

ビッグデータ教材

－ 第2章 －

資料

本ビッグデータ基礎教材 第2章資料は、文部科学省の教育政策推進事業委託費による委託事業として、一般社団法人全国専門学校情報教育協会が実施した令和6年度「地方やデジタル分野における専修学校理系転換等推進事業」の成果物です。

目次

2-1. データ分析の全体プロセス	2
2-2. ビジネス課題の明確化と分析設計	7
2-3. データの種類と尺度	12
2-4. 質的データと量的データの違い	16
2-5. 基本統計量：平均値と中央値	21
2-6. 基本統計量：分散と標準偏差	26
2-7. 基本統計量：最頻値とパーセンタイル	30
2-8. データの分布とヒストグラム	35
2-9. 箱ひげ図による外れ値の検出	40
2-10. 散布図による 2 変数の関係把握	44
2-11. 相関係数の計算と解釈	49
2-12. クロス集計表の作成と分析	54
2-13. データクレンジングの基本	58
2-14. 欠損値の処理方法	63
2-15. 外れ値の処理方法	68
2-16. データの正規化と標準化	72
2-17. カテゴリデータのエンコーディング	77
2-18. 時系列データの基本的な扱い方	82
2-19. データのサンプリング手法	86
2-20. 仮説検定の基本的な考え方	91
2-21. t 検定の実施方法	96
2-22. カイ二乗検定の実施方法	100
2-23. 単回帰分析の基礎	105
2-24. 重回帰分析の基礎	110
2-25. クラスタリング分析の基礎	114
2-26. 主成分分析の基礎	119
2-27. データ分析結果の解釈方法	124
2-28. 分析結果の可視化と報告	128

ビッグデータ基礎

- 第1章 ビッグデータとその活用
- 第2章 データ分析基礎**
- 第3章 スクリプト言語による分析
- 第4章 ツールとモニタリング
- 第5章 列指向ストレージによる高速化
- 第6章 データマートの基本構造
- 第7章 大規模分散処理



Mirai no tobira

第2章 データ分析基礎

ビッグデータ基礎

●学習の目的

データ分析の基本的な考え方と手法を習得し、実務で必要となる統計的な概念を理解する。データの収集から前処理、基礎的な分析手法まで、データ分析の全体像を把握し、実践的なスキルの基礎を身につける。



第2章 データ分析基礎

ビッグデータ基礎

- | | |
|------------------------|------------------------|
| 2-1. データ分析の全体プロセス | 2-15. 外れ値の処理方法 |
| 2-2. ビジネス課題の明確化と分析設計 | 2-16. データの正規化と標準化 |
| 2-3. データの種類と尺度 | 2-17. カテゴリデータのエンコーディング |
| 2-4. 質的データと量的データの違い | 2-18. 時系列データの基本的な扱い方 |
| 2-5. 基本統計量：平均値と中央値 | 2-19. データのサンプリング手法 |
| 2-6. 基本統計量：分散と標準偏差 | 2-20. 仮説検定の基本的な考え方 |
| 2-7. 基本統計量：最頻値とパーセンタイル | 2-21. t検定の実施方法 |
| 2-8. データの分布とヒストグラム | 2-22. カイ二乗検定の実施方法 |
| 2-9. 箱ひげ図による外れ値の検出 | 2-23. 単回帰分析の基礎 |
| 2-10. 散布図による2変数の関係把握 | 2-24. 重回帰分析の基礎 |
| 2-11. 相関係数の計算と解釈 | 2-25. クラスター分析の基礎 |
| 2-12. クロス集計表の作成と分析 | 2-26. 主成分分析の基礎 |
| 2-13. データクレンジングの基本 | 2-27. データ分析結果の解釈方法 |
| 2-14. 欠損値の処理方法 | 2-28. 分析結果の可視化と報告 |

●学習内容の概要

データ分析の進め方、基本統計量の理解、データの可視化手法、相関分析などの基礎的な統計手法について学習する。また、データクレンジングや前処理の重要性についても理解を深める。

第2章 データ分析基礎

ビッグデータ基礎

はじめに

データ分析は、ビジネスや研究において重要な意思決定を支える基盤となっています。しかし、データから意味のある洞察を得るためには、適切な分析手法の理解と、体系的なアプローチが必要です。

本章では、データ分析の基礎となる考え方と手法について学んでいきます。まず、分析の全体プロセスを理解し、ビジネス課題をデータ分析の視点で捉える方法を学びます。次に、基本統計量の計算から、データの可視化、相関分析まで、基礎的な統計手法を習得します。

第2章 データ分析基礎

ビッグデータ基礎

はじめに

また、実務で重要となるデータクレンジングや前処理の技術、欠損値や外れ値への対処方法についても詳しく解説します。さらに、回帰分析やクラスタリングなど、より発展的な分析手法の基礎についても学習します。

データ分析は、単なる数値の計算ではありません。データの特性を理解し、適切な手法を選択し、結果を正しく解釈する能力が求められます。この章で学ぶ基礎的な知識は、より高度なデータ分析技術を習得するための重要な土台となります。

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

プロセス管理と品質確保



ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

データ分析の基本ステップ

ビジネス課題の定義から施策実行までのデータ分析の基本的な流れ

課題定義フェーズ

ビジネス課題を分析可能な形に落とし込む方法目標設定のプロセス

データ収集フェーズ

必要なデータの特特定と収集方法の選定、データソースの評価

データ分析は、複数の重要なステップで構成される体系的なプロセスです。まず「課題定義」から始まり、「データ収集」「データ準備」「分析」「評価・適用」という基本的なステップを経て進められます。これらのステップを順序立てて実行することで、効果的な分析が可能となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

データ分析の基本ステップ

ビジネス課題の定義から施策実行までのデータ分析の基本的な流れ

課題定義フェーズ

ビジネス課題を分析可能な形に落とし込む方法目標設定のプロセス

データ収集フェーズ

必要なデータの特特定と収集方法の選定、データソースの評価

課題定義フェーズでは、ビジネス上の課題を分析可能な形に具体化します。例えば「売上を増やしたい」という漠然とした課題を、「どの商品の、どの顧客層の、どの時期の売上を、どの程度増やすのか」といった具体的な目標に落とし込みます。この段階での明確な定義が、その後の分析の方向性を決定づけます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

データ分析の基本ステップ

ビジネス課題の定義から施策実行までのデータ分析の基本的な流れ

課題定義フェーズ

ビジネス課題を分析可能な形に落とし込む方法目標設定のプロセス

データ収集フェーズ

必要なデータの特特定と収集方法の選定、データソースの評価

データ収集フェーズでは、設定した課題を解決するために必要なデータを特定し、収集します。社内システムのデータベース、外部データ、新規に収集するデータなど、利用可能なデータソースを評価し、最適な収集方法を決定します。データの品質や信頼性の確認も、この段階で重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

データ準備フェーズ

データクレンジング、加工、統合など分析前の準備作業

分析フェーズ

探索的データ分析から高度な統計分析まで分析手法の選択と実行

評価・適用フェーズ

分析結果の評価とビジネスへの適用プロセス

データ準備フェーズは、収集したデータを分析可能な状態に整えるプロセスです。データのクリーニング、欠損値や異常値の処理、データ形式の統一化、必要な加工や集計など、多岐にわたる準備作業を行います。この作業の質が、分析結果の信頼性に大きく影響します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

データ準備フェーズ

データクレンジング、加工、統合など分析前の準備作業

分析フェーズ

探索的データ分析から高度な統計分析まで分析手法の選択と実行

評価・適用フェーズ

分析結果の評価とビジネスへの適用プロセス

分析フェーズでは、準備したデータに対して適切な分析手法を適用します。まず探索的なデータ分析で全体像を把握し、その後、統計的な分析や機械学習など、目的に応じた手法を選択して分析を進めます。必要に応じて、複数の手法を組み合わせることもあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

データ準備フェーズ

データクレンジング、加工、統合など分析前の準備作業

分析フェーズ

探索的データ分析から高度な統計分析まで分析手法の選択と実行

評価・適用フェーズ

分析結果の評価とビジネスへの適用プロセス

評価・適用フェーズでは、得られた分析結果を評価し、実際のビジネスへの適用を検討します。分析結果の統計的な有意性、ビジネス上の実現可能性、期待される効果など、多角的な評価を行います。また、具体的な施策への落とし込みや、実行計画の立案も行います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

データ分析プロセスの全体像

プロセス管理と品質確保



ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

プロセス管理と品質確保

プロジェクト計画

データ分析プロジェクトの計画立案タイムライン設定の方法

リソース配分

人材、ツール、時間などのリソース配分と管理方法

品質管理ポイント

各フェーズでの品質チェックポイントと管理方法

プロジェクト計画は、分析の成否を左右する重要な要素です。目的、スコープ、期間、必要なリソース、期待される成果などを明確に定義し、具体的なタイムラインを設定します。特に、中間マイルストーンの設定により、進捗の可視化と管理が容易になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

プロセス管理と品質確保

プロジェクト計画

データ分析プロジェクトの計画立案タイムライン設定の方法

リソース配分

人材、ツール、時間などのリソース配分と管理方法

品質管理ポイント

各フェーズでの品質チェックポイントと管理方法

リソース配分では、プロジェクトに必要な様々なリソースを適切に管理します。データサイエンティストやドメインエキスパートなどの人材、分析ツールやコンピューティングリソース、予算、時間など、必要なリソースを特定し、効率的な配分を計画します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

プロセス管理と品質確保

プロジェクト計画

データ分析プロジェクトの計画立案タイムライン設定の方法

リソース配分

人材、ツール、時間などのリソース配分と管理方法

品質管理ポイント

各フェーズでの品質チェックポイントと管理方法

品質管理は各フェーズで重要です。データ収集段階でのデータ品質チェック、前処理段階での整合性確認、分析段階での妥当性検証など、それぞれの段階で適切な品質管理ポイントを設定します。チェックリストの活用や、定期的なレビューの実施も効果的です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

プロセス管理と品質確保

ドキュメント管理

分析プロセスの記録と文書化の重要性と方法

進捗管理手法

マイルストーン管理やアジャイル手法効果的な進捗管理

成果物管理

分析結果やモデル成果物の管理と版管理

ドキュメント管理も忘れてはならない重要な要素です。分析の目的、使用したデータ、前処理の内容、分析手法、パラメータ設定、結果の解釈など、プロセス全体を通じて適切な文書化を行います。これにより、分析の再現性が確保され、後の改善や展開も容易になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

プロセス管理と品質確保

ドキュメント管理

分析プロセスの記録と文書化の重要性と方法

進捗管理手法

マイルストーン管理やアジャイル手法効果的な進捗管理

成果物管理

分析結果やモデル成果物の管理と版管理

進捗管理では、従来型のウォーターフォール型管理やアジャイル型の管理など、プロジェクトの特性に応じた手法を選択します。特に、データ分析プロジェクトでは、試行錯誤が必要なことも多いため、柔軟な進捗管理が求められます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-1. データ分析の全体プロセス

第2章 データ分析基礎

プロセス管理と品質確保

ドキュメント管理

分析プロセスの記録と文書化の重要性と方法

進捗管理手法

マイルストーン管理やアジャイル手法効果的な進捗管理

成果物管理

分析結果やモデル成果物の管理と版管理

成果物管理では、分析結果やモデル、ドキュメントなどの成果物を適切に管理します。バージョン管理を徹底し、どの時点のデータを使用して、どのような設定で分析を行ったのかを追跡可能にします。また、成果物の保管場所や共有方法についても、ルールを定めて運用します。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

分析設計の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

課題の構造化

複雑なビジネス課題を具体的な分析課題に分解する方法

KPIの設定

課題解決の成果を測定するための適切なKPIの設定方法

ステークホルダー分析

関係者の特定と要求事項の整理、合意形成のプロセス

課題の構造化は、複雑なビジネス課題を扱いやすい単位に分解するプロセスです。例えば、「顧客満足度の向上」という大きな課題を、「商品の品質改善」「配送時間の短縮」「カスタマーサポートの改善」といった具体的な要素に分解します。この構造化により、データ分析の焦点を絞ることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

課題の構造化

複雑なビジネス課題を具体的な分析課題に分解する方法

KPIの設定

課題解決の成果を測定するための適切なKPIの設定方法

ステークホルダー分析

関係者の特定と要求事項の整理、合意形成のプロセス

KPIの設定は、課題解決の効果を測定するために重要です。例えば、「売上向上」という課題に対して、「月間売上高」「顧客単価」「リピート率」などの具体的な指標を設定します。KPIは測定可能で、目標達成度が明確に判断できるものを選ぶ必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

課題の構造化

複雑なビジネス課題を具体的な分析課題に分解する方法

KPIの設定

課題解決の成果を測定するための適切なKPIの設定方法

ステークホルダー分析

関係者の特定と要求事項の整理、合意形成のプロセス

ステークホルダー分析では、プロジェクトに関係する全ての人の要求と期待を整理します。経営層、現場担当者、システム部門、顧客など、それぞれの立場からの要求を把握し、優先順位をつけながら、合意形成を図っていきます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

現状分析

現状の業務プロセスやデータの流れを把握し問題点を特定する方法

仮説の構築

データ分析によって検証すべき仮説の立て方と優先順位付け

目標値の設定

定量的な目標値の設定方法達成基準の明確化

現状分析では、既存の業務プロセスやデータの流れを詳細に調査します。どのようなデータがどこで発生し、どのように処理され、どこで活用されているのか、その全体像を把握します。この過程で、非効率な部分や改善すべき点が明らかになります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

現状分析

現状の業務プロセスやデータの流れを把握し問題点を特定する方法

仮説の構築

データ分析によって検証すべき仮説の立て方と優先順位付け

目標値の設定

定量的な目標値の設定方法達成基準の明確化

仮説の構築は、データ分析の方向性を定める重要なステップです。業務知識や過去の経験を基に、「なぜその問題が起きているのか」「どうすれば改善できるのか」という仮説を立てます。複数の仮説を整理し、検証の優先順位を決定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

現状分析

現状の業務プロセスやデータの流れを把握し問題点を特定する方法

仮説の構築

データ分析によって検証すべき仮説の立て方と優先順位付け

目標値の設定

定量的な目標値の設定方法達成基準の明確化

目標値の設定では、具体的な数値目標を定めます。例えば、「顧客満足度を3ヶ月以内に10%向上させる」といった形で、期間と達成レベルを明確にします。目標値は、意欲的でありながらも現実的に達成可能な水準に設定することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

ビジネス課題の明確化

分析設計の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

分析設計の実践

分析手法の選定

課題解決に適した分析手法の選定基準と決定プロセス

必要データの特定

分析に必要なデータ項目の洗い出しデータソースの検討方法

分析スケジュール

データ収集から結果検証までの具体的なスケジュール策定

分析手法の選定では、課題の性質と得られるデータの特性を考慮します。例えば、顧客の離反予測には機械学習の分類モデル、需要予測には時系列分析、顧客セグメンテーションにはクラスタリング分析というように、目的に応じた適切な手法を選びます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

分析設計の実践

分析手法の選定

課題解決に適した分析手法の選定基準と決定プロセス

必要データの特定

分析に必要なデータ項目の洗い出しデータソースの検討方法

分析スケジュール

データ収集から結果検証までの具体的なスケジュール策定

必要データの特定では、選定した分析手法に基づいて、具体的に必要なデータ項目を洗い出します。社内システムのデータ、外部データ、新規に収集が必要なデータなど、データソースごとに利用可能性を評価します。データの粒度や期間なども、この段階で検討します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

分析設計の実践

分析手法の選定

課題解決に適した分析手法の選定基準と決定プロセス

必要データの特定

分析に必要なデータ項目の洗い出しデータソースの検討方法

分析スケジュール

データ収集から結果検証までの具体的なスケジュール策定

分析スケジュールでは、データ収集から結果検証までの工程を時系列で整理します。データの入手にかかる時間、前処理の工数、分析作業の期間、結果の検証時間など、各工程に必要な時間を見積もり、具体的なスケジュールを策定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

