（CDO およびデータスチュワード）の協働における主要な役割を示している。AI とデータは相互依存しているため、その担当者も緊密に連携することが求められる。

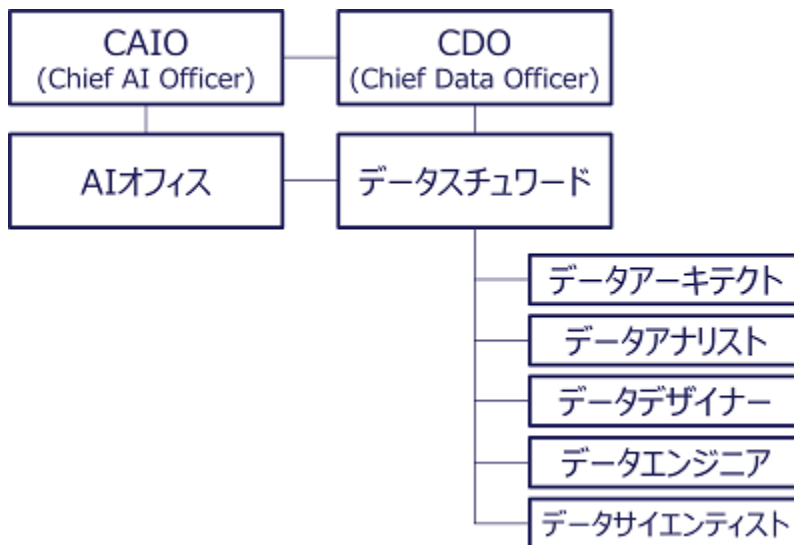


図 30. AI チームとデータガバナンスチームの連携

データを利用する事業部門でも、データ品質を含むリテラシー教育が必要である。

下図は、人材育成とリテラシー向上に向けた循環的なプロセスを示している。学習コンテンツは、体系的に整理された知識の集約として知識習得を支える。その知識はユースケースを通じて応用され、実践的な経験へとつながる。こうした経験はコミュニティ内で共有され、現場におけるさらなる行動を促す。現場での行動は、新たな学習コンテンツの創出や改善に寄与し、これによりサイクルが一巡する。スキル標準は、このプロセス全体を支える基盤として機能し、すべての段階において一貫性と継続的な発展を支える共通の土台を提供する。

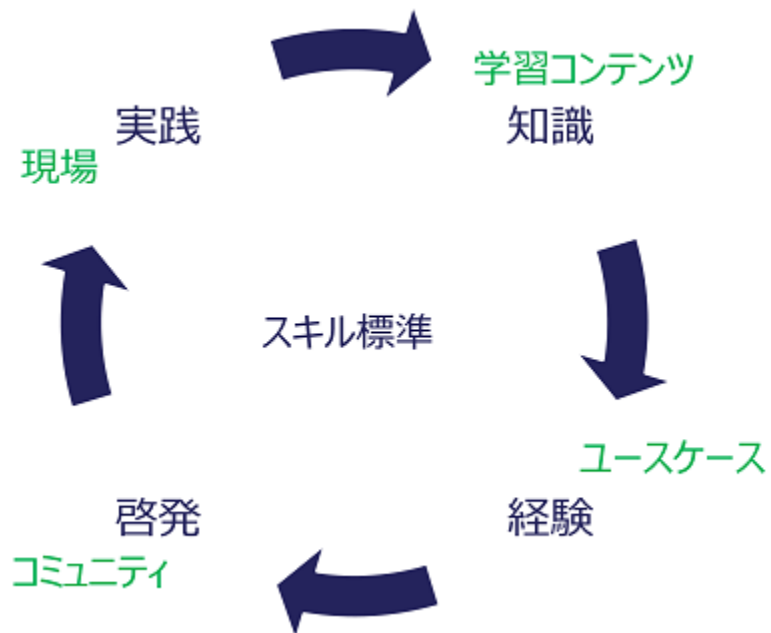


図 31. 人材育成サイクル

4.3 ゲートウェイビュー（品質特性）

4.3.1 概要

データ品質を維持するには、プロセスの境界やトランザクションの節目で確認する必要がある。センサーデータの場合は、定期点検がそれに当たることがある。確認の際には、データ品質を損なう可能性のあるデータ種別や事象を想定して、各特性を確認すべきである。現場の負荷を増やさないため、品質確認は可能な限り自動化すべきである。

4.3.1.1 本節の構成

ISO/IEC 25012 および ISO/IEC 5259-2 は、データ品質特性を、データ内容に関わる特性、システム面に関わる特性、その両面に関わる特性、AI 文脈に固有の追加特性に分類している。以下では、この分類に従って各特性を説明し、評価観点と不