分析設計の実践

リスク分析

想定されるリスクの特定対応策の検討方法

実行計画の策定

具体的な分析作業の計画立案実行体制の整備

評価指標の設計

分析結果を評価するための指標と基準の設定方法

リスク分析では、プロジェクト遂行上の潜在的な問題を特定します。データの品質リスク、スケジュールリスク、技術的なリスク、組織的なリスクなど、様々な観点からリスクを洗い出し、それぞれに対する対応策を検討します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

分析設計の実践

リスク分析

想定されるリスクの特定対応策の検討方法

実行計画の策定

具体的な分析作業の計画立案実行体制の整備

評価指標の設計

分析結果を評価するための指標と基準の設定方法

実行計画の策定では、具体的な作業計画を立案します。担当者の割り当て、必要なツールの準備、環境の整備、進捗管理の方法など、実務的な側面まで含めて検討します。特に、チーム内での役割分担と責任範囲を明確にすることが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-2. ビジネス課題の明確化と分析設計

第2章 データ分析基礎

分析設計の実践

リスク分析

想定されるリスクの特定対応策の検討方法

実行計画の策定

具体的な分析作業の計画立案実行体制の整備

評価指標の設計

分析結果を評価するための指標と基準の設定方法

評価指標の設計では、分析結果の妥当性を判断する基準を定めます。統計的な評価指標（精度、再現率など）に加えて、ビジネス上の評価指標（コスト削減額、収益改善額など）も設定します。これらの指標に基づいて、分析の成否を判断します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

データ尺度の活用



ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

データ尺度の種類

名義尺度、順序尺度、間隔尺度、比率尺度の4つの基本尺度

名義尺度の特徴

カテゴリを表すデータの特徴適切な分析手法

順序尺度の特徴

順序関係を持つデータの特徴適切な分析手法

データ尺度は、大きく4つの種類に分類されます。最も基本的な「名義尺度、順序関係」を持つ「順序尺度、等間隔性」を持つ「間隔尺度」、そして絶対的なゼロ点を持つ「比率尺度」です。これらの尺度は、データの性質と可能な演算の範囲を規定する重要な概念です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

データ尺度の種類

名義尺度、順序尺度、間隔尺度、比率尺度の4つの基本尺度

名義尺度の特徴

カテゴリを表すデータの特徴適切な分析手法

順序尺度の特徴

順序関係を持つデータの特徴適切な分析手法

名義尺度は、データを分類するためのラベルとして使用されます。例えば、性別、血液型、商品カテゴリなどが該当します。これらのデータでは、同じか異なるかの判断のみが可能で、大小や加減算の概念は存在しません。

ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

データ尺度の種類

名義尺度、順序尺度、間隔尺度、比率尺度の4つの基本尺度

名義尺度の特徴

カテゴリを表すデータの特徴適切な分析手法

順序尺度の特徴

順序関係を持つデータの特徴適切な分析手法

順序尺度は、データ間の順序関係を表現できる尺度です。例えば、満足度（非常に満足・やや満足・やや不満・非常に不満）や、学校の成績（A・B・C・D）などが該当します。大小関係は判断できますが、その差の大きさを数値的に扱うことはできません。

ビッグデータ基礎 1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

間隔尺度の特徴

間隔の等間隔性を持つデータの特徴適切な分析手法

比率尺度の特徴

絶対的なゼロ点を持つデータの特徴適切な分析手法

尺度の制約

各尺度で可能な演算と統計処理の制約

間隔尺度は、データ間の差が等間隔である尺度です。温度（摂氏・華氏）や知能指数（IQ）などが該当します。データ間の差の計算は可能ですが、比の計算は意味を持ちません。例えば、20°Cと10°Cの差は10°Cですが、20°Cが10°Cの2倍という解釈は適切ではありません。

ビッグデータ基礎 1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

間隔尺度の特徴

間隔の等間隔性を持つデータの特徴適切な分析手法

比率尺度の特徴

絶対的なゼロ点を持つデータの特徴適切な分析手法

尺度の制約

各尺度で可能な演算と統計処理の制約

比率尺度は、絶対的なゼロ点を持つ尺度です。身長、体重、売上金額、年齢などが該当します。この尺度では、加減乗除のすべての演算が可能で、「AはBの2倍」といった比の解釈も意味を持ちます。

ビッグデータ基礎 1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

間隔尺度の特徴

間隔の等間隔性を持つデータの特徴適切な分析手法

比率尺度の特徴

絶対的なゼロ点を持つデータの特徴適切な分析手法

尺度の制約

各尺度で可能な演算と統計処理の制約

各尺度には、可能な統計処理に制約があります。例えば、名義尺度では最頻値しか求められませんが、比率尺度ではあらゆる統計量の計算が可能です。この制約を理解し、適切な統計処理を選択することが重要です。

NEXT

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の基本概念

データ尺度の活用



2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の活用

尺度の選択

分析目的に応じた適切なデータ尺度の選択方法

尺度変換

異なる尺度間のデータ変換方法と注意点

適切な統計手法

データ尺度に応じた適切な統計分析手法の選択

尺度の選択は、分析の目的に応じて慎重に行う必要があります。例えば、顧客満足度を測定する場合、単純な2値（満足・不満）の名義尺度よりも、5段階評価などの順序尺度の方が、より詳細な分析が可能になります。また、正確な数値測定が必要な場合は、間隔尺度や比率尺度を選択します。

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の活用

尺度の選択

分析目的に応じた適切なデータ尺度の選択方法

尺度変換

異なる尺度間のデータ変換方法と注意点

適切な統計手法

データ尺度に応じた適切な統計分析手法の選択

尺度変換は、データの性質を変えず、より適切な形式に変換する作業です。例えば、年齢（比率尺度）を年代（順序尺度）に変換したり、連続的な満足度スコア（間隔尺度）を満足・不満の2値（名義尺度）に変換したりします。ただし、より詳細な尺度から粗い尺度への変換は可能ですが、その逆は一般的にできません。

ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の活用

尺度の選択

分析目的に応じた適切なデータ尺度の選択方法

尺度変換

異なる尺度間のデータ変換方法と注意点

適切な統計手法

データ尺度に応じた適切な統計分析手法の選択

統計手法の選択では、データの尺度を慎重に考慮します。名義尺度データにはカイ二乗検定、順序尺度データにはノンパラメトリック検定、間隔・比率尺度データにはt検定や分散分析といった、それぞれの尺度に適した手法を選択します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の活用

可視化方法

各尺度のデータに適した可視化手法の選択

データ収集設計

目的に応じた尺度を考慮したデータ収集方法の設計

品質管理

データ尺度を考慮した品質管理の方法

可視化方法も、データの尺度に応じて適切に選択する必要があります。名義尺度データには円グラフや棒グラフ、順序尺度データには積み上げ棒グラフ、間隔・比率尺度データにはヒストグラムや散布図といった、データの特性を効果的に表現できる図を選びます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の活用

可視化方法

各尺度のデータに適した可視化手法の選択

データ収集設計

目的に応じた尺度を考慮したデータ収集方法の設計

品質管理

データ尺度を考慮した品質管理の方法

データ収集の設計では、将来の分析ニーズも考えます。例えば、年齢データを収集する際、最初から年代区分（順序尺度）で収集するのではなく、実年齢（比率尺度）で収集しておけば、後で必要に応じて様々な区分に変換できます。

2-3. データの種類と尺度

第2章 データ分析基礎

データ尺度の活用

可視化方法

各尺度のデータに適した可視化手法の選択

データ収集設計

目的に応じた尺度を考慮したデータ収集方法の設計

品質管理

データ尺度を考慮した品質管理の方法

品質管理では、各尺度の特性に応じたチェック方法を実装します。例えば、名義尺度データでは定義された値以外が含まれていないか、比率尺度データでは負の値や異常に大きい値がないかなど、尺度の特性に応じた適切なチェックを行います。

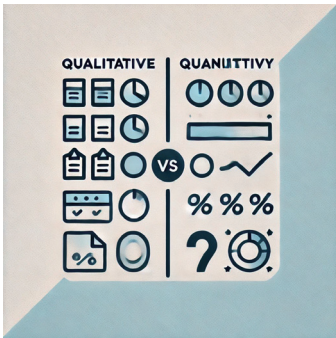
NEXT

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

量的データの理解



2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

質的データの定義

カテゴリカルデータの特徴と種類、データ収集方法

名義データ

性別や血液型分類を表す質的データの特徴と扱い方

順序データ

満足度や学歴順序を持つ質的データの特徴と扱い方

質的データとは、数値では表現できない特性や属性を表すカテゴリ化可能なデータのことで、例えば、性別、職業、商品カテゴリ、好みなどが該当します。これらのデータは、分類や順序づけによって情報を表現します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

質的データの定義

カテゴリカルデータの特徴と種類、データ収集方法

名義データ

性別や血液型分類を表す質的データの特徴と扱い方

順序データ

満足度や学歴順序を持つ質的データの特徴と扱い方

名義データは、単純な分類を表す質的データです。例えば、性別、血液型、婚姻状態などが該当します。これらのデータでは、カテゴリ間に大小関係はなく、ただ「異なるという区別のみが意味を持ちます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

質的データの定義

カテゴリカルデータの特徴と種類、データ収集方法

名義データ

性別や血液型分類を表す質的データの特徴と扱い方

順序データ

満足度や学歴順序を持つ質的データの特徴と扱い方

順序データは、カテゴリ間に順序関係のある質的データです。例えば、満足度（非常に満足・やや満足・普通・やや不満・非常に不満）や、学歴（小学校・中学校・高校・大学）などが該当します。カテゴリ間の順序は明確ですが、その間隔は必ずしも等しくありません。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

質的データの収集

アンケートやインタビューによる質的データの収集方法

コーディング

質的データのカテゴリー化とコード化の方法

質的データの記述

質的データの要約と記述統計の方法

質的データの収集では、主にアンケートやインタビューが用いられます。アンケートでは選択肢の設計が重要で、カテゴリに漏れや重複がないように注意が必要です。インタビューでは、回答を適切にカテゴリ化する技術が求められます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

質的データの収集

アンケートやインタビューによる質的データの収集方法

コーディング

質的データのカテゴリー化とコード化の方法

質的データの記述

質的データの要約と記述統計の方法

コーディングは、質的データを分析可能な形に変換するプロセスです。例えば、職業を「会社員」「公務員」「自営業」などのカテゴリに分類し、それぞれに数値コードを割り当てます。このコード化により、データの整理や分析が容易になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

質的データの収集

アンケートやインタビューによる質的データの収集方法

コーディング

質的データのカテゴリー化とコード化の方法

質的データの記述

質的データの要約と記述統計の方法

質的データの記述では、頻度や比率が主な指標となります。例えば、各カテゴリの出現頻度、構成比率、最頻値などを求めます。これらの指標を用いて、データの分布や特徴を把握します。

ビッグデータ基礎

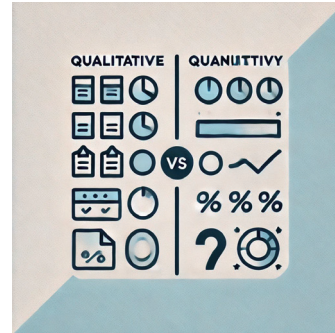
1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

質的データの理解

量的データの理解



ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

量的データの理解

量的データの定義

数値で表されるデータの特徴と種類、データ収集方法

離散データ

個数や回数とびとびの値を取るデータの特徴と扱い方

連続データ

身長や重量連続的な値を取るデータの特徴と扱い方

量的データとは、数値で表現される測定可能なデータのことです。例えば、身長、体重、年齢、売上金額、温度などが該当します。これらのデータは、大きさの比較や四則演算が可能という特徴を持ちます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

量的データの理解

量的データの定義

数値で表されるデータの特徴と種類、データ収集方法

離散データ

個数や回数とびとびの値を取るデータの特徴と扱い方

連続データ

身長や重量連続的な値を取るデータの特徴と扱い方

離散データは、とびとびの値のみを取るデータです。例えば、家族の人数、商品の個数、来店回数などが該当します。これらのデータは、必ず整数値を取り、その間の値は存在しません。例えば、家族の人数が2.5人ということはありません。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

量的データの理解

量的データの定義

数値で表されるデータの特徴と種類、データ収集方法

離散データ

個数や回数とびとびの値を取るデータの特徴と扱い方

連続データ

身長や重量連続的な値を取るデータの特徴と扱い方

連続データは、ある範囲内のどの値でも取り得るデータです。例えば、身長、体重、時間、距離などが該当します。これらのデータは、測定精度の範囲内で任意の値を取ることができます。例えば、身長は171.234cmというように、理論的には無限の精度で表現できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

量的データの理解

測定と精度

量的データの測定方法と測定精度の考え方

量的データの変換

スケーリングや標準化データ変換の方法

データ型の選択

分析目的に応じた適切なデータ型の選択方法

測定と精度は、量的データを扱う上で重要な概念です。測定には必ず誤差が伴い、その精度は測定機器や測定方法に依存します。例えば、体重計の精度が100g単位なら、それ以上の細かい値は意味を持ちません。適切な精度での測定と記録が重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

量的データの理解

測定と精度

量的データの測定方法と測定精度の考え方

量的データの変換

スケーリングや標準化データ変換の方法

データ型の選択

分析目的に応じた適切なデータ型の選択方法

量的データの変換では、スケーリングや標準化などの処理を行います。例えば、異なる単位のデータを比較可能にするための標準化や、分布の歪みを補正するための対数変換などが用いられます。これらの変換により、より適切な分析が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-4. 質的データと量的データの違い

第2章 データ分析基礎

量的データの理解

測定と精度

量的データの測定方法と測定精度の考え方

量的データの変換

スケールや標準化データ変換の方法

データ型の選択

分析目的に応じた適切なデータ型の選択方法

データ型の選択では、分析の目的や必要な精度を考慮します。例えば、年齢データを扱う場合、実年齢（連続データ）で収集するか、年代区分（離散データ）で収集するかを、分析目的に応じて選択します。また、保存や処理の効率性も考慮に入れる必要があります。

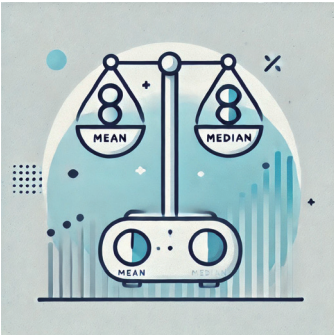
NEXT

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

中央値の理解



2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

算術平均の概念

データの合計を個数で割る算術平均の基本的な考え方と計算方法

加重平均の活用

重要度や規模を考慮した加重平均の計算方法と適用場面

幾何平均の特徴

成長率や収益率比率データに適した幾何平均

算術平均は、最も一般的な平均値の計算方法です。データの合計をデータ数で割ることで求められます。例えば、5人の試験点数が80点、75点、90点、85点、70点の場合、合計400点を5で割って80点が算術平均となります。日常的に使用される「平均」とは、通常この算術平均を指します。

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

算術平均の概念

データの合計を個数で割る算術平均の基本的な考え方と計算方法

加重平均の活用

重要度や規模を考慮した加重平均の計算方法と適用場面

幾何平均の特徴

成長率や収益率比率データに適した幾何平均

加重平均は、データの重要度や規模を考慮した平均値です。例えば、科目別の成績評価で、試験の配点が異なる場合に使用します。中間試験30%、期末試験50%、課題20%というように、重みをつけて計算します。また、企業の平均給与を、従業員数で重みづけて計算する場合なども、加重平均が適しています。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

算術平均の概念

データの合計を個数で割る算術平均の基本的な考え方と計算方法

加重平均の活用

重要度や規模を考慮した加重平均の計算方法と適用場面

幾何平均の特徴

成長率や収益率比率データに適した幾何平均

幾何平均は、比率や変化率を扱う際に適した平均値です。例えば、3年間の売上成長率が10%、15%、5%の場合、これらの平均的な成長率を求めるのに使用します。算術平均では正確な値が得られず、幾何平均を使用する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

調和平均の応用

速度や単価特定の場面で使用される調和平均

平均値の特性

平均値の数学的性質データ分析における利点と欠点

外れ値の影響

平均値が外れ値に影響される性質その対処方法

調和平均は、速度や単価などの「率の逆数」を扱う際に使用します。例えば、往路50km/hと復路30km/hで行った場合の平均速度は、算術平均の40km/hではなく、調和平均の37.5km/hが正しい値となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

調和平均の応用

速度や単価特定の場面で使用される調和平均

平均値の特性

平均値の数学的性質データ分析における利点と欠点

外れ値の影響

平均値が外れ値に影響される性質その対処方法

平均値の特性として、すべてのデータを一つの値で代表できる点が利点です。また、数学的な処理が容易で、様々な統計分析に利用できます。しかし、データの分布の形状を考慮していないため、偏りのあるデータでは実態を正確に表現できないこともあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

調和平均の応用

速度や単価特定の場面で使用される調和平均

平均値の特性

平均値の数学的性質データ分析における利点と欠点

外れ値の影響

平均値が外れ値に影響される性質その対処方法

外れ値の影響も重要な考慮点です。例えば、月収データで一人だけ極端に高額な値がある場合、平均値は大きく上方に引っ張られてしまいます。このような場合、外れ値を除外して計算するか、中央値など他の代表値の使用を検討する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

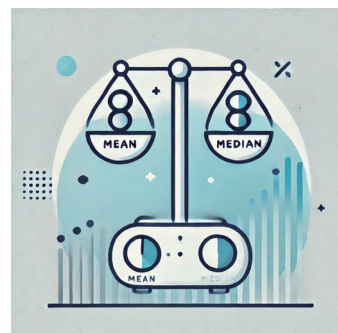


2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

平均値の基本

中央値の理解



ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

中央値の理解

中央値の定義

データを順序づけた際の中央に位置する値の特徴と求め方

中央値の特性

外れ値の影響を受けにくい中央値の性質と活用場面

平均値との比較

平均値と中央値の違いそれぞれの適用場面

中央値は、データを小さい順に並べた時の中央に位置する値です。データ数が奇数の場合は中央の値そのもの、偶数の場合は中央に位置する2つの値の平均を取ります。例えば、1、3、4、7、9の5個のデータでは4が中央値、1、3、4、7、9、12の6個のデータでは $(4+7) \div 2 = 5.5$ が中央値となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

中央値の理解

中央値の定義

データを順序づけた際の中央に位置する値の特徴と求め方

中央値の特性

外れ値の影響を受けにくい中央値の性質と活用場面

平均値との比較

平均値と中央値の違いそれぞれの適用場面

中央値の最大の特徴は、外れ値の影響を受けにくい点です。例えば、100、120、150、200、1000というデータで、最後の1000は極端な外れ値ですが、中央値は150のままで変化しません。一方、平均値は314と大きく影響を受けてしまいます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

中央値の理解

中央値の定義

データを順序づけた際の中央に位置する値の特徴と求め方

中央値の特性

外れ値の影響を受けにくい中央値の性質と活用場面

平均値との比較

平均値と中央値の違いそれぞれの適用場面

平均値との比較では、それぞれの特性を理解することが重要です。平均値はすべてのデータを考慮する反面、外れ値の影響を受けやすく、歪んだ分布では実態を表現しにくくなります。中央値は外れ値の影響を受けにくい反面、データの総量的な情報は失われます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

中央値の理解

四分位数との関係

中央値と四分位数の関係性

代表値の選択

データの性質に応じた適切な代表値の選択方法

実務での活用

ビジネス場面における中央値の効果的な活用方法

四分位数との関係も重要です。中央値は、データを4等分する四分位数の第2四分位数（Q2）に相当します。四分位数と合わせて使用することで、データの分布の特徴をより詳細に把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

中央値の理解

四分位数との関係

中央値と四分位数の関係性

代表値の選択

データの性質に応じた適切な代表値の選択方法

実務での活用

ビジネス場面における中央値の効果的な活用方法

代表値の選択では、データの性質と分析の目的を考慮します。例えば、所得のような右に歪んだ分布では中央値が適していますが、テストの平均点のような教育評価では算術平均が一般的です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-5. 基本統計量：平均値と中央値

第2章 データ分析基礎

中央値の理解

四分位数との関係

中央値と四分位数の関係性

代表値の選択

データの性質に応じた適切な代表値の選択方法

実務での活用

ビジネス場面における中央値の効果的な活用方法

実務での活用例として、不動産価格の分析があります。特に住宅価格は高額物件の影響で平均値が高くなりますが、中央値を使用することで、より一般的な価格水準を把握することができます。また、企業の給与分析でも、中央値を用いることで実態をより適切に表現できることがあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

標準偏差の活用



ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

分散の定義

データのばらつきを示す分散の基本的な概念と計算方法

偏差の概念

平均値からのズレを表す偏差の意味と計算方法

分散の性質

分散の数学的性質データ分析における意味

分散は、データが平均値からどれだけばらついているかを示す指標です。具体的には、各データと平均値との差（偏差）を二乗して平均を取った値です。例えば、テストの点数が60点、70点、80点、90点、100点の場合、平均値は80点で、各データの偏差を二乗して平均を取ることで分散を求めることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

分散の定義

データのばらつきを示す分散の基本的な概念と計算方法

偏差の概念

平均値からのズレを表す偏差の意味と計算方法

分散の性質

分散の数学的性質データ分析における意味

偏差は、各データから平均値を引いた値です。例えば、平均値が80点の場合、60点の偏差は-20、90点の偏差は+10となります。偏差の合計は必ずゼロになるという特徴があります。この偏差を二乗して平均を取ることで、データのばらつきの大きさを数値化できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

分散の定義

データのばらつきを示す分散の基本的な概念と計算方法

偏差の概念

平均値からのズレを表す偏差の意味と計算方法

分散の性質

分散の数学的性質データ分析における意味

分散の性質として、すべての偏差を二乗するため、必ず正の値となります。また、データの単位を二乗した単位となるため、元のデータとは単位が異なります。例えば、身長分散は cm^2 、体重分散は kg^2 となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

不偏分散

母集団の分散を推定する際の不偏分散

共分散

2つの変数間の関係性を示す共分散

分散の活用

データ分析における分散の活用方法と解釈

不偏分散は、標本から母集団の分散を推定する際に使用します。標本の分散は母集団の分散を過小評価する傾向があるため、データ数から1を引いた値で割ることで、より正確な推定値を得ることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

不偏分散

母集団の分散を推定する際の不偏分散

共分散

2つの変数間の関係性を示す共分散

分散の活用

データ分析における分散の活用方法と解釈

共分散は、2つの変数の関係性を示す指標です。2つの変数それぞれの偏差の積の平均として計算されます。正の値は両変数が同じ方向に変動する傾向を、負の値は逆方向に変動する傾向を示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

不偏分散

母集団の分散を推定する際の不偏分散

共分散

2つの変数間の関係性を示す共分散

分散の活用

データ分析における分散の活用方法と解釈

分散の活用例として、品質管理における製品のばらつきの評価、投資におけるリスク評価、実験データのばらつきの分析などがあります。大きな分散は、データのばらつきが大きいことを示し、管理や予測が難しい状況を示唆します。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

分散の理解

標準偏差の活用



ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

標準偏差の活用

標準偏差の定義

分散の平方根として定義される標準偏差の概念と計算方法

標準化の方法

標準偏差を用いたデータの標準化プロセス

正規分布との関係

標準偏差と正規分布の関係性

標準偏差は、データのばらつきを元のデータと同じ単位で表現できる指標です。分散の平方根を取ることで、例えば身長であればcm、体重であればkgという元の単位でばらつきを理解することができます。このため、分散よりも直感的に理解しやすい特徴があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

標準偏差の活用

標準偏差の定義

分散の平方根として定義される標準偏差の概念と計算方法

標準化の方法

標準偏差を用いたデータの標準化プロセス

正規分布との関係

標準偏差と正規分布の関係性

標準化は、異なるデータを比較可能な形に変換するプロセスです。各データから平均値を引き、標準偏差で割ることで、平均0、標準偏差1の標準化変数に変換できます。これにより、単位やスケールの異なるデータを公平に比較することが可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

標準偏差の活用

標準偏差の定義

分散の平方根として定義される標準偏差の概念と計算方法

標準化の方法

標準偏差を用いたデータの標準化プロセス

正規分布との関係

標準偏差と正規分布の関係性

正規分布との関係も重要です。正規分布では、平均値から標準偏差の ± 1 倍の範囲に約68%のデータが、 ± 2 倍の範囲に約95%のデータが、 ± 3 倍の範囲に約99.7%のデータが含まれます。この性質は、データの分布を理解する上で重要な指標となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

標準偏差の活用

信頼区間

標準偏差を用いた信頼区間の算出方法

変動係数

平均値の異なるデータを比較する際の変動係数

実務での解釈

ビジネスデータ分析における標準偏差の解釈方法

信頼区間の算出では、標準偏差を用いて推定値の精度を評価します。例えば、平均値の95%信頼区間は、標準偏差を用いて計算することができます。これにより、推定値の信頼性を定量的に評価することが可能です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

標準偏差の活用

信頼区間

標準偏差を用いた信頼区間の算出方法

変動係数

平均値の異なるデータを比較する際の変動係数

実務での解釈

ビジネスデータ分析における標準偏差の解釈方法

変動係数は、標準偏差を平均値で割った値で、相対的なばらつきを示します。平均値の異なるデータのばらつきを比較する際に有用です。例えば、年収の高い層と低い層でのばらつきの比較などに使用されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-6. 基本統計量：分散と標準偏差

第2章 データ分析基礎

標準偏差の活用

信頼区間

標準偏差を用いた信頼区間の算出方法

変動係数

平均値の異なるデータを比較する際の変動係数

実務での解釈

ビジネスデータ分析における標準偏差の解釈方法

実務での解釈では、標準偏差を用いて異常値の検出や品質管理を行います。例えば、製造ラインでの製品品質のばらつきを評価したり、サービスの提供時間のばらつきを分析したりする際に活用されます。また、予測モデルの精度評価にも標準偏差が重要な役割を果たします。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

パーセンタイルの活用



ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

最頻値の定義

データセットで最も頻繁に出現する値の特徴と求め方

離散データでの最頻値

カテゴリカルデータや離散的な数値での最頻値の扱い

連続データでの最頻値

ヒストグラムを用いた連続データの最頻値の求め方

最頻値は、データセットの中で最も多く出現する値のことです。例えば、試験の点数が60、70、70、80、90点の場合、70点が2回出現して最も多いため、70点が最頻値となります。この指標は、データの中で最も一般的な値を知るのに役立ちます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

最頻値の定義

データセットで最も頻繁に出現する値の特徴と求め方

離散データでの最頻値

カテゴリカルデータや離散的な数値での最頻値の扱い

連続データでの最頻値

ヒストグラムを用いた連続データの最頻値の求め方

離散データでの最頻値は、比較的簡単に求めることができます。例えば、サイコロを10回振って出た目が1,2,2,3,2,4,5,6,2,3の場合、2が4回出現して最も多いため、2が最頻値となります。カテゴリカルデータでも同様に、最も出現頻度の高いカテゴリを最頻値として特定できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

最頻値の定義

データセットで最も頻繁に出現する値の特徴と求め方

離散データでの最頻値

カテゴリカルデータや離散的な数値での最頻値の扱い

連続データでの最頻値

ヒストグラムを用いた連続データの最頻値の求め方

連続データでの最頻値は、通常ヒストグラムを使って求めます。データを適切な区間に分割し、最も度数の高い区間の中心値を最頻値とします。ただし、区間の取り方によって結果が変わる可能性があるため、注意が必要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

複数の最頻値

データが複数の最頻値を持つ場合の解釈と扱い方

代表値としての特徴

最頻値の利点と限界、適切な使用場面

分布の形状理解

最頻値を通じたデータ分布の特徴把握

複数の最頻値が存在する場合があります。例えば、1,2,2,3,3,4のデータでは、2と3がそれぞれ2回ずつ出現するため、両方が最頻値となります。このような分布は双峰性と呼ばれ、データに異なる特性を持つ集団が混在している可能性を示唆します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

複数の最頻値

データが複数の最頻値を持つ場合の解釈と扱い方

代表値としての特徴

最頻値の利点と限界、適切な使用場面

分布の形状理解

最頻値を通じたデータ分布の特徴把握

代表値としての最頻値は、特に質的データや離散的なデータの分析に適しています。例えば、アンケートの回答で最も多かった選択肢や、店舗で最も売れている商品を特定する際に有用です。ただし、すべてのデータが異なる値を持つ場合、最頻値は意味をなさないという限界もあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

複数の最頻値

データが複数の最頻値を持つ場合の解釈と扱い方

代表値としての特徴

最頻値の利点と限界、適切な使用場面

分布の形状理解

最頻値を通じたデータ分布の特徴把握

分布の形状理解に関しては、最頻値は平均値や中央値と合わせて考察することが重要です。これら3つの代表値の位置関係から、分布の歪みや特徴を理解することができます。

ビッグデータ基礎

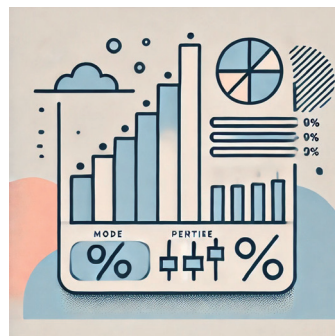
1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

最頻値の基本

パーセンタイルの活用



ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

パーセンタイルの活用

パーセンタイルの概念

データを100等分する際の位置を示すパーセンタイル

四分位数の計算

25、50、75パーセンタイルにあたる四分位数の求め方

十分位数の活用

10等分点による詳細なデータ分布の把握方法

パーセンタイルは、データを小さい順に並べた時の位置を百分率で表したものです。例えば、第25パーセンタイルは、データの下から25%に位置する値を示します。この概念により、データの分布を詳細に把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

パーセンタイルの活用

パーセンタイルの概念

データを100等分する際の位置を示すパーセンタイル

四分位数の計算

25、50、75パーセンタイルにあたる四分位数の求め方

十分位数の活用

10等分点による詳細なデータ分布の把握方法

四分位数は、特に重要なパーセンタイルです。第1四分位数（25パーセンタイル）、第2四分位数（50パーセンタイル、中央値）、第3四分位数（75パーセンタイル）の3つの値により、データを4つの等分に分割します。これにより、データの散らばり具合を把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

パーセンタイルの活用

パーセンタイルの概念

データを100等分する際の位置を示すパーセンタイル

四分位数の計算

25、50、75パーセンタイルにあたる四分位数の求め方

十分位数の活用

10等分点による詳細なデータ分布の把握方法

十分位数は、データを10等分する値です。例えば、第1十分位数は下から10%の位置、第9十分位数は下から90%の位置を示します。これにより、より細かなデータの分布状況を把握することができます。特に、大規模なデータセットの分布を理解する際に有用です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

パーセンタイルの活用

百分位数の解釈

細かな順位付けに用いる百分位数の活用方法

異常値の検出

パーセンタイルを用いた異常値の検出方法

実務での応用

成績評価やベンチマーク実務でのパーセンタイルの活用

百分位数は、最も細かい区分を提供します。例えば、テストの成績で上位3%に入るための得点や、身長の発育曲線で95パーセンタイルの値など、詳細な位置づけが必要な場合に使用されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

パーセンタイルの活用

百分位数の解釈

細かな順位付けに用いる百分位数の活用方法

異常値の検出

パーセンタイルを用いた異常値の検出方法

実務での応用

成績評価やベンチマーク実務でのパーセンタイルの活用

異常値の検出では、パーセンタイルが重要な役割を果たします。例えば、第1四分位数から第3四分位数までの範囲（四分位範囲）の1.5倍を超えて外れている値を異常値として検出する方法がよく使われます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-7. 基本統計量：最頻値とパーセンタイル

第2章 データ分析基礎

パーセンタイルの活用

百分位数の解釈

細かな順位付けに用いる百分位数の活用方法

異常値の検出

パーセンタイルを用いた異常値の検出方法

実務での応用

成績評価やベンチマーク実務でのパーセンタイルの活用

実務での応用例として、学校でのテストの成績評価、企業での給与水準の設定、製品の品質管理における規格値の設定などがあります。例えば、上位10%に入る社員の給与水準を把握したり、製品の寸法が規格の95パーセンタイル内に収まっているかを確認したりする際に活用されます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