適切な例を示す。加えて、センサーデータについてのデータ品質特性にも言及する。

- 評価観点
 - 各特性が満たされているかを評価するために用いる観点の例を示す。
- 不適切な例
 - その特性が適切に満たされていない代表的な状況例を示す。

4.3.2 データ固有の品質特性

4.3.2.1 正確性

正確性は、データが現実世界の値を正しく反映していることを示す。高い正確性は、AI が信頼できる結果を出し、誤った予測や不適切な意思決定を避けるうえで重要である。

- 評価観点
 - データが検証済みソースと一致している
 - エラー率が許容閾値内である
 - 外れ値に妥当な説明がある、または補正されている
- 不適切な例
 - 学習データセットの誤ラベル
 - 数値データの単位不一致
 - テキストデータのスペルミス

4.3.2.2 完全性

完全性は、必要なすべてのデータポイントが存在するかを測る。欠損値は AI 予測を歪め、システムの有効性を損なう可能性がある。

- 評価観点
 - 重要な項目がすべて埋まっている
 - Null 値または欠損率が最小限である

- 完全性確認が AI モデル要件と整合している
- 不適切な例
 - 利用者プロフィールにおける人口統計属性の欠落
 - 一部のみ存在する取引記録
 - 主要入力変数における Null 値

4.3.2.3 一貫性

一貫性は、データ値がデータセット間および時間軸で統一されていることを示す。不整合は、AI がパターンを誤って解釈する原因となる。

- 評価観点
 - 形式や単位が統一されている
 - データセットが時間を通じて同期している
 - 統合するデータ内に相反する記録がない
- 不適切な例
 - データセット内で異なる日付形式
 - 詳細が矛盾する重複レコード
 - 統合データセット内の不一致記録

4.3.2.4 信憑性

信憑性は、データソースの信頼性を測る。信頼できるソースは、AI モデルの妥当性と信頼を高める。

- 評価観点
 - 検証済みで評判のあるデータ起源を採用している
 - データ来歴が文書化されている
 - 査読済みまたは認証済みのソースを利用している
- 不適切な例
 - 不明または未検証のソース由来データ
 - 偽造または改ざんされたデータセット

- 検証されていない利用者生成コンテンツ

4.3.2.5 最新性

最新性は、データが想定用途に対して十分に新しいかを示す。古いデータは、時代遅れのインサイトや行動につながる可能性がある。

- 評価観点
 - データが最新の文脈と整合している
 - 時間的網羅性が適切である
 - 定期的更新が維持されている
- 不適切な例
 - 金融モデルにおける古い株価情報
 - 気候予測における古い気象データ
 - リアルタイムアプリケーションにおける過去の利用者嗜好

4.3.3 固有かつシステム依存のデータ品質特性

4.3.3.1 アクセシビリティ

アクセシビリティは、認可された利用者やシステムが障壁なくデータへアクセスできることを示す。適切な保管、堅牢な API、ユニバーサルデザイン原則を含む。アクセシビリティは、スムーズなデータフローを可能にすることで、AI モデル学習と活用を支える。

- 評価観点
 - データがプラットフォームやデバイスをまたいでアクセス可能である
 - API が十分に文書化され、エラーが少ない
 - あらゆる利用者層に対するアクセシビリティ標準を満たしている
- 不適切な例
 - 独自形式に閉じ込められたデータ
 - API ドキュメントの欠如

- ADA (Americans with Disabilities Act) などアクセシビリティ基準への非適合

4.3.3.2 標準適合性

標準適合性は、データ管理が GDPR や CCPA のような法令、規制、業界標準に従っていることを示す。適切な標準適合性は、法的リスクを回避し、プライバシーを保護し、信頼を構築する。これは AI システムの信頼性と社会的受容にとって重要である。

- 評価観点
 - 関連するすべての法的基準を満たしている
 - 適切なデータに関する同意取得の仕組みを実装している
 - 定期的にコンプライアンスチェックが行われている
- 不適切な例
 - 同意なしでのデータ収集
 - 管轄ごとのプライバシー法を無視
 - データ取扱いの監査証跡がない

4.3.3.3 機密性

機密性は、機微データが不正アクセスや漏えいから保護されていることを示す。重要な対策には、暗号化、アクセス制御、匿名化技術が含まれる。機密性を維持することは、信頼を守り、データ誤用に伴うリスクを低減する。

- 評価観点
 - 保管時および転送時のデータを強力に暗号化している
 - ロールベースのアクセス制御システムを利用する
 - セキュリティプロトコルは定期的に更新している
- 不適切な例
 - 機微データを平文で保存
 - 利用者同意なしに私的データを共有