分布の理解と解釈



ビッグデータ基礎

1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

ヒストグラムの定義

データの分布を視覚化するヒストグラムの基本概念と作成方法

階級の設定

データを区分する階級の幅と数の適切な決定方法

度数分布表

ヒストグラム作成の基となる度数分布表の作成方法

ヒストグラムは、連続的なデータの分布を視覚化するための棒グラフです。横軸にデータの値の範囲（階級）、縦軸に各階級に含まれるデータの個数（度数）をとり、データの分布状況を表現します。例えば、100人の身長データを5cm間隔で区切り、各区間に何人含まれるかを表示することで、身長の分布を視覚的に理解することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

ヒストグラムの定義

データの分布を視覚化するヒストグラムの基本概念と作成方法

階級の設定

データを区分する階級の幅と数の適切な決定方法

度数分布表

ヒストグラム作成の基となる度数分布表の作成方法

階級の設定は、ヒストグラムの見え方に大きく影響します。一般的な目安として、データ数の平方根程度の階級数が推奨されます。例えば、100個のデータなら10程度の階級数が適切です。また、階級の幅は、データの特性や分析の目的に応じて決定します。狭すぎると細かすぎる変動が見えてしまい、広すぎると重要な特徴を見逃す可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

ヒストグラムの定義

データの分布を視覚化するヒストグラムの基本概念と作成方法

階級の設定

データを区分する階級の幅と数の適切な決定方法

度数分布表

ヒストグラム作成の基となる度数分布表の作成方法

度数分布表は、ヒストグラム作成の基礎となります。各階級の範囲、その階級に含まれるデータの個数（度数）、必要に応じて相対度数や累積度数を表形式で整理します。これにより、データの分布状況を数値的に把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

相対度数

データ数の違いを考慮した相対度数ヒストグラム

累積度数

累積相対度数と累積度数曲線の作成方法

グラフの読み方

ヒストグラムからデータの特徴を読み取る方法

相対度数は、各階級の度数を全データ数で割った値です。これにより、データ数の異なる複数の集団を比較することが可能になります。例えば、クラスの人数が異なる2つのクラスのテスト結果を比較する際に有用です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

相対度数

データ数の違いを考慮した相対度数ヒストグラム

累積度数

累積相対度数と累積度数曲線の作成方法

グラフの読み方

ヒストグラムからデータの特徴を読み取る方法

累積度数は、各階級までの度数の合計を表します。これを相対値で表したものが累積相対度数で、データの値がある値以下である割合を示します。累積相対度数を折れ線グラフで表したものを累積度数曲線と呼び、特定の値以下のデータの割合を視覚的に把握することができます。

ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

相対度数

データ数の違いを考慮した相対度数ヒストグラム

累積度数

累積相対度数と累積度数曲線の作成方法

グラフの読み方

ヒストグラムからデータの特徴を読み取る方法

ヒストグラムからは、データの分布の特徴を読み取ることができます。分布の中心傾向、データのばらつき、分布の形状（対称性や歪み）、外れ値の存在など、データの重要な特徴を視覚的に確認することができます。

ビッグデータ基礎 1 2 3 4 5 6 7

NEXT

2-8. データの分布とヒストグラム

第2章 データ分析基礎

ヒストグラムの基礎

分布の理解と解釈



ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

分布の理解と解釈

分布の形状

正規分布、歪んだ分布、双峰性分布典型的な分布形状

分布の中心

平均値、中央値、最頻値の分布上での位置関係

分布の広がり

標準偏差と分布の広がりとの関係

分布の形状には、いくつかの典型的なパターンがあります。最も基本的なのが釣鐘型の正規分布で、多くの自然現象や社会現象で観察されます。また、右に裾が長い分布（右に歪んだ分布：所得など）や、左に裾が長い分布（左に歪んだ分布）、2つの山を持つ双峰性分布（異なる集団が混在）などがあります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

分布の理解と解釈

分布の形状

正規分布、歪んだ分布、双峰性分布典型的な分布形状

分布の中心

平均値、中央値、最頻値の分布上での位置関係

分布の広がり

標準偏差と分布の広がりとの関係

分布の中心を示す代表値として、平均値、中央値、最頻値があります。正規分布では、これら3つの値は一致しますが、歪んだ分布では異なる値を示します。例えば、右に歪んだ分布では、平均値は中央値より右側（大きい方）に位置します。

ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

分布の理解と解釈

分布の形状

正規分布、歪んだ分布、双峰性分布典型的な分布形状

分布の中心

平均値、中央値、最頻値の分布上での位置関係

分布の広がり

標準偏差と分布の広がりとの関係

分布の広がり、標準偏差によって特徴づけられます。正規分布の場合、平均値から標準偏差の ± 1 倍の範囲に約68%、 ± 2 倍の範囲に約95%のデータが含まれます。このような性質は、データのばらつきを理解する上で重要な指標となります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

分布の理解と解釈

歪度と尖度

分布の歪みと尖りを表す指標

外れ値の確認

分布から外れ値を視覚的に確認する方法

分布の比較

複数のデータセットの分布を比較する方法

歪度は分布の非対称性を、尖度は分布の尖り具合を示す指標です。歪度が正の値なら右に、負の値なら左に歪んでいることを示します。尖度が大きいと、分布は正規分布より尖っており、小さいと扁平な形状となります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

分布の理解と解釈

歪度と尖度

分布の歪みと尖りを表す指標

外れ値の確認

分布から外れ値を視覚的に確認する方法

分布の比較

複数のデータセットの分布を比較する方法

外れ値は、分布の主要な部分から大きく離れた値です。ヒストグラムでは、主要な分布から離れた位置に孤立して現れる小さな山として観察されます。これらの値が本当に異常値なのか、あるいは重要な情報を含んでいるのかを、注意深く検討する必要があります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-8. データの分布とヒストグラム

第2章 データ分析基礎

分布の理解と解釈

歪度と尖度

分布の歪みと尖りを表す指標

外れ値の確認

分布から外れ値を視覚的に確認する方法

分布の比較

複数のデータセットの分布を比較する方法

分布の比較では、複数のデータセットの分布を重ねて表示することで、その違いを視覚的に把握することができます。例えば、男女の身長分布を重ねて表示することで、平均値の違いやばらつきの違いを直感的に理解することができます。

ビッグデータ基礎 1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

外れ値の検出と処理



ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

箱ひげ図の要素

四分位数、中央値、外れ値を示す箱ひげ図の基本構造

四分位範囲

データの散らばりを示すIQR（四分位範囲）の計算方法

ひげの範囲

データの広がりを示すひげの設定方法と意味

箱ひげ図は、データの分布を5つの要素で表現します。箱の下端が第1四分位数（Q1）、上端が第3四分位数（Q3）、箱の中の線が中央値を示します。箱から上下に伸びる「ひげ」は、外れ値を除いたデータの範囲を表し、点で表示される値が外れ値となります。これらの要素により、データの分布状況を簡潔に把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

箱ひげ図の要素

四分位数、中央値、外れ値を示す箱ひげ図の基本構造

四分位範囲

データの散らばりを示すIQR（四分位範囲）の計算方法

ひげの範囲

データの広がりを示すひげの設定方法と意味

四分位範囲（IQR）は、箱の長さに相当し、第3四分位数から第1四分位数を引いた値です。例えば、Q1が10、Q3が30の場合、IQRは20となります。この値は、データの中心的な50%がどの程度散らばっているかを示す指標となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

箱ひげ図の要素

四分位数、中央値、外れ値を示す箱ひげ図の基本構造

四分位範囲

データの散らばりを示すIQR（四分位範囲）の計算方法

ひげの範囲

データの広がりを示すひげの設定方法と意味

ひげの範囲は、通常、 $Q1-1.5 \times IQR$ から $Q3+1.5 \times IQR$ までの範囲内にあるデータの最小値と最大値まで伸びます。この範囲を超えるデータは外れ値として点で表示されます。これにより、データの全体的な広がりや異常値の存在を同時に確認できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

箱の解釈

箱の長さや位置から読み取れる分布の特徴

対称性の確認

データの分布の対称性を箱ひげ図から判断する方法

複数グループの比較

複数のデータセットを箱ひげ図で比較する方法

箱の解釈では、箱の長さや位置関係から重要な情報が読み取れます。箱が長ければデータのばらつきが大きく、短ければ集中していることを示します。また、中央値の位置が箱の中で偏っていれば、分布の歪みを示唆します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

箱の解釈

箱の長さや位置から読み取れる分布の特徴

対称性の確認

データの分布の対称性を箱ひげ図から判断する方法

複数グループの比較

複数のデータセットを箱ひげ図で比較する方法

対称性の確認は、中央値を中心とした箱とひげの形状から判断できます。中央値が箱の中央にあり、上下のひげの長さが同じような場合は、データが対称的に分布していることを示します。一方、どちらかに偏りがある場合は、分布の歪みを示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

箱の解釈

箱の長さや位置から読み取れる分布の特徴

対称性の確認

データの分布の対称性を箱ひげ図から判断する方法

複数グループの比較

複数のデータセットを箱ひげ図で比較する方法

複数グループの比較は、箱ひげ図を並べて表示することで容易に行えます。例えば、男女の身長データや、部門別の売上データなど、グループ間の水準やばらつきの違いを視覚的に比較することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

箱ひげ図の構造

外れ値の検出と処理



ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

外れ値の検出と処理

外れ値の定義

箱ひげ図における外れ値の判定基準

外れ値の種類

マイルドな外れ値と極端な外れ値の区別

検出方法

IQRを用いた外れ値の具体的な検出方法

外れ値は、箱ひげ図では一般的に、 $Q1 - 1.5 \times IQR$ より小さい値、または $Q3 + 1.5 \times IQR$ より大きい値として定義されます。この基準は、データが正規分布に従う場合、約99.3%のデータがこの範囲に含まれるに基づいています。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

外れ値の検出と処理

外れ値の定義

箱ひげ図における外れ値の判定基準

外れ値の種類

マイルドな外れ値と極端な外れ値の区別

検出方法

IQRを用いた外れ値の具体的な検出方法

外れ値は、その程度によって2種類に分類されます。1.5×IQRを超えて3×IQRまでの範囲にある値は「マイルドな外れ値」、3×IQRを超える値は「極端な外れ値」として区別されます。箱ひげ図では、これらを異なる記号で表示することもあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

外れ値の検出と処理

外れ値の定義

箱ひげ図における外れ値の判定基準

外れ値の種類

マイルドな外れ値と極端な外れ値の区別

検出方法

IQRを用いた外れ値の具体的な検出方法

検出方法は、具体的な数値計算により行います。例えば、Q1が10、Q3が30、IQRが20の場合、下側の境界は $10 - 1.5 \times 20 = -20$ 、上側の境界は $30 + 1.5 \times 20 = 60$ となり、この範囲を超える値が外れ値として検出されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

外れ値の検出と処理

外れ値の評価

検出された外れ値の正当性を評価する方法

処理方法の選択

外れ値に対する適切な処理方法の選択

分析への影響

外れ値が統計分析に与える影響とその対処法

外れ値の評価では、単に統計的な基準だけでなく、データの性質や発生状況を考慮する必要があります。測定ミスや入力ミスによる異常値なのか、それとも重要な情報を含む特異値なのかを、慎重に判断する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

外れ値の検出と処理
外れ値の評価 検出された外れ値の正当性を評価する方法
処理方法の選択 外れ値に対する適切な処理方法の選択
分析への影響 外れ値が統計分析に与える影響とその対処法

処理方法の選択には、いくつかのオプションがあります。明らかな誤りであれば削除、正当なデータであれば保持、分析の目的によっては平均値や中央値で置換するなど、適切な方法を選択します。特に、外れ値が分析結果に大きな影響を与える場合は、慎重な判断が必要です。

2-9. 箱ひげ図による外れ値の検出

第2章 データ分析基礎

外れ値の検出と処理
外れ値の評価 検出された外れ値の正当性を評価する方法
処理方法の選択 外れ値に対する適切な処理方法の選択
分析への影響 外れ値が統計分析に与える影響とその対処法

分析への影響として、外れ値は特に平均値や標準偏差に大きく影響します。このため、外れ値の存在が疑われる場合は、中央値や四分位数など、外れ値の影響を受けにくい指標の使用を検討する必要があります。また、分析手法によっては、外れ値に対してロバストな手法を選択することも重要です。

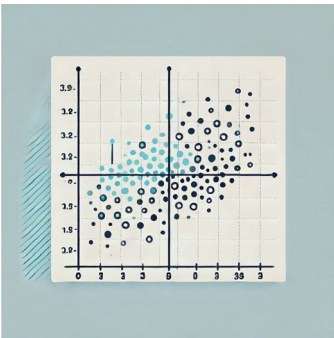
NEXT

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

散布図の解釈



2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

散布図の定義

2つの変数の関係を平面上の点で表現する散布図の基本概念

作図の方法

軸の設定、スケールの選択、プロットの方法散布図の作成手順

変数の選択

散布図で表示する変数の組み合わせの選び方

散布図は、2つの変数の関係を平面上の点で表現するグラフです。横軸に一方の変数、縦軸にもう一方の変数を取り、各データを点としてプロットします。例えば、身長と体重の関係、広告費と売上高の関係など、2つの数値データの関連性を視覚的に表現することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

散布図の定義

2つの変数の関係を平面上の点で表現する散布図の基本概念

作図の方法

軸の設定、スケールの選択、プロットの方法散布図の作成手順

変数の選択

散布図で表示する変数の組み合わせの選び方

作図の方法では、まず適切な軸の設定が重要です。データの最小値と最大値を考慮して軸の範囲を決め、読み取りやすい目盛りの間隔を設定します。また、変数の性質に応じて、通常の線形スケールか対数スケールかを選択します。プロットする点の大きさや形状も、見やすさを考慮して決定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

散布図の定義

2つの変数の関係を平面上の点で表現する散布図の基本概念

作図の方法

軸の設定、スケールの選択、プロットの方法散布図の作成手順

変数の選択

散布図で表示する変数の組み合わせの選び方

変数の選択では、分析の目的に応じて適切な組み合わせを選びます。例えば、因果関係を疑う変数の組み合わせや、相互に影響し合うと考えられる変数の組み合わせなどです。また、同じ単位や似たような尺度の変数を組み合わせると、解釈がしやすくなります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

点の分布パターン

直線的、曲線的、ランダム典型的な分布パターン

密度の表現

データ点が重なる場合の表示方法と密度の表現方法

グラフの装飾

タイトル、軸ラベル、凡例効果的な表示要素

点の分布パターンには、いくつかの典型的なものがあります。正の相関を示す右上がりの分布、負の相関を示す右下がりの分布、曲線的な関係を示す分布、特定のパターンを持たないランダムな分布などです。これらのパターンから、変数間の関係性の特徴を読み取ることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

点の分布パターン

直線的、曲線的、ランダム典型的な分布パターン

密度の表現

データ点が重なる場合の表示方法と密度の表現方法

グラフの装飾

タイトル、軸ラベル、凡例効果的な表示要素

密度の表現は、多くのデータ点が重なる場合に特に重要です。透明度を調整したり、密度を色の濃さで表現したり、バブルチャートのように点の大きさで頻度を表現したりする方法があります。これにより、データの集中度合いを視覚的に把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

点の分布パターン

直線的、曲線的、ランダム典型的な分布パターン

密度の表現

データ点が重なる場合の表示方法と密度の表現方法

グラフの装飾

タイトル、軸ラベル、凡例効果的な表示要素

グラフの装飾は、散布図の理解しやすさを高めます。適切なタイトル、軸ラベル（単位を含む）、必要に応じて凡例やグリッド線を追加します。また、重要なデータ点にラベルを付けたり、特定の領域に注釈を加えたりすることで、より効果的な表現が可能になります。

ビッグデータ基礎

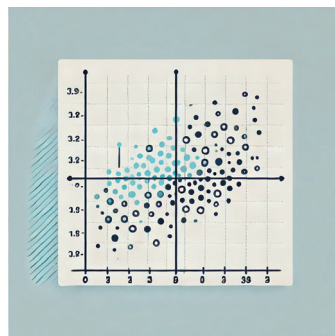
1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の基礎

散布図の解釈



ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の解釈

傾向の読み取り

散布図から2変数間の関係性を読み取る方法

相関の強さ

点の散らばり具合から相関の強さを判断する方法

因果関係の検討

相関と因果関係の違いを理解し、解釈する方法

傾向の読み取りでは、点の分布の方向性や形状から、2変数の関係性を把握します。右上がりの分布は正の相関（一方が増加すると他方も増加）、右下がりの分布は負の相関（一方が増加すると他方は減少）を示します。また、直線的か曲線的かといった関係の形状も重要な情報となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の解釈

傾向の読み取り

散布図から2変数間の関係性を読み取る方法

相関の強さ

点の散らばり具合から相関の強さを判断する方法

因果関係の検討

相関と因果関係の違いを理解し、解釈する方法

相関の強さは、点がどの程度直線的なパターンに近づいているかで判断します。点が直線上に近く集まっているほど相関が強く、ばらついているほど相関が弱いことを示します。ただし、相関の強さは目視による主観的な判断だけでなく、相関係数などの数値指標と合わせて評価することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の解釈

傾向の読み取り

散布図から2変数間の関係性を読み取る方法

相関の強さ

点の散らばり具合から相関の強さを判断する方法

因果関係の検討

相関と因果関係の違いを理解し、解釈する方法

因果関係の検討では、注意が必要です。散布図で相関関係が見られても、必ずしも因果関係があるとは限りません。例えば、アイスクリームの売上と熱中症の発生件数には正の相関がありますが、これは気温という第三の要因の影響によるものです。

ビッグデータ基礎 1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の解釈

グループの識別

散布図上での異なるグループの識別方法

外れ値の確認

散布図を用いた外れ値の視覚的な確認方法

時系列変化

時間経過による関係性の変化を散布図で表現する方法

グループの識別では、点の色や形を変えることで、異なるグループを区別します。例えば、男女別のデータを異なる色でプロットしたり、年代別にマーカーの形を変えたりすることで、グループごとの特徴を比較することができます。

ビッグデータ基礎 1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の解釈

グループの識別

散布図上での異なるグループの識別方法

外れ値の確認

散布図を用いた外れ値の視覚的な確認方法

時系列変化

時間経過による関係性の変化を散布図で表現する方法

外れ値の確認は、全体的な分布パターンから大きく外れた点を見つけることで行います。これらの点 genuinely 異常値なのか、あるいは重要な情報を含んでいるのかを、データの文脈に基づいて判断する必要があります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-10. 散布図による2変数の関係把握

第2章 データ分析基礎

散布図の解釈

グループの識別

散布図上での異なるグループの識別方法

外れ値の確認

散布図を用いた外れ値の視覚的な確認方法

時系列変化

時間経過による関係性の変化を散布図で表現する方法

時系列変化を表現する場合、点をつなぐ線を追加したり、時間の経過を矢印で示したりします。また、アニメーション機能を使用して、時間経過による変化を動的に表現することもできます。これにより、関係性の時間的な変化を視覚的に理解することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

相関分析の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

相関係数の定義

ピアソンの積率相関係数の数学的定義と意味

計算方法

共分散と標準偏差を用いた相関係数の具体的な計算方法

値の範囲

-1から+1までの値が持つ意味と解釈

相関係数は、2つの変数間の直線的な関係の強さを-1から+1の値で表す統計量です。最もよく使われるのはピアソンの積率相関係数で、2つの変数の共分散を、それぞれの標準偏差の積で割って算出します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

相関係数の定義

ピアソンの積率相関係数の数学的定義と意味

計算方法

共分散と標準偏差を用いた相関係数の具体的な計算方法

値の範囲

-1から+1までの値が持つ意味と解釈

計算方法を具体的に見ていきましょう。まず、2つの変数それぞれについて、各データと平均値との差（偏差）を求めます。次に、2つの変数の偏差の積を求め、その平均を計算して共分散を得ます。最後に、この共分散を各変数の標準偏差の積で割ることで、相関係数が求められます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

相関係数の定義

ピアソンの積率相関係数の数学的定義と意味

計算方法

共分散と標準偏差を用いた相関係数の具体的な計算方法

値の範囲

-1から+1までの値が持つ意味と解釈

相関係数の値は必ず-1から+1の間に収まります。+1は完全な正の相関（一方が増加すると他方も比例して増加）、-1は完全な負の相関（一方が増加すると他方は比例して減少）、0は無相関（直線的な関係がない）を示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

強さの判断

相関係数の値から関係の強さを判断する基準

符号の意味

正の相関と負の相関の違いと解釈

限界と注意点

相関係数の使用における注意点と限界

相関の強さの判断には、一般的な目安があります。絶対値が0.7以上で強い相関、0.4から0.7で中程度の相関、0.2から0.4で弱い相関、0.2未満でほとんど相関なし、とされることが多いです。ただし、この基準は分野や目的によって異なることがあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

強さの判断

相関係数の値から関係の強さを判断する基準

符号の意味

正の相関と負の相関の違いと解釈

限界と注意点

相関係数の使用における注意点と限界

符号の意味は、変数間の関係の方向を示します。正の相関は、一方の変数が増加すると他方も増加する関係を表します。例えば、身長と体重の関係です。負の相関は、一方が増加すると他方が減少する関係を表します。例えば、気温と暖房の使用時間の関係です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

強さの判断

相関係数の値から関係の強さを判断する基準

符号の意味

正の相関と負の相関の違いと解釈

限界と注意点

相関係数の使用における注意点と限界

相関係数を使用する際の注意点もいくつかあります。まず、相関は因果関係を示すものではありません。また、外れ値の影響を受けやすく、非線形な関係は適切に評価できません。さらに、変数の単位や尺度の影響を受けないという利点がある一方で、実際の変化の大きさは反映されないという限界もあります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関係数の基本

相関分析の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関分析の実践

順位相関係数

スピアマンの順位相関係数の特徴と使用場面

偏相関係数

第三の変数の影響を除いた相関関係の分析方法

相関行列

複数変数間の相関関係を一覧表示する方法

順位相関係数は、データを順位に変換して計算する相関係数です。スピアマンの順位相関係数が代表的で、データの分布が正規分布に従わない場合や、外れ値の影響を軽減したい場合に使用されます。例えば、テストの得点と学習時間の関係を分析する際に有効です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関分析の実践

順位相関係数

スピアマンの順位相関係数の特徴と使用場面

偏相関係数

第三の変数の影響を除いた相関関係の分析方法

相関行列

複数変数間の相関関係を一覧表示する方法

偏相関係数は、第三の変数の影響を制御した上での2変数間の相関を測ります。例えば、年齢の影響を除いた上での身長と体重の関係を知りたい場合に使用します。これにより、より純粋な2変数間の関係を把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関分析の実践

順位相関係数

スピアマンの順位相関係数の特徴と使用場面

偏相関係数

第三の変数の影響を除いた相関関係の分析方法

相関行列

複数変数間の相関関係を一覧表示する方法

相関行列は、3つ以上の変数間の相関関係を一覧表示する方法です。行列の各要素に相関係数を配置することで、すべての変数の組み合わせにおける相関関係を一目で把握することができます。ヒートマップなどの可視化手法と組み合わせることで、より直感的な理解が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関分析の実践

検定と評価

相関係数の統計的有意性の検定方法

可視化の方法

相関関係の効果的な可視化手法

実務での解釈

ビジネスコンテキストでの相関係数の解釈方法

相関係数の統計的検定では、計算された相関係数が統計的に有意かどうかを評価します。帰無仮説「真の相関係数は0である」に対して、p値を計算して判断します。ただし、サンプルサイズが大きい場合、弱い相関でも統計的に有意となることに注意が必要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関分析の実践

検定と評価

相関係数の統計的有意性の検定方法

可視化の方法

相関関係の効果的な可視化手法

実務での解釈

ビジネスコンテキストでの相関係数の解釈方法

可視化の方法としては、散布図に回帰直線を追加したり、相関行列をヒートマップで表示したりする方法があります。また、変数間の関係を線の太さや色で表現するネットワーク図なども、複雑な相関関係を理解するのに役立ちます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-11. 相関係数の計算と解釈

第2章 データ分析基礎

相関分析の実践

検定と評価

相関係数の統計的有意性の検定方法

可視化の方法

相関関係の効果的な可視化手法

実務での解釈

ビジネスコンテキストでの相関係数の解釈方法

実務での解釈では、相関係数の統計的な意味だけでなく、ビジネス的な重要性も考慮する必要があります。例えば、マーケティングでは0.3程度の弱い相関でも実務的に重要な示唆を持つ場合があります。また、相関関係から因果関係を安易に推測せず、業務知識や他の分析結果と組み合わせて総合的に判断することが重要です。

ビッグデータ基礎

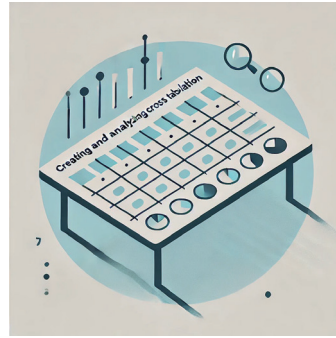
1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

クロス集計の分析



ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

クロス集計の概念

2つの変数の関係を表形式で表現するクロス集計の基本

表の作成方法

行変数と列変数の選択、度数のカウント方法

比率の計算

行比率、列比率、全体比率の計算方法と解釈

クロス集計は、2つのカテゴリ変数の組み合わせごとの度数を表形式で表現する方法です。例えば、性別と商品の購買有無、年齢層と職業など、2つの質的変数の関係を把握するのに適しています。これにより、変数間の関連性を直感的に理解することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

クロス集計の概念

2つの変数の関係を表形式で表現するクロス集計の基本

表の作成方法

行変数と列変数の選択、度数のカウント方法

比率の計算

行比率、列比率、全体比率の計算方法と解釈

表の作成では、まず行と列に配置する変数を選択します。一般的に、分析の主対象となる変数を行に、説明変数を列に配置します。次に、各セルに該当するデータの度数をカウントします。例えば、男性で商品を購入した人数、女性で商品を購入しなかった人数といった具合です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

クロス集計の概念

2つの変数の関係を表形式で表現するクロス集計の基本

表の作成方法

行変数と列変数の選択、度数のカウント方法

比率の計算

行比率、列比率、全体比率の計算方法と解釈

比率の計算には、行比率、列比率、全体比率の3種類があります。行比率は行方向の合計を100%とした場合の割合、列比率は列方向の合計を100%とした場合の割合、全体比率は全データ数を100%とした場合の割合です。分析の目的に応じて適切な比率を選択します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

期待度数

独立性を仮定した場合の期待度数の計算方法

残差分析

観測度数と期待度数の差から関係性を分析する方法

表の体裁

見やすいクロス集計表の作成方法と表示規則

期待度数は、2つの変数が独立である場合に期待される度数です。行の周辺度数と列の周辺度数の積を全体の度数で割ることで計算されます。この期待度数と実際の観測度数を比較することで、変数間の関連性を評価することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

期待度数

独立性を仮定した場合の期待度数の計算方法

残差分析

観測度数と期待度数の差から関係性を分析する方法

表の体裁

見やすいクロス集計表の作成方法と表示規則

残差分析では、観測度数から期待度数を引いた値（残差）を分析します。特に、残差を標準化した調整済み残差を用いることで、どのセルが有意に多いか少ないかを統計的に判断することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

期待度数

独立性を仮定した場合の期待度数の計算方法

残差分析

観測度数と期待度数の差から関係性を分析する方法

表の体裁

見やすいクロス集計表の作成方法と表示規則

表の体裁も重要です。適切な行・列の見出し、単位の明示、合計行・列の追加、有効数字の統一など、読み手が理解しやすいよう工夫が必要です。また、重要なセルの強調表示や色分けなども、結果の解釈を助ける効果があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計表の基礎

クロス集計の分析



ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計の分析

カイ二乗検定

変数間の独立性を検定する方法

クラメールの連関係数

関連の強さを数値化する方法

オッズ比

2値変数間の関連性を測るオッズ比の計算と解釈

カイ二乗検定は、2つの変数間に有意な関連があるかどうかを統計的に検定する方法です。観測度数と期待度数の差を基に計算されるカイ二乗統計量を用いて、変数間の独立性を検定します。p値が有意水準（一般的に0.05）より小さければ、変数間に関連があると判断できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計の分析

カイ二乗検定

変数間の独立性を検定する方法

クラメールの連関係数

関連の強さを数値化する方法

オッズ比

2値変数間の関連性を測るオッズ比の計算と解釈

クラメールの連関係数は、関連の強さを0から1の値で表す指標です。カイ二乗統計量を基に計算され、値が大きいほど強い関連があることを示します。この指標は、サンプルサイズの影響を調整しているため、異なる調査間での比較も可能です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計の分析

カイ二乗検定

変数間の独立性を検定する方法

クラメールの連関係数

関連の強さを数値化する方法

オッズ比

2値変数間の関連性を測るオッズ比の計算と解釈

オッズ比は、 2×2 のクロス集計表で特に有用な指標です。例えば、男女別の商品購入の有無を分析する場合、男性の購入オッズと女性の購入オッズの比を計算します。オッズ比が1より大きければ正の関連、1より小さければ負の関連があることを示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計の分析

条件付き確率

クロス集計表から条件付き確率を算出する方法

多重クロス集計

3つ以上の変数による分析方法

視覚化手法

クロス集計結果の効果的な図示方法

条件付き確率は、一方の変数の値が与えられた時の、他方の変数の確率を表します。例えば、商品を購入した人の中での男性の割合や、男性の中での商品購入者の割合などを計算することで、より詳細な関係性を把握できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計の分析

条件付き確率

クロス集計表から条件付き確率を算出する方法

多重クロス集計

3つ以上の変数による分析方法

視覚化手法

クロス集計結果の効果的な図示方法

多重クロス集計は、3つ以上の変数を組み合わせた分析です。例えば、性別×年齢層×職業のように、3次元以上の表を作成します。ただし、次元が増えるほど解釈が複雑になるため、目的に応じて適切な変数の組み合わせを選択する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-12. クロス集計表の作成と分析

第2章 データ分析基礎

クロス集計の分析

条件付き確率

クロス集計表から条件付き確率を算出する方法

多重クロス集計

3つ以上の変数による分析方法

視覚化手法

クロス集計結果の効果的な図示方法

視覚化手法としては、モザイクプロット、バブルチャート、積み上げ棒グラフなどがあります。特にモザイクプロットは、面積の大きさで度数を、色の濃さで残差を表現できるため、クロス集計の結果を効果的に可視化することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

クレンジングの実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

データクレンジングの目的

分析前のデータ品質向上とエラー除去の重要性

クレンジングの手順

データ確認、エラー検出、修正の基本的な流れ

データ品質の次元

完全性、正確性、一貫性、適時性品質の評価基準

データクレンジングは、分析の信頼性を確保するための重要な工程です。生のデータには様々なエラーや不整合が含まれており、これらを放置したまま分析を行うと、誤った結論を導く可能性があります。そのため、データの品質を向上させ、信頼できる分析基盤を整えることが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

データクレンジングの目的

分析前のデータ品質向上とエラー除去の重要性

クレンジングの手順

データ確認、エラー検出、修正の基本的な流れ

データ品質の次元

完全性、正確性、一貫性、適時性品質の評価基準

クレンジングの基本的な手順は、まずデータの全体像を把握することから始まります。データの件数、項目の種類、値の分布などを確認し、次に具体的なエラーの検出を行います。発見されたエラーは、定められたルールに従って修正や除去を行います。この一連の作業を体系的に実施することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

データクレンジングの目的

分析前のデータ品質向上とエラー除去の重要性

クレンジングの手順

データ確認、エラー検出、修正の基本的な流れ

データ品質の次元

完全性、正確性、一貫性、適時性品質の評価基準

データ品質は複数の次元から評価します。完全性（必要なデータがすべて揃っているか）、正確性（値は正しいか）、一貫性（データ間に矛盾はないか）、適時性（データは最新か）などが主要な評価基準となります。これらの基準に基づいて、データの品質レベルを総合的に判断します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