- 脆弱または古いアクセス制御

4.3.3.4 効率性

データ管理における効率性とは、品質を損なうことなく処理時間や資源使用を最小化することである。効率的なデータは、AI 学習とデプロイの高速化、コスト削減、拡張性向上につながる。

- 評価観点
 - データ保管構造を最適化する
 - データアクセス時の遅延を最小限にする
 - 効率的な抽出・変換・格納（ETL）パイプラインを構築する
- 不適切な例
 - 冗長なデータ処理ステップ
 - API 呼び出しの大きな遅延
 - 単純作業に対する過剰な資源消費

4.3.3.5 精度

精度は、データに含まれる詳細度や粒度の程度を示す。精度が高いほど、より正確で詳細な情報を取得できる。

- 評価観点
 - データ値が十分な詳細度で記録されている
 - 単位と尺度が明確に指定されている
 - 時間情報および空間情報が適切な粒度で記録されている
- 不適切な例
 - 粗すぎる緯度経度の値
 - 精度水準が不一致で、計算が困難な値
 - 過度に高い解像度により、処理や保管の負荷が増大

4.3.3.6 追跡可能性

追跡可能性は、データの起源、変換、利用を追跡し、アカウントビリティと再現性を確保する。詳細なログとメタデータは透明性を高め、デバッグやコンプライアンス確保にとって重要である。

- 評価観点
 - データリネージ文書が包括的である
 - 変換ログが維持されている
 - 追跡のための一意 ID がある
- 不適切な例
 - データセットのソース詳細欠落
 - 記録されていないデータ変更
 - 一貫しないデータセットのバージョン管理

4.3.3.7 理解性

理解性は、データが人にも機械にも正しく解釈できることを示す。明確なラベル、メタデータ、直感的な構造は利用性を高め、効果的な AI 学習に重要である。

- 評価観点
 - データが明確かつ十分にラベル付けされている
 - メタデータがスキーマ標準と整合している
 - 一貫性があり論理的なデータ構造となっている
- 不適切な例
 - 曖昧または欠落したラベル
 - 複雑で文書化されていない構造
 - 誤解を招くメタデータ注記

4.3.4 システム依存のデータ品質

4.3.4.1 可用性

可用性は、必要なときにいつでも AI データへアクセスできることを示す。信頼できるシステムは停止時間を最小化し、継続運用を確保する。これはリアルタイム AI アプリケーションにとって重要である。

- 評価観点
 - 稼働率がサービス水準合意（SLA）を満たしている
 - 単一障害点を防ぐ冗長構成となっている
 - アクセシビリティ問題に対する定期監視とアラートがある
- 不適切な例
 - 頻繁なサーバ停止によるデータアクセス中断
 - バックアップがなく、ハードウェア障害時にデータへアクセスが不可能
 - 重要作業中のアクセス問題解消の遅延

4.3.4.2 移植性

移植性とは、AI データをプラットフォーム、システム、環境間で互換性の問題なく円滑に移転できる能力を指す。柔軟性と適応性を確保する。

- 評価観点
 - データが広く受け入れられている形式（例：CSV、JSON）で保管されている
 - データ交換に標準化 API を利用している
 - データ移行に十分な文書がある
- 不適切な例
 - 専用ソフトを必要とする独自形式
 - システム間で一貫しないデータ構造
 - 移行時の誤解を招くメタデータ欠如

4.3.4.3 回復性

回復性は、予期しない障害や故障の後に、AI データを迅速かつ正確に復旧するシステムの能力に焦点を当てる。データ損失を最小限に抑えることを目的とする。

- 評価観点
 - 複数の復旧ポイントを持つ定期バックアップを取得している
 - 定期的に試験される災害復旧計画がある
 - 単一障害点を避ける冗長構成のストレージとなっている
- 不適切な例
 - 古いバックアップによる不可逆なデータ損失
 - 多大な手作業を要する復旧プロセス
 - 復旧計画が試験されておらず、遅延が発生

4.3.5 AI/ML 向け追加データ品質特性

4.3.5.1 監査可能性

監査可能性は、データの追跡可能性と、データ収集・利用プロセスをレビューできる能力を確保する。アカウントビリティおよび倫理的・法的基準への適合を可能にする。追跡可能性と異なり、監査可能性は、データ取扱いがレビューおよび検証可能かに焦点を当てる。

- 評価観点
 - データソースの明確な文書化がされている
 - データ収集プロセスを明確に定義している
 - データリネージ記録が利用可能である
- 不適切な例
 - データソースに関するメタデータ欠如
 - 曖昧なデータ来歴
 - 主要なデータプロセスに関するログへのアクセス不可

4.3.5.2 均衡性

均衡性は、データがすべての関連カテゴリや結果を適切な比率で表現し、AI モデルのバイアスを最小化することを示す。

- 評価観点
 - カテゴリの分布が利用目的に照らして適切に均衡している
 - 過剰サンプリングまたは不足サンプリングを回避している
 - データセット間で一貫性がある
- 不適切な例
 - データセットにおける性別の偏り
 - 地理的地域の偏った表現
 - ある年齢層の過剰表現

4.3.5.3 多様性

多様性は、AI システムにおける汎化能力を向上させるため、データセットが幅広い視点、シナリオおよびバリエーションを含んでいることを確保するものである。

- 評価観点
 - 異なる文化的文脈への網羅性がある
 - 多様なシナリオと人口統計属性を含んでいる
 - 言語表現に幅がある
- 不適切な例
 - 少数派方言の除外
 - 多言語環境での均質なデータ
 - 多様な環境要因の無視

4.3.5.4 有効性

有効性は、そのデータセットが特定の AI タスクに利用可能かを示す。データがタスク要件を満たしているかに焦点を当てる。

- 評価観点
 - 画像解像度など、十分な入力品質がある
 - カテゴリごとに十分なサンプル数がある
 - 利用可能なラベルまたはアノテーションがある
- 不適切な例
 - 少なすぎる主要サンプル数
 - タスク上重要な入力の欠如
 - 入力ノイズが過剰

4.3.5.5 識別可能性

識別可能性は、データから個人を識別できるかどうかを示す。これは直接的にも、属性の組み合わせによっても起こり得る。機密性と異なり、識別可能性は、データから個人が識別できるかに焦点を当てる。

- 評価観点
 - 直接識別子が存在しない
 - 連結による識別リスクが低い
 - 再識別リスクの評価が行われている
- 不適切な例
 - 識別可能な個人情報の保持
 - 不十分なデータマスキング
 - 可逆的な匿名化手法

4.3.5.6 関連性

関連性は、データが想定用途に適合しているかを示す。意味のある情報を含み、不要な情報を除くべきである。有効性と異なり、関連性は、そのデータが想定用途にとって意味を持つかに焦点を当てる。

- 評価観点
 - AI モデルの対象領域に関連する特徴量が含まれている

- 無関係な変数が除外されている
- 冗長な情報が避けられている
- 不適切な例
 - 対象外の記録の混在
 - 無関係な特徴量
 - 重要でない変数に過度に注目

4.3.5.7 代表性

代表性は、データが AI モデルが遭遇する現実世界の条件や状況を反映しているかを示す。本番条件との乖離が大きいと、モデル性能が低下する可能性がある。代表性は、現実世界データが必ずしも均等分布しないため、均衡性とは異なる場合がある。

- 評価観点
 - 対象母集団の特性と整合している
 - 想定条件への網羅性がある
 - サンプルングバイアスを回避している
- 不適切な例
 - 全国調査で都市部人口を過剰サンプリング
 - まれだが重要な条件を無視
 - 地理的カバレッジが不完全

4.3.5.8 類似性

類似性は、サンプルが過度に似通っていないかを示す。過剰な類似性は汎化性能を低下させる。

- 評価観点
 - 十分なサンプル変動がある
 - 近似重複が限定的である
 - 密集したパターン集中がない

- 不適切な例
 - 重複サンプルの過多
 - すべての測定が同じ装置由来
 - 似た文書の焼き直しが大半

4.3.5.9 適時性

適時性は、データが利用に間に合うタイミングで利用可能になるかを示す。たとえ正しいデータでも、到着が遅すぎれば価値を失う可能性がある。最新性と異なり、適時性はデータが「十分に新しいか」ではなく「利用時点に間に合うか」に焦点を当てる。

- 評価観点
 - データが期限内に到着する
 - 到着遅延が許容範囲内である
 - 更新頻度が適切である
- 不適切な例
 - 必要な意思決定の締切後に届くデータ
 - ばらつきが大きい遅延
 - 遅すぎる更新サイクル

4.3.6 センサーデータ品質特性

4.3.6.1 センサーデータ

センサーデータは、温度、湿度、位置、加速度などの物理量をセンサーによって測定し、デジタル形式で取得されたデータを指す。主にリアルタイムで収集・蓄積され、装置制御や分析、サービス提供に活用されており、IoT やスマートシステムの普及が進む現代社会において不可欠な情報基盤となっている。