エラーの種類

入力ミス、形式不一致、重複一般的なエラーの分類

データの検証

論理チェック、範囲チェック基本的な検証方法

自動化の方法

データクレンジング作業の自動化手法

エラーには様々な種類があります。単純な入力ミス（タイプミスなど）、データ形式の不一致（日付形式が混在するなど）、重複データの存在、論理的な矛盾（生年月日と年齢が合わないなど）、異常値の混入などが代表的です。それぞれのエラーに対して、適切な対処方法を選択する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

エラーの種類

入力ミス、形式不一致、重複一般的なエラーの分類

データの検証

論理チェック、範囲チェック基本的な検証方法

自動化の方法

データクレンジング作業の自動化手法

データの検証では、様々なチェック方法を組み合わせます。値の範囲チェック（数値が妥当な範囲内かどうか）、論理チェック（データ間の関係は正しいか）、形式チェック（指定された形式に従っているか）などを実施します。このような多角的なチェックにより、エラーを漏れなく検出することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

エラーの種類

入力ミス、形式不一致、重複一般的なエラーの分類

データの検証

論理チェック、範囲チェック基本的な検証方法

自動化の方法

データクレンジング作業の自動化手法

自動化は、大規模データのクレンジングには不可欠です。プログラミングやデータ処理ツールを活用して、定型なチェックや修正を自動化します。ただし、完全な自動化は難しく、人による確認と判断も必要です。効率性と正確性のバランスを考慮した適切な自動化レベルを設定することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

データクレンジングの概要

クレンジングの実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

クレンジングの実践

形式の統一

日付、数値、文字列データ形式の統一方法

重複の排除

重複データの検出と除去の具体的な方法

表記揺れの修正

同じ内容の異なる表記を統一する方法

形式の統一は、データクレンジングの基本です。日付データは「yyyy/mm/dd」のような統一された形式に変換し、数値データは小数点の位置や桁数を揃え、文字列データは全角・半角や大文字・小文字を統一します。このような標準化により、後続の処理がスムーズになります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

クレンジングの実践

形式の統一

日付、数値、文字列データ形式の統一方法

重複の排除

重複データの検出と除去の具体的な方法

表記揺れの修正

同じ内容の異なる表記を統一する方法

重複データの排除は、分析結果の歪みを防ぐために重要です。完全な重複（すべての項目が同じ）と部分的な重複（一部の項目のみが同じ）を区別し、適切な処理を行います。例えば、顧客データの場合、名前と住所が微妙に異なる同一人物のレコードを特定し、統合する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

クレンジングの実践

形式の統一

日付、数値、文字列データ形式の統一方法

重複の排除

重複データの検出と除去の具体的な方法

表記揺れの修正

同じ内容の異なる表記を統一する方法

表記揺れの修正は、特に文字列データで重要です。例えば、「株式会社」「(株)」「株」など、同じ内容を表す異なる表記を統一形式に変換します。辞書やルールベースの置換処理を活用することで、効率的に修正を行うことができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

クレンジングの実践

誤入力の修正

明らかな誤入力を検出し修正する方法

データの結合

複数のデータソースを結合する際の注意点

品質チェック

クレンジング後の品質確認方法

誤入力の修正では、明らかな間違いを検出し、修正します。例えば、年齢が「234歳」となっているデータや、都道府県が定義された値以外になっているデータなどです。修正の根拠と方法を明確に記録し、トレーサビリティを確保することも重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

クレンジングの実践

誤入力の修正

明らかな誤入力を検出し修正する方法

データの結合

複数のデータソースを結合する際の注意点

品質チェック

クレンジング後の品質確認方法

データの結合では、異なるソースのデータを統合する際の注意点があります。キーとなる項目の定義や形式の違い、タイムスタンプのズレ、重複データの扱いなどを慎重に確認します。また、結合後のデータの整合性チェックも忘れずに行います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-13. データクレンジングの基本

第2章 データ分析基礎

クレンジングの実践

誤入力の修正

明らかな誤入力を検出し修正する方法

データの結合

複数のデータソースを結合する際の注意点

品質チェック

クレンジング後の品質確認方法

品質チェックでは、クレンジング後のデータが期待通りの品質レベルに達しているかを確認します。基本統計量の確認、サンプリングチェック、論理チェックなどを実施し、必要に応じて追加のクレンジングを行います。また、クレンジングの過程と結果を文書化し、後から確認できるようにすることも重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損値の処理



ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損の種類

完全無作為欠損、無作為欠損、非無作為欠損の違い

欠損パターン

単一変数の欠損、複数変数の欠損欠損パターンの分析方法

欠損の影響

分析結果に与える欠損の影響

欠損値には、発生メカニズムによって3つの種類があります。完全無作為欠損（MCAR）は、欠損が完全にランダムに発生する場合です。無作為欠損（MAR）は、他の観測された変数に依存して欠損が発生する場合です。非無作為欠損（MNAR）は、欠損する値自体に依存して欠損が発生する場合です。例えば、高所得者が収入の回答を避ける傾向があるような場合が該当します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損の種類

完全無作為欠損、無作為欠損、非無作為欠損の違い

欠損パターン

単一変数の欠損、複数変数の欠損欠損パターンの分析方法

欠損の影響

分析結果に与える欠損の影響

欠損パターンの分析では、欠損の発生状況を体系的に把握します。単一の変数のみに欠損がある場合は比較的単純ですが、複数の変数に欠損がある場合は、その組み合わせパターンを分析する必要があります。例えば、ある質問に答えなかった人が他の特定の質問にも答えない傾向があるかどうかを確認します。

ビッグデータ基礎 1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損の種類

完全無作為欠損、無作為欠損、非無作為欠損の違い

欠損パターン

単一変数の欠損、複数変数の欠損欠損パターンの分析方法

欠損の影響

分析結果に与える欠損の影響

欠損値は分析結果に重大な影響を与える可能性があります。サンプルサイズの減少による統計的検出力の低下、推定値のバイアス、相関関係の歪みなどが生じる可能性があります。特に、欠損が無作為でない場合、結果の解釈に注意が必要です。

ビッグデータ基礎 1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損率の計算

変数ごとや全体の欠損率の算出方法

欠損メカニズム

欠損が発生する原因と背景の分析方法

対処方針の決定

欠損値への対処方針を決定するための判断基準

欠損率の計算は、変数ごとと全体の両方について行います。変数ごとの欠損率は、その変数の欠損数を全サンプル数で割って算出します。全体の欠損率は、データセット全体での欠損値の割合を示します。一般的に、欠損率が5%未満なら比較的小さな問題とされますが、20%を超えると深刻な問題となる可能性があります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損率の計算

変数ごとや全体の欠損率の算出方法

欠損メカニズム

欠損が発生する原因と背景の分析方法

対処方針の決定

欠損値への対処方針を決定するための判断基準

欠損メカニズムの分析では、なぜデータが欠損しているのかを調査します。技術的な問題（センサーの故障など）、回答者の意図的な拒否、データ入力ミス、測定不能など、様々な原因が考えられます。この理解は、適切な対処方法を選択する上で重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損値の処理



ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の理解

欠損率の計算

変数ごとや全体の欠損率の算出方法

欠損メカニズム

欠損が発生する原因と背景の分析方法

対処方針の決定

欠損値への対処方針を決定するための判断基準

対処方針の決定では、欠損の種類、パターン、率、そしてメカニズムを総合的に考慮します。完全無作為欠損の場合は単純な処理で十分かもしれませんが、非無作為欠損の場合はより慎重な対応が必要です。また、分析の目的や重要性によっても、適切な対処方法は異なってきます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の処理

リストワイズ除去

欠損のあるケースを完全に除去する方法と影響

ペアワイズ除去

分析ごとに利用可能なケースを使用する方法

平均値代入

欠損値を平均値で補完する方法とその特徴

リストワイズ除去は、欠損値を含むケース（行）を完全に分析から除外する方法です。実装が簡単で解釈も容易ですが、多くのデータを失う可能性があります。例えば、100個の変数があり、各変数に5%の欠損がランダムにある場合、完全なケースは極めて少なくなってしまいます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の処理

リストワイズ除去

欠損のあるケースを完全に除去する方法と影響

ペアワイズ除去

分析ごとに利用可能なケースを使用する方法

平均値代入

欠損値を平均値で補完する方法とその特徴

ペアワイズ除去は、各分析に必要な変数のペアについて、両方の値が存在するケースのみを使用する方法です。相関係数の計算などでよく使われます。リストワイズ除去よりもデータの損失は少なくなりますが、異なる分析で異なるサンプルを使用することになり、結果の整合性に問題が生じる可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の処理

リストワイズ除去

欠損のあるケースを完全に除去する方法と影響

ペアワイズ除去

分析ごとに利用可能なケースを使用する方法

平均値代入

欠損値を平均値で補完する方法とその特徴

平均値代入は、欠損値を変数の平均値で置き換える方法です。実装が簡単で直感的にわかりやすい一方、データのばらつきを過小評価し、変数間の関係を歪める可能性があります。より洗練された方法として、グループごとの平均値を使用する条件付き平均値代入もあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の処理

回帰代入

他の変数を用いた予測値での補完方法

多重代入法

複数の補完値を生成する高度な方法

補完の評価

欠損値処理の妥当性を評価する方法

回帰代入は、他の変数を説明変数として欠損値を予測する方法です。例えば、年齢と収入の関係から、欠損している収入を推定します。より正確な推定が可能ですが、予測モデルの精度に結果が依存し、不確実性を過小評価する傾向があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の処理

回帰代入

他の変数を用いた予測値での補完方法

多重代入法

複数の補完値を生成する高度な方法

補完の評価

欠損値処理の妥当性を評価する方法

多重代入法は、統計的手法を用いて複数の補完値を生成し、それぞれで分析を行って結果を統合する方法です。不確実性を適切に考慮できる利点がありますが、計算が複雑で解釈も難しくなります。高度な統計ソフトウェアが必要となることも多いです。

ビッグデータ基礎

1 2 3 4 5 6 7

2-14. 欠損値の処理方法

第2章 データ分析基礎

欠損値の処理

回帰代入

他の変数を用いた予測値での補完方法

多重代入法

複数の補完値を生成する高度な方法

補完の評価

欠損値処理の妥当性を評価する方法

補完の評価では、処理の妥当性を確認します。感度分析（異なる処理方法で結果がどの程度変わるか）の実施、補完値の分布の確認、理論的な妥当性の検討などを行います。特に重要な分析では、複数の方法を試して結果の安定性を確認することが推奨されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

外れ値への対処



ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

外れ値の定義

統計的な観点から見た外れ値の定義と種類

検出方法

標準偏差法、四分位範囲法主要な検出方法

グラフによる確認

箱ひげ図、散布図視覚的な外れ値確認方法

外れ値とは、データの主要な分布から大きく離れた値のことです。統計的な観点からは、データの一般的な変動範囲を超えて極端に大きいまたは小さい値として定義されます。ただし、その判断基準は分野や状況によって異なり、必ずしも誤りというわけではありません。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

外れ値の定義

統計的な観点から見た外れ値の定義と種類

検出方法

標準偏差法、四分位範囲法主要な検出方法

グラフによる確認

箱ひげ図、散布図視覚的な外れ値確認方法

外れ値の検出方法には、いくつかの標準的なアプローチがあります。標準偏差法では、平均値から標準偏差の何倍か（一般的には ± 3 倍）離れた値を外れ値とみなします。四分位範囲法では、第1四分位数から第3四分位数までの範囲（IQR）の1.5倍を超えて外れた値を検出します。これらの方法は、データの分布特性に応じて使い分けます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

外れ値の定義

統計的な観点から見た外れ値の定義と種類

検出方法

標準偏差法、四分位範囲法主要な検出方法

グラフによる確認

箱ひげ図、散布図視覚的な外れ値確認方法

グラフによる確認は、外れ値の視覚的な検出に効果的です。箱ひげ図では、箱から伸びるひげの範囲を超える点として外れ値が表示されます。散布図では、全体的な分布パターンから著しく離れた点として確認できます。これらの視覚的手法により、データの全体像を把握しながら外れ値を特定できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

多変量の外れ値

複数変数を考慮した外れ値検出方法

異常値と外れ値

測定誤差による異常値と真の外れ値の区別

検出基準の設定

業務知識を踏まえた適切な検出基準の設定方法

多変量データの場合、複数の変数を同時に考慮した外れ値検出が必要です。マハラノビス距離やクック距離などの統計量を用いて、多次元空間での異常な観測値を検出します。例えば、身長と体重のように関連する複数の変数で、通常とは異なるパターンを示すデータを特定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

多変量の外れ値

複数変数を考慮した外れ値検出方法

異常値と外れ値

測定誤差による異常値と真の外れ値の区別

検出基準の設定

業務知識を踏まえた適切な検出基準の設定方法

異常値と外れ値の区別も重要です。測定機器の故障や入力ミスによる明らかな誤りは異常値として扱い、通常は修正や除去の対象となります。一方、稀少だが実際に存在する現象を示す外れ値は、重要な情報を含んでいる可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

多変量の外れ値

複数変数を考慮した外れ値検出方法

異常値と外れ値

測定誤差による異常値と真の外れ値の区別

検出基準の設定

業務知識を踏まえた適切な検出基準の設定方法

検出基準の設定では、統計的な基準に加えて、業務知識や専門家の判断を考慮します。例えば、製造工程の品質管理では、製品の仕様限界値を考慮した基準を設定します。また、季節変動や特殊なイベントの影響なども考慮に入れる必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値の特定

外れ値への対処



ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値への対処

除去の判断

外れ値を除去すべきかどうかの判断基準

修正方法

外れ値の値を修正する際の具体的な方法

影響の評価

外れ値が分析結果に与える影響の評価方法

外れ値の除去判断は慎重に行う必要があります。明らかな測定エラーや入力ミスによる異常値は除去の対象となりますが、実際の現象を反映した外れ値は、重要な情報源となる可能性があります。例えば、顧客の購買データにおける大口取引は、重要な顧客セグメントを示している可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値への対処

除去の判断

外れ値を除去すべきかどうかの判断基準

修正方法

外れ値の値を修正する際の具体的な方法

影響の評価

外れ値が分析結果に与える影響の評価方法

修正方法には、いくつかのアプローチがあります。異常値の場合は、正しい値が分かれば置き換え、不明な場合は欠損値として扱います。外れ値の場合は、ウインザー化（一定の範囲を超える値を境界値で置き換える）やトリミング（極端な値を除去する）などの手法が用いられます。

ビッグデータ基礎 1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値への対処

除去の判断

外れ値を除去すべきかどうかの判断基準

修正方法

外れ値の値を修正する際の具体的な方法

影響の評価

外れ値が分析結果に与える影響の評価方法

影響の評価では、外れ値が分析結果にどの程度影響を与えるかを確認します。外れ値を含めた場合と除去した場合の結果を比較したり、感度分析を行ったりすることで、その影響度を評価します。影響が大きい場合は、より慎重な対処が必要です。

ビッグデータ基礎 1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値への対処

ロバスト統計

外れ値の影響を受けにくい統計手法の適用

記録と報告

外れ値の処理内容の記録と報告方法

予防措置

将来的な外れ値発生を防ぐための対策

ロバスト統計手法は、外れ値の影響を受けにくい分析手法です。例えば、平均値の代わりに中央値を使用したり、通常の回帰分析の代わりにロバスト回帰を使用したりします。これらの手法により、外れ値を除去せずに安定した分析結果を得ることができます。

ビッグデータ基礎 1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値への対処

ロバスト統計

外れ値の影響を受けにくい統計手法の適用

記録と報告

外れ値の処理内容の記録と報告方法

予防措置

将来的な外れ値発生を防ぐための対策

処理内容の記録と報告は重要です。どの値を外れ値と判断したか、その理由は何か、どのように処理したかを詳細に記録します。また、分析結果を報告する際には、外れ値の処理方法とその影響について明確に説明する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-15. 外れ値の処理方法