センサーデータは 1 件あたりは小さいが、リアルタイムに集約すると大量データとなり、その多くは装置やサービスに組み込まれて安全性にも直接関わる。また、センサーは屋外や機器内部など多様な環境で利用され、設置条件（測定位置の高

さ、使用デバイス、測定方法など)の違いによってデータにばらつきが生じるほか、外部環境の影響も受けやすい。このため、センサー種別や利用条件に応じたデータ収集管理が不可欠である。一方で、多数のセンサーが存在する場合には、冗長性を活用し、個別の故障を周辺センサーや時系列データで補完できる場合もある。

センサーデータの品質特性は、ISO/IEC 25012 および ISO/IEC 5259-2 で整理された品質特性をセンサーデータ向けに読み替えて適用できる。加えて、収集するデータの書式や意味的情報以外にも、データを発生させるセンサーデバイスの状態による品質を把握しておく必要がある。

以上を踏まえ、前項までに示した品質特性の中でも特に次の 4 特性が重点的な品質管理対象となる。

- 正確性
- 完全性
- 一貫性
- 精度

4.3.6.2 デバイス依存品質測定量

前述の通り、センサーデータの品質にはデバイスの状態が大きく影響するため、品質評価においてはデバイスに注目した品質測定量を定義することが有効である。「センシングデータの品質レベル評価のためのガイドライン」では、それらを以下の通りに定義している。

| 区分 | デバイス依存の品質測定 | |
|------|-------------|---------------------------------------|
| | 量 | 説明 |
| 設計情報 | デバイスの情報 | デバイスに入力された物理量（光、音など）の計測原理、処理方式等の把握レベル |
| | 故障のしにくさ | デバイスの稼働レベル |

| 区分 | デバイス依存の品質測定量 | 説明 |
|-------|--------------|----------------------|
| | 耐久性 | 寿命部品の低下レベル |
| | セキュリティの対策 | セキュリティ対策の実施レベル |
| | 通信の安定性 | 通信が途絶、遅延なく動作するレベル |
| 設置・調整 | 設置方法の適切さ | 条件にあった適切な設置の実施レベル |
| 運用・保守 | システムの安定稼働 | 安定稼働の計画レベル |
| | システムの環境監視 | 設置状況の把握レベル |
| | アップデートの適切さ | 適切なソフトウェアバージョンの運用レベル |

出典：一般社団法人データ社会推進協議会（DSA），2024，ホワイトペーパー“センシングデータの品質レベル評価のためのガイドライン策定に向けた検討”

4.3.6.3 センサーデータの品質に影響を与える異常

センサーデータは時間経過に伴って変化し、データ特性も変わることがあるため、補正措置が必要になる。ISO 8000-210 では、データ分析者がデータ品質の評価や異常検知・補正を行う上で有効な情報として、異常の種類を「単独センサー」と「複数センサー」に分けて整理している。

(1) 単独センサーにおける異常

単一センサーのデータ系列に現れる典型的な異常。

- オフセット
 - 真値からの一定のずれ。
- ドリフト

- 時間とともに生じる緩やかな変化。
- トリム
 - 誤差を補正するための調整。
- スパイク
 - 突発的で短時間の跳ね上がり。
- ノイズ
 - データ中のランダムな変動。
- データ損失
 - データポイントの欠落または空白。
- データ不足
 - 収集データ量の不足。
- シフト
 - ベースラインの急激な変化。
- 急減または急増
 - 急激な低下または増加。
- スタック
 - 同一値の繰り返し出力。
- 有界振動
 - 規則的で限定的な変動。
- 不規則な頻度
 - 不規則なデータ間隔（サンプリング周期の乱れ）。
- 解像度の差異
 - 期待仕様との乖離を含むデータ粒度の問題。
- 不正確なタイムスタンプ
 - 記録時刻のずれ。
- レイテンシー
 - 事象発生から記録までの遅延。

(2) 複数センサーにおける異常

複数の関連センサー間の関係として現れる異常。

- 不一致
 - 同一対象を測定しているはずのセンサー間で値に有意な差異がある状態。
- 関係違反
 - センサー間で成立すべき関係や制約に違反（例：高温側より低温側の値が高い）。
- 時刻不整合
 - 同期しているはずのセンサー間で時刻情報が不一致。

4.3.6.4 異常の背景要因および運用管理項目

以下は ISO 8000-210 では直接の異常類型として定義されていないが、センサーの運用および複数センサーの統合においてデータ品質に大きく影響する要因である。なお、本項には前述の異常と概念的に対応する内容が含まれるが、ここではそれらを「異常そのもの」ではなく、その発生要因および運用上の管理対象として整理している。

1. 配置・空間管理
 - センサー位置の不適切さ：設置位置が不適切な場合、測定対象を正しく表現できない状態が生じる。
2. ノイズと干渉
 - 環境ノイズ：周辺環境から不要な信号を受け、データ品質の低下が生じる。
 - クロストーク：センサー間の干渉により誤った読取りが生じる。
3. 解像度とサンプリングに関する課題
 - 異なる粒度：センサーごとに記録精度や解像度が異なる状態が生じる。

- サンプリングレート差：センサー間でサンプリング頻度が一致しない状態が生じる。
4. 故障および健全性管理
- センサードリフト：経年劣化により測定値の変動が生じる。
 - スタック・停止センサー：同一値の繰り返しや測定停止の状態が生じる。
 - データ欠落：故障や通信エラーによりデータの欠損が生じる。
5. 校正管理
- バイアス誤差：校正不良によりオフセットやスケーリング誤差が生じる。
 - センサー間の比較可能性の確保が不十分な状態が生じる。
6. 冗長性とデータ統合
- 重複データ：センサーのカバレッジ重複により冗長なデータが生じる。
 - 矛盾データ：センサー間で相反する値が生じる。
7. 統合上の課題
- 異なるプロトコル：通信方式の違いによりデータ統合が困難な状態が生じる。
 - 異種データ形式：データ形式の違いにより標準化が必要な状態が生じる。
8. 環境影響
- 温度、湿度、圧力：環境条件の違いによりセンサー性能の変動が生じる。

4.3.6.5 クラウド・エッジ・IoT における処理

IoT デバイスやエッジで収集されたデータは、エッジ AI などのデバイス上、データ集約ポイント、あるいは大量データ処理のためにクラウドへ集められて処理されることがあり、それぞれの場所でデータ品質対策が必要となる。

- クラウド：大量データ処理の過程で、異常データや偏ったデータを検出する。
- 集約ポイント：地域単位などでデータを集約する。必要に応じてデータ変換や統合を行い、一部の指示をエッジへ返すこともある。
- エッジ：データクレンジング、認識、匿名化などの処理をエッジ側で行う。センサー固有のオフセット補正などを行う場合もある。



図 32. クラウド、エッジ、IoT をまたぐデータ処理

5 結び

5.1 メッセージ

人々の関心は AI の先端技術に向きがちですが、AI を社会で持続可能かつ安全に活用するためには、データ品質を適切に管理することが重要です。

常に「Garbage in, Garbage out.」の原則を念頭に置き、AI の価値を最大化しましょう。

6 文書情報

6.1 発行主体・作成体制

AI セーフティ・インスティテュート (AISI) は、AI および AI を活用したイノベーションの加速を支える基盤として、安全性確保を担う政府の取組である。情報処理推進機構 (IPA) は、デジタル技術の中核とする政府系機関であり、AISI の活動に参画している。

本書は、AISI の標準チームと、IPA デジタル基盤センターのデータ専門家チームが共同で作成し、AISI のデータ品質サブワーキンググループからも貴重な貢献を受けている。

本書は継続的に更新される文書であり、読者からのフィードバックを歓迎する。

6.2 参考文献

- ISO/IEC 25012: SQuaRE, Data quality model
- ISO/IEC 25024: SQuaRE, Measurement of data quality
- ISO 8000: Data quality
- ISO/IEC 5259: Data quality for analytics and machine learning (ML)
- ISO/IEC 8183: Data life cycle framework
- ISO/IEC 38505-1: Information technology - Governance of IT - Governance of data
- ISO 19157: Geographic information, Data quality
- DAMA-DMBOK (2nd edition, 2017), DAMA International
- データ連携基盤を通して提供されるデータの品質管理ガイドブック, 内閣府
 - https://www.chisou.go.jp/tiiki/kokusentoc/supercity/supercity_230926_guidebook.html
- ホワイトペーパー「センシングデータの品質レベル評価のためのガイドライン策定に向けた検討」, 一般社団法人データ社会推進協議会 (DSA)
 - <https://data-society-alliance.org/survey-research/data-quality-evaluation-standards/>
- 機械学習品質マネジメントガイドライン 第4版, 国立研究開発法人産業技術総合研究所 (AIST)
 - <https://www.digiarc.aist.go.jp/publication/aiqm/>