第2章 データ分析基礎

外れ値への対処

ロバスト統計

外れ値の影響を受けにくい統計手法の適用

記録と報告

外れ値の処理内容の記録と報告方法

予防措置

将来的な外れ値発生を防ぐための対策

予防措置としては、データ収集段階でのチェック機能の強化、測定機器の定期的な校正、入力時のバリデーション強化などが考えられます。また、過去の外れ値事例を分析し、発生パターンを理解することで、より効果的な予防策を講じることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

データの標準化の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

正規化の概念

データを特定の範囲内に収める正規化の基本的な考え方と目的

最小最大正規化

データを0から1の範囲に変換する方法とその特徴

小数スケーリング

小数点の位置を調整する正規化手法

正規化とは、異なるスケールや範囲を持つデータを、特定の範囲内に収めるための変換処理です。例えば、0から100までの得点と、0から1000までの得点を比較する場合、同じスケールに変換することで、公平な比較が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

正規化の概念

データを特定の範囲内に収める正規化の基本的な考え方と目的

最小最大正規化

データを0から1の範囲に変換する方法とその特徴

小数スケーリング

小数点の位置を調整する正規化手法

最小最大正規化は、データを0から1の範囲に変換する最も一般的な方法です。データの最小値を0、最大値を1に対応させ、その間の値を比例的に変換します。例えば、テストの得点を正規化することで、異なる満点のテスト結果を统一的に比較できるようになります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

正規化の概念

データを特定の範囲内に収める正規化の基本的な考え方と目的

最小最大正規化

データを0から1の範囲に変換する方法とその特徴

小数スケーリング

小数点の位置を調整する正規化手法

小数スケーリングは、データの桁数を調整する簡単な正規化方法です。例えば、1000以上の値を持つデータを10や100で割ることで、扱いやすい範囲に収めます。為替レートや株価など、大きな数値を扱う際によく使用されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

対数変換

データの分布を調整する対数変換の方法と効果

べき変換

データの非線形性に対処するべき変換

正規化の選択

データの特성에応じた適切な正規化手法の選択方法

対数変換は、データの分布を調整する強力な手法です。特に、データが広い範囲に分布している場合や、右に歪んだ分布を示す場合に効果的です。例えば、企業の売上高や人口データなど、桁数の異なる値を扱う際によく使用されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

対数変換

データの分布を調整する対数変換の方法と効果

べき変換

データの非線形性に対処するべき変換

正規化の選択

データの特性に適切な正規化手法の選択方法

べき変換は、データの非線形性に対処するために使用します。データを特定の指数でべき乗することで、分布の形状を調整します。対数変換よりも柔軟に分布を変形できる利点がありますが、適切な指数の選択が重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

対数変換

データの分布を調整する対数変換の方法と効果

べき変換

データの非線形性に対処するべき変換

正規化の選択

データの特性に適切な正規化手法の選択方法

正規化手法の選択は、データの特性と分析の目的に応じて行います。データの分布の形状、外れ値の有無、値の範囲などを考慮します。また、変換後のデータの解釈のしやすさも、選択の重要な基準となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの正規化の基礎

データの標準化の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの標準化の実践

標準化の概念

平均0、分散1に変換する標準化の基本的な考え方

Z得点化

標準正規分布に従うようにデータを変換する方法

ロバスト標準化

外れ値の影響を受けにくい標準化手法

標準化は、データの平均を0、標準偏差を1に変換する処理です。これにより、異なる単位や尺度で測定されたデータを、共通の基準で比較できるようになります。例えば、身長 (cm) と体重 (kg) のような異なる単位のデータを、同じ基準で分析することが可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの標準化の実践

標準化の概念

平均0、分散1に変換する標準化の基本的な考え方

Z得点化

標準正規分布に従うようにデータを変換する方法

ロバスト標準化

外れ値の影響を受けにくい標準化手法

Z得点化は、各データから平均値を引き、標準偏差で割る処理です。この変換により、データは標準正規分布に従うような形に変換されます。例えば、テストの得点をZ得点化することで、異なる科目間の成績を比較したり、偏差値を算出したりすることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの標準化の実践

標準化の概念

平均0、分散1に変換する標準化の基本的な考え方

Z得点化

標準正規分布に従うようにデータを変換する方法

ロバスト標準化

外れ値の影響を受けにくい標準化手法

ロバスト標準化は、外れ値の影響を受けにくい手法です。平均値の代わりに中央値を、標準偏差の代わりに四分位範囲を使用することで、極端な値の影響を軽減します。特に、外れ値を含むデータセットの処理に適しています。

ビッグデータ基礎 1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの標準化の実践

尺度の統一

異なる尺度のデータを比較可能にする方法

標準化の影響

分析結果に対する標準化の影響

元データへの復元

標準化されたデータを元の尺度に戻す方法

尺度の統一は、異なる測定尺度のデータを比較可能な形に変換します。例えば、5段階評価と7段階評価のアンケート結果を比較する際に、標準化によって統一的な評価が可能になります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの標準化の実践

尺度の統一

異なる尺度のデータを比較可能にする方法

標準化の影響

分析結果に対する標準化の影響

元データへの復元

標準化されたデータを元の尺度に戻す方法

標準化の影響は、分析手法によって異なります。主成分分析や因子分析などの多変量解析では、標準化が必須となります。一方、単純な平均値の比較などでは、元のスケールのまま分析の方が解釈しやすい場合もあります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-16. データの正規化と標準化

第2章 データ分析基礎

データの標準化の実践

尺度の統一

異なる尺度のデータを比較可能にする方法

標準化の影響

分析結果に対する標準化の影響

元データへの復元

標準化されたデータを元の尺度に戻す方法

元データへの復元は、標準化の逆変換によって行います。標準化されたデータに標準偏差を掛け、平均値を足すことで、元のスケールに戻すことができます。この復元可能性は、分析結果を元の単位で解釈する必要がある場合に重要となります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

エンコーディング手法



ビッグデータ基礎

1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

カテゴリデータの種類

名義尺度と順序尺度のカテゴリデータの特徴

エンコードの必要性

機械学習や統計分析におけるカテゴリデータ変換の重要性

エンコード方式の種類

様々なエンコーディング手法の概要と特徴

カテゴリデータには、大きく分けて2種類があります。名義尺度は、性別や血液型のように単純な分類を表すデータで、カテゴリ間に順序関係はありません。一方、順序尺度は、満足度評価（非常に満足、やや満足など）のように、カテゴリ間に順序関係があるデータです。

ビッグデータ基礎

1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

カテゴリデータの種類

名義尺度と順序尺度のカテゴリデータの特徴

エンコードの必要性

機械学習や統計分析におけるカテゴリデータ変換の重要性

エンコード方式の種類

様々なエンコーディング手法の概要と特徴

エンコードが必要となる理由は、多くの統計手法や機械学習アルゴリズムが数値データを前提としているためです。例えば、「男性」「女性」という文字列データは、そのままでは計算に使用できません。これらを数値に変換することで、様々な分析手法を適用することが可能になります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

カテゴリデータの種類

名義尺度と順序尺度のカテゴリデータの特徴

エンコードの必要性

機械学習や統計分析におけるカテゴリデータ変換の重要性

エンコード方式の種類

様々なエンコーディング手法の概要と特徴

エンコード方式には様々な手法があり、データの性質や分析の目的に応じて選択します。単純な数値への置き換え、バイナリ変数への展開、統計量を用いた変換など、それぞれに特徴と適用場面があります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

前処理の重要性

エンコーディング前のデータクリーニングと準備

カーディナリティ

カテゴリの数が多い場合の考慮事項

欠損カテゴリ

新しいカテゴリや欠損値の処理方法

前処理の重要性は、エンコーディングの成否を左右します。まず、表記揺れの修正（「男性」「Male」「M」などの統一）や、誤入力の修正を行います。また、カテゴリの粒度が適切かどうかの検討も重要です。例えば、都道府県データを地域ブロックにまとめるなどの調整が必要な場合もあります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

前処理の重要性

エンコーディング前のデータクリーニングと準備

カーディナリティ

カテゴリの数が多の場合の考慮事項

欠損カテゴリ

新しいカテゴリや欠損値の処理方法

カーディナリティ（カテゴリの数）が高い場合は特別な考慮が必要です。例えば、商品コードのように数千、数万のカテゴリがある場合、単純なエンコーディングでは次元が爆発的に増加し、分析が困難になる可能性があります。このような場合は、カテゴリの集約や次元削減手法の適用を検討します。

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

前処理の重要性

エンコーディング前のデータクリーニングと準備

カーディナリティ

カテゴリの数が多の場合の考慮事項

欠損カテゴリ

新しいカテゴリや欠損値の処理方法

欠損カテゴリの処理も重要な課題です。訓練データには存在しないが、将来的に発生する可能性のあるカテゴリや、欠損値をどのように扱うかを事前に決めておく必要があります。一般的には、「その他」や「不明」といった特別なカテゴリを用意します。

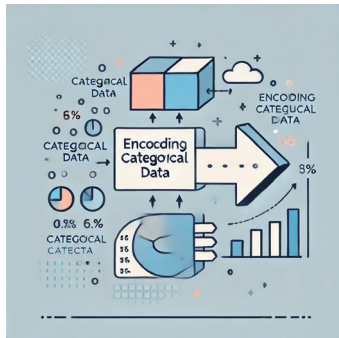
NEXT

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

カテゴリデータの基本

エンコーディング手法



2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

エンコーディング手法

ラベルエンコーディング

カテゴリに連番を割り当てる方法とその適用場面

ワンホットエンコーディング

カテゴリをダミー変数に変換する方法と特徴

ターゲットエンコーディング

目的変数の統計量を用いたエンコーディング

ラベルエンコーディングは、各カテゴリに連番を割り当てる最も単純な方法です。例えば、「赤」「青」「緑」を0、1、2に変換します。実装が簡単で、メモリ効率が良いという利点がありますが、数値の大小関係が意味を持たない名義尺度データには適していません。

ビッグデータ基礎 1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

エンコーディング手法

ラベルエンコーディング

カテゴリに連番を割り当てる方法とその適用場面

ワンホットエンコーディング

カテゴリをダミー変数に変換する方法と特徴

ターゲットエンコーディング

目的変数の統計量を用いたエンコーディング

ワンホットエンコーディングは、各カテゴリを0と1のバイナリ変数の組み合わせに変換します。例えば、「赤」「青」「緑」は、それぞれ[1,0,0]、[0,1,0]、[0,0,1]となります。カテゴリ間の大小関係を作らないため、名義尺度データの処理に適していますが、カテゴリ数が多いと変数が増えすぎる欠点があります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

エンコーディング手法

ラベルエンコーディング

カテゴリに連番を割り当てる方法とその適用場面

ワンホットエンコーディング

カテゴリをダミー変数に変換する方法と特徴

ターゲットエンコーディング

目的変数の統計量を用いたエンコーディング

ターゲットエンコーディングは、目的変数との関連性を利用する手法です。各カテゴリについて、そのカテゴリに属するサンプルの目的変数の平均値などを割り当てます。例えば、商品カテゴリごとの平均購買額を用いてエンコードします。過学習のリスクがあるため、適切な検証が必要です。

ビッグデータ基礎 1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

エンコーディング手法

頻度エンコーディング

出現頻度を用いたエンコーディング手法

順序エンコーディング

順序のあるカテゴリデータの変換方法

ハッシュエンコーディング

高次元カテゴリデータの次元削減方法

頻度エンコーディングは、各カテゴリの出現頻度を用いる手法です。例えば、全体の30%を占めるカテゴリには0.3を割り当てます。カテゴリの一般性や希少性を数値化できる利点がありますが、頻度の違いが実際の意味を反映していない場合もあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

エンコーディング手法

頻度エンコーディング

出現頻度を用いたエンコーディング手法

順序エンコーディング

順序のあるカテゴリデータの変換方法

ハッシュエンコーディング

高次元カテゴリデータの次元削減方法

順序エンコーディングは、順序尺度データに適した手法です。例えば、「低」「中」「高」を[0,0]、[1,0]、[1,1]のように累積的なバイナリ値に変換します。カテゴリ間の順序関係を保持しながら、過度な数値的意味付けを避けることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-17. カテゴリデータのエンコーディング

第2章 データ分析基礎

エンコーディング手法

頻度エンコーディング

出現頻度を用いたエンコーディング手法

順序エンコーディング

順序のあるカテゴリデータの変換方法

ハッシュエンコーディング

高次元カテゴリデータの次元削減方法

ハッシュエンコーディングは、高次元カテゴリデータの次元削減に使用します。ハッシュ関数を使用して、大量のカテゴリを少数の次元に圧縮します。情報の損失が発生する可能性がありますが、高カーディナリティデータの処理に有効です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

時系列データの処理



ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

時系列データの構造

時間軸に沿ったデータの基本的な構造と特徴

時間単位の選択

分析目的に応じた適切な時間単位の設定方法

トレンド成分

長期的な傾向を示すトレンドの抽出方法

時系列データは、時間の経過に沿って記録された連続的なデータです。例えば、日々の売上高、毎時間の気温、毎月の製造量など、時間軸に沿って値が変化するデータを指します。時系列データの特徴として、データ点間に時間的な順序と依存関係があることが挙げられます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

時系列データの構造

時間軸に沿ったデータの基本的な構造と特徴

時間単位の選択

分析目的に応じた適切な時間単位の設定方法

トレンド成分

長期的な傾向を示すトレンドの抽出方法

時間単位の選択は、分析の目的に応じて重要な決定となります。例えば、小売店の売上分析では、時間帯別、日次、週次、月次など、様々な時間単位が考えられます。より細かい単位では詳細な変動が観察できますが、ノイズも多くなります。逆に、大きな単位では全体的な傾向は把握しやすくなりますが、細かな変動は見えなくなります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

時系列データの構造

時間軸に沿ったデータの基本的な構造と特徴

時間単位の選択

分析目的に応じた適切な時間単位の設定方法

トレンド成分

長期的な傾向を示すトレンドの抽出方法

トレンド成分は、データの長期的な変化傾向を示します。例えば、企業の売上高が年々増加している傾向や、気温が徐々に上昇している傾向などです。このトレンドは、線形的な変化だけでなく、曲線的な変化を示すこともあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

季節性成分

周期的なパターンを示す季節性の分析方法

ノイズ成分

不規則な変動を示すノイズの扱い方

欠損値の特徴

時系列データ特有の欠損値の問題と対処法

季節性成分は、一定の周期で繰り返されるパターンを指します。例えば、小売業での年末の売上増加、夏季の電力消費量の増加など、カレンダーや季節に関連した周期的な変動が該当します。この季節性を理解することは、将来予測や異常検知に重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

季節性成分

周期的なパターンを示す季節性の分析方法

ノイズ成分

不規則な変動を示すノイズの扱い方

欠損値の特徴

時系列データ特有の欠損値の問題と対処法

ノイズ成分は、トレンドや季節性では説明できない不規則な変動です。これらの変動は、一時的な要因や偶然によって生じる変動を含みます。ノイズを適切に扱うことで、より正確な分析が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

季節性成分

周期的なパターンを示す季節性の分析方法

ノイズ成分

不規則な変動を示すノイズの扱い方

欠損値の特徴

時系列データ特有の欠損値の問題と対処法

時系列データの欠損値には特有の課題があります。例えば、センサーの故障による欠測、休日のデータ欠損、不定期なサンプリングによるギャップなどです。これらの欠損値は、時系列の連続性を考慮しながら、適切な補完方法を選択する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの特徴

時系列データの処理



ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの処理

データの整形

時系列インデックスの作成と日付処理

リサンプリング

データの時間間隔を変更する方法

移動平均

データの平滑化と傾向抽出の方法

データの整形では、まず適切な時系列インデックスを作成します。日付や時刻を統一された形式で扱えるように変換し、時間順に並べ替えます。また、タイムゾーンの調整や日付演算の処理も重要です。例えば、営業日の計算や期間の計算などが必要になることがあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの処理

データの整形

時系列インデックスの作成と日付処理

リサンプリング

データの時間間隔を変更する方法

移動平均

データの平滑化と傾向抽出の方法

リサンプリングは、データの時間間隔を変更する処理です。例えば、分単位のデータを時間単位に集約したり、日次データを週次や月次に変換したりします。この際、平均値、合計値、最大値など、適切な集約方法を選択する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの処理