6.3 改版履歴

- 1.02 版 (2026-05-14)
 - アクセシビリティ、利用性、保守性を高めるため、文書形式をプレゼンテーションベースからテキストベースへ変換した。
 - 見出しや階層を含む文書構造を再編し、一貫性と明確性を確保した。
 - テキストベース形式での可読性向上のため、説明を調整・追加した。また、リスト、表、図などを調整した。
 - 全体の内容や意図を変えない範囲で、一貫性維持のための軽微な修正、記述の正確性向上を行った。
 - データ品質特性に関する記述の正確性を改善した。
 - 上記更新に加え、機械翻訳による日本語参考版を公開した。手動修正は、意味上不可欠な不整合への対応のみに最小限とした。
- 1.01 版 (2025-12-02)
 - 全体の内容に影響を与えない範囲で、正確性と可読性向上のため表現を見直し、ページ順を調整した。
 - 本編前に「本ガイドブックについて」を追加した。
 - 1.00 版で「IV.実装 (Implementation)」として掲載されていた内容は限定的であったため、該当部分をセクション III に統合した。
 - 今後の更新では、「実装」セクションを再編・拡充する予定である。
- 1.00 版 (2025-03-31)
 - 初版を公開した。

6.4 日本語参考訳について

本ガイドブックは、英語版が正式版であり、日本語版は英語版の内容理解を補助するための参考訳です。内容の理解の一助としてご活用いただけますが、正確な情報については以下の原文もあわせてご確認ください。なお、ISO/IEC 等の標準に由来する用語の日本語訳は、本資料内での理解を助けるために便宜的に付したも

のであり、必ずしも各標準における正式な日本語訳を意味するものではありません。

https://aisi.go.jp/output/output_framework/data_quality_management_guidebook/

7 付録 1

7.1 機械学習品質マネジメントガイドライン 第 4 版

「機械学習品質マネジメントガイドライン」は、AI 品質向上とリスク低減を目的として、機械学習システムの品質基準および目標水準を定義している。主に、AI 搭載製品・サービスの提供者およびシステム開発者を対象としており、品質に関する理解の共有と可視化を通じて、適切な受発注と高品質システムの評価を促進することを目指している。

このガイドラインでは、品質を「利用時品質」「外部品質」「内部品質」の 3 層で整理し、内部品質の向上が外部品質の達成につながり、それによって利用時品質の実現が可能になると位置づけている。以下の 2 つの図は、本ガイドラインにおける品質達成の構造と、内部品質の構造および特性を示している。データ品質は内部品質の一部として扱われている。