データの整形

時系列インデックスの作成と日付処理

リサンプリング

データの時間間隔を変更する方法

移動平均

データの平滑化と傾向抽出の方法

移動平均は、データの短期的な変動を滑らかにし、全体的な傾向を把握するための手法です。例えば、7日移動平均を計算することで、週単位での傾向が見やすくなります。移動平均の期間は、分析の目的に応じて適切に設定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの処理

季節調整

季節性を除去してトレンドを分析する方法

ラグ特徴量

過去の値を特徴量として利用する方法

時系列の可視化

効果的な時系列データの可視化手法

季節調整は、季節性の影響を除去してトレンドを分析する手法です。例えば、小売業の売上データから季節的な変動を除去することで、真の成長トレンドを把握することができます。X-12-ARIMAなどの標準的な季節調整法が広く使用されています。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの処理

季節調整

季節性を除去してトレンドを分析する方法

ラグ特徴量

過去の値を特徴量として利用する方法

時系列の可視化

効果的な時系列データの可視化手法

ラグ特徴量は、過去の値を説明変数として利用する手法です。例えば、1日前、1週間前、1ヶ月前の値を特徴量として使用することで、時系列的な依存関係を考慮した分析が可能になります。ただし、適切なラグ期間の選択が重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-18. 時系列データの基本的な扱い方

第2章 データ分析基礎

時系列データの処理

季節調整

季節性を除去してトレンドを分析する方法

ラグ特徴量

過去の値を特徴量として利用する方法

時系列の可視化

効果的な時系列データの可視化手法

時系列の可視化では、線グラフが基本となりますが、様々な工夫が可能です。トレンドラインの追加、季節性の強調表示、異常値のハイライト、信頼区間の表示など、データの特徴を効果的に伝える表現方法を選択します。また、インタラクティブな可視化ツールを使用することで、詳細な分析も容易になります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

サンプリング手法の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

サンプリングの目的

母集団からの標本抽出の意義と目的

サンプルサイズ

適切なサンプルサイズの決定方法と考慮点

代表性の確保

母集団の特徴を反映したサンプリング方法

サンプリングの目的は、母集団全体を調査することが困難な場合に、一部のデータを抽出して母集団の特性を推定することです。適切なサンプリングにより、時間とコストを節約しながら、信頼性の高い分析が可能となります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

サンプリングの目的

母集団からの標本抽出の意義と目的

サンプルサイズ

適切なサンプルサイズの決定方法と考慮点

代表性の確保

母集団の特徴を反映したサンプリング方法

サンプルサイズの決定は、求める精度と信頼水準に基づいて行います。大きなサンプルサイズは精度を向上させますが、コストと時間も増加します。必要な精度、許容される誤差範囲、予算制約などを考慮して、適切なサイズを決定する必要があります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

サンプリングの目的

母集団からの標本抽出の意義と目的

サンプルサイズ

適切なサンプルサイズの決定方法と考慮点

代表性の確保

母集団の特徴を反映したサンプリング方法

代表性の確保は、サンプリングの信頼性を左右する重要な要素です。母集団の特徴をバランスよく反映したサンプルを抽出するため、適切な抽出方法の選択と、必要に応じた層化や重み付けを行います。

ビッグデータ基礎 1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

バイアスの回避

サンプリングバイアスの種類と防止方法

信頼区間

サンプリング誤差と信頼区間の関係

コスト考慮

効率的なサンプリングのためのコスト考慮

バイアスの回避は、サンプリングの精度を確保するために不可欠です。選択バイアス、非回答バイアス、生存バイアスなど、様々な種類のバイアスが存在します。これらを認識し、適切な対策を講じることで、より信頼性の高いサンプリングが可能となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

バイアスの回避

サンプリングバイアスの種類と防止方法

信頼区間

サンプリング誤差と信頼区間の関係

コスト考慮

効率的なサンプリングのためのコスト考慮

信頼区間の設定では、サンプリング誤差を考慮した推定の範囲を定めます。信頼水準（例：95%）と許容誤差を設定し、その範囲内で母集団の特性を推定します。サンプルサイズが大きくなるほど、信頼区間は狭くなり、より正確な推定が可能となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

バイアスの回避

サンプリングバイアスの種類と防止方法

信頼区間

サンプリング誤差と信頼区間の関係

コスト考慮

効率的なサンプリングのためのコスト考慮

コスト考慮は、実践的なサンプリング計画に不可欠です。データ収集コスト、人的リソース、時間的制約などを考慮し、効率的なサンプリング方法を選択します。また、段階的なサンプリングや予備調査の実施により、コストを最適化することも重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリングの基礎

サンプリング手法の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリング手法の実践

単純無作為抽出

完全にランダムなサンプリング手法の実施方法

層化抽出

母集団を層に分けてサンプリングする方法

系統抽出

一定間隔でサンプルを抽出する方法

単純無作為抽出は、最も基本的なサンプリング手法です。乱数を使用して、母集団から完全にランダムにサンプルを抽出します。この方法は、実装が容易で偏りが少ないという利点がありますが、母集団の特性によっては代表性が確保されない可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリング手法の実践

単純無作為抽出

完全にランダムなサンプリング手法の実施方法

層化抽出

母集団を層に分けてサンプリングする方法

系統抽出

一定間隔でサンプルを抽出する方法

層化抽出は、母集団を特定の特性に基づいて層に分割し、各層から適切な比率でサンプルを抽出する方法です。例えば、年齢層や地域などで層を分け、各層の大きさに応じた数のサンプルを抽出します。これにより、母集団の構造を反映したサンプリングが可能となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリング手法の実践

単純無作為抽出

完全にランダムなサンプリング手法の実施方法

層化抽出

母集団を層に分けてサンプリングする方法

系統抽出

一定間隔でサンプルを抽出する方法

系統抽出は、一定の間隔でサンプルを抽出する方法です。例えば、母集団を順番に並べ、k番目ごとにサンプルを選択します。実装が簡単で、均等な分布が得られやすいという利点がありますが、周期性のあるデータでは偏りが生じる可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリング手法の実践

クラスタ抽出

グループ単位でサンプリングする方法

不均衡データの処理

アンダーサンプリングとオーバーサンプリング

サンプリング評価

抽出されたサンプルの品質評価方法

クラスタ抽出は、母集団をグループ（クラスタ）に分け、グループ単位でサンプリングを行う方法です。例えば、地域ごとにグループ化し、選択された地域内のすべてのデータを収集します。調査の効率性が高まりますが、クラスタ間の違いが大きい場合は精度が低下する可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリング手法の実践

クラスタ抽出

グループ単位でサンプリングする方法

不均衡データの処理

アンダーサンプリングとオーバーサンプリング

サンプリング評価

抽出されたサンプルの品質評価方法

不均衡データの処理では、クラス間でサンプル数に大きな差がある場合の対処方法を考えます。アンダーサンプリングは多数クラスのサンプル数を減らし、オーバーサンプリングは少数クラスのサンプル数を増やす方法です。これらを組み合わせることで、バランスの取れたデータセットを作成できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-19. データのサンプリング手法

第2章 データ分析基礎

サンプリング手法の実践

クラスタ抽出

グループ単位でサンプリングする方法

不均衡データの処理

アンダーサンプリングとオーバーサンプリング

サンプリング評価

抽出されたサンプルの品質評価方法

サンプリング評価では、抽出されたサンプルの品質を確認します。母集団との特性の比較、統計的検定、クロスバリデーションなどの手法を用いて、サンプルの代表性と信頼性を評価します。また、必要に応じてサンプリング方法の調整や再サンプリングを行います。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の基本概念

仮説検定の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の基本概念

仮説検定の目的

データに基づいて統計的な判断を行う仮説検定の基本的な考え方

帰無仮説と対立仮説

検定における二つの仮説の設定方法と意味

有意水準

判断基準となる有意水準の設定と解釈

仮説検定は、データに基づいて統計的な判断を行うための体系的な方法です。例えば、新薬の効果や広告施策の効果など、何らかの効果や差異が本当に存在するのかわ、データを用いて科学的に判断します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の基本概念

仮説検定の目的

データに基づいて統計的な判断を行う仮説検定の基本的な考え方

帰無仮説と対立仮説

検定における二つの仮説の設定方法と意味

有意水準

判断基準となる有意水準の設定と解釈

仮説検定では、二つの仮説を設定します。帰無仮説は「差がない」「効果がない」という、否定したい仮説です。対立仮説は「差がある」「効果がある」という、証明したい仮説です。例えば、新薬の効果を検証する場合、帰無仮説は「新薬には効果がない」、対立仮説は「新薬には効果がある」となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の基本概念

仮説検定の目的

データに基づいて統計的な判断を行う仮説検定の基本的な考え方

帰無仮説と対立仮説

検定における二つの仮説の設定方法と意味

有意水準

判断基準となる有意水準の設定と解釈

有意水準は、検定結果を判断する基準となる値です。一般的に5% (0.05) や1% (0.01) が使用されます。これは、帰無仮説が正しいのに誤って棄却してしまう確率の許容限界を示します。有意水準が厳しいほど、より確実な証拠が必要になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の基本概念

p値の概念

統計的有意性を示すp値の意味と解釈方法

検定統計量

仮説の判定に用いる検定統計量の計算方法

検定力

検定の精度を示す検定力とサンプルサイズの関係

p値は、観測されたデータが帰無仮説のもとでどれくらい起こりにくいかを示す値です。例えば、p値が0.03の場合、帰無仮説が正しければ3%の確率でしか観測されないような結果が得られたことを意味します。p値が有意水準より小さい場合、帰無仮説を棄却します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

仮説検定の基本概念	
p値の概念	統計的有意性を示すp値の意味と解釈方法
検定統計量	仮説の判定に用いる検定統計量の計算方法
検定力	検定の精度を示す検定力とサンプルサイズの関係

検定統計量は、データから計算される値で、仮説の判定に使用されます。例えば、t検定ではt値、カイ二乗検定ではカイ二乗値を用います。この値が特定の基準値を超えると、帰無仮説を棄却することになります。

2-20. 仮説検定の基本的な考え方

仮説検定の基本概念	
p値の概念	統計的有意性を示すp値の意味と解釈方法
検定統計量	仮説の判定に用いる検定統計量の計算方法
検定力	検定の精度を示す検定力とサンプルサイズの関係

検定力は、実際に差や効果が存在する場合に、それを検出できる確率を示します。検定力は主にサンプルサイズに依存し、サンプルサイズが大きいほど、小さな差も検出できるようになります。ただし、必要以上に大きなサンプルサイズは、些細な差も有意と判定してしまう可能性があります。

NEXT

2-20. 仮説検定の基本的な考え方

仮説検定の基本概念

仮説検定の実践



2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の実践

両側検定と片側検定

検定の方向性による違いと選択基準

第一種の誤り

帰無仮説を誤って棄却してしまう誤りのリスク

第二種の誤り

帰無仮説を誤って採択してしまう誤りのリスク

両側検定と片側検定は、検定の方向性に関する選択です。両側検定は「差がある」ことを検証する場合に用い、片側検定は「大きい」「小さい」という特定の方向性を検証する場合に使用します。例えば、新薬が従来薬より効果が高いことを証明したい場合は片側検定を選択します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の実践

両側検定と片側検定

検定の方向性による違いと選択基準

第一種の誤り

帰無仮説を誤って棄却してしまう誤りのリスク

第二種の誤り

帰無仮説を誤って採択してしまう誤りのリスク

第一種の誤りは、帰無仮説が実際には正しいのに、誤って棄却してしまう誤りです。例えば、実際には効果のない薬を「効果がある」と判断してしまうような場合です。この誤りの確率は有意水準によって制御されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の実践

両側検定と片側検定

検定の方向性による違いと選択基準

第一種の誤り

帰無仮説を誤って棄却してしまう誤りのリスク

第二種の誤り

帰無仮説を誤って採択してしまう誤りのリスク

第二種の誤りは、帰無仮説が実際には誤っているのに、誤って採択してしまう誤りです。例えば、実際には効果のある薬を「効果がない」と判断してしまうような場合です。この誤りの確率は検定力と関連しています。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の実践

多重検定

複数の検定を行う際の注意点と対処方法

効果量

統計的有意性と実践的重要性の違い

結果の解釈

検定結果の正しい解釈と報告方法

多重検定は、複数の検定を同時に行う場合の問題です。検定回数が増えると、偶然による有意な結果が得られる確率が高くなります。これに対しては、ボンフェローニ補正など、有意水準を調整する方法が用いられます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の実践

多重検定

複数の検定を行う際の注意点と対処方法

効果量

統計的有意性と実践的重要性の違い

結果の解釈

検定結果の正しい解釈と報告方法

効果量は、統計的有意性とは別に、実践的な重要性を示す指標です。大規模なサンプルでは些細な差でも統計的に有意になりやすいため、効果量を併せて報告することで、より適切な判断が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-20. 仮説検定の基本的な考え方

第2章 データ分析基礎

仮説検定の実践

多重検定

複数の検定を行う際の注意点と対処方法

効果量

統計的有意性と実践的重要性の違い

結果の解釈

検定結果の正しい解釈と報告方法

結果の解釈では、p値や検定統計量だけでなく、効果量や信頼区間も含めて総合的に判断することが重要です。また、統計的有意性が実践的な重要性を必ずしも意味しないことや、因果関係の推論には慎重である必要があることにも注意が必要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

t検定の実施手順



ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

t検定の種類

1 標本、対応のある 2 標本、独立した 2 標本の t 検定

前提条件

正規性や等分散性 t 検定の適用条件

検定統計量

t 値の計算方法と自由度の考え方

t 検定には主に 3 つの種類があります。1 標本の t 検定は、ある集団の平均値が特定の値と異なるかを検証します。例えば、新製品の重量が規格値と異なるかどうかを確認する場合に使用します。対応のある 2 標本の t 検定は、同じ対象の前後差を比較します。例えば、トレーニング前後の成績の変化を検証する場合です。独立した 2 標本の t 検定は、異なる 2 群の平均値を比較します。例えば、新薬と従来薬の効果の差を比較する場合に使用します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

t検定の種類

1 標本、対応のある 2 標本、独立した 2 標本の t 検定

前提条件

正規性や等分散性 t 検定の適用条件

検定統計量

t 値の計算方法と自由度の考え方

t 検定の前提条件として、データが正規分布に従うことと、比較する群の分散が等しいこと（等分散性）が重要です。ただし、サンプルサイズが十分大きい場合は、中心極限定理により、正規性の仮定は緩和されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

t検定の種類

1 標本、対応のある2 標本、独立した2 標本のt検定

前提条件

正規性や等分散性t検定の適用条件

検定統計量

t値の計算方法と自由度の考え方

検定統計量であるt値は、群間の平均値の差を標準誤差で割って算出します。自由度は、サンプルサイズから計算され、t分布の形状を決定する重要なパラメータとなります。例えば、独立した2標本のt検定では、両群のサンプルサイズの合計から2を引いた値が自由度となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

有意水準の設定

適切な有意水準の選択と判断基準

効果量の計算

Cohen's dなどの効果量の計算方法と解釈

検定力分析

必要なサンプルサイズの決定方法

有意水準は、一般的に5%や1%に設定されます。ただし、検定の重要性や結果の影響度に応じて、より厳格な値を設定することもあります。両側検定か片側検定かの選択も、研究の目的に応じて適切に行う必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

有意水準の設定

適切な有意水準の選択と判断基準

効果量の計算

Cohen's dなどの効果量の計算方法と解釈

検定力分析

必要なサンプルサイズの決定方法

効果量は、差の大きさを標準化した指標です。Cohen's dは最も一般的な効果量の一つで、平均値の差を標準偏差で割って算出します。一般的に、0.2が小さな効果、0.5が中程度の効果、0.8が大きな効果とされます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

有意水準の設定

適切な有意水準の選択と判断基準

効果量の計算

Cohen's dなどの効果量の計算方法と解釈

検定力分析

必要なサンプルサイズの決定方法

検定力分析では、事前に必要なサンプルサイズを決定します。検出したい効果量、目標とする検定力（一般的に0.8）、有意水準を指定することで、必要なサンプルサイズを