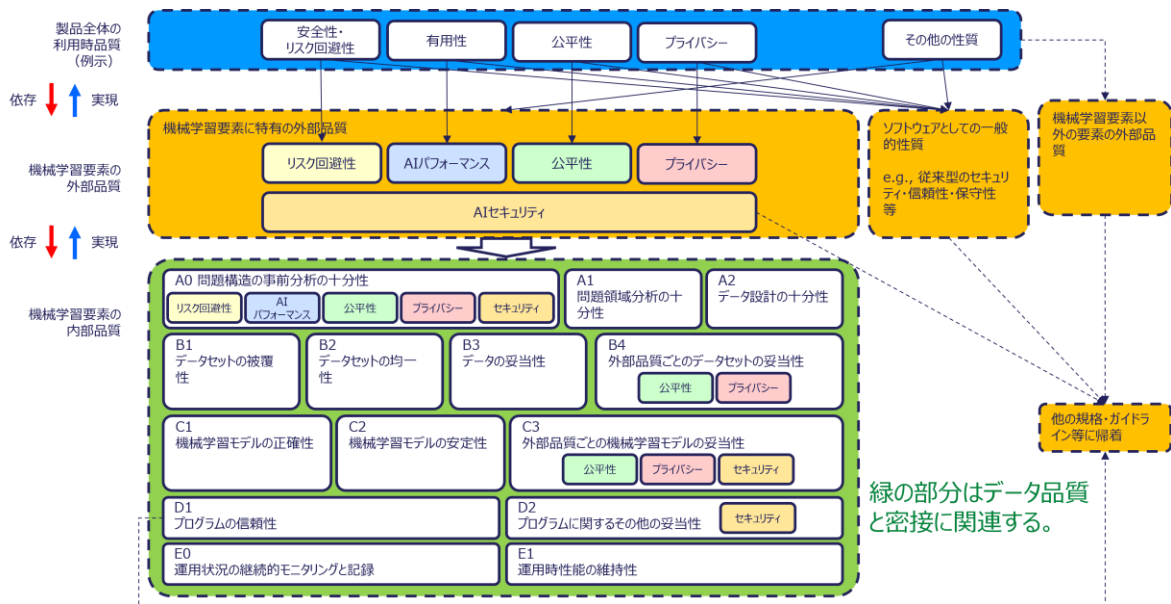


図 33. 製品品質達成の構造

データは AI における中核要素である。本ガイドラインは、機械学習を推進する際に考慮すべきデータ品質マネジメント上の課題と対応について指針を提供する。

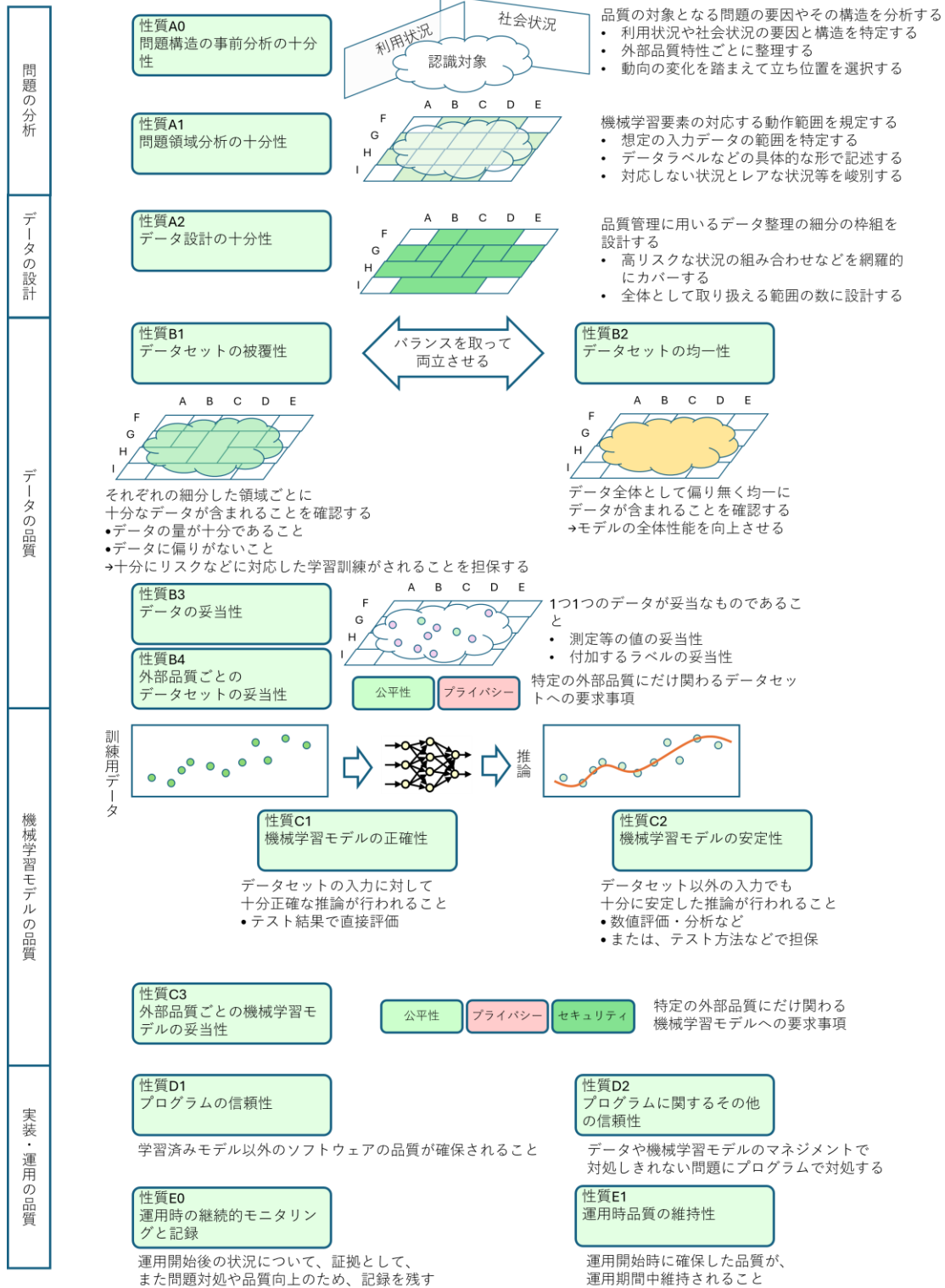


図 34. 内部品質の特性

出典：国立研究開発法人産業技術総合研究所（AIST），2023, “機械学習品質マネジメントガイドライン”

8 免責事項

本ガイドブックに記載した情報の正確性および信頼性の確保には最大限努めているが、その内容について明示または黙示を問わず、いかなる保証も行わない。いずれの組織も、本ガイドブックの利用に起因して生じたいかなる損失、損害、請求についても責任を負わない。すべての情報は現状有姿で提供されており、その内容に基づいて行われる一切の行為については、読者が全面的な責任を負うものとする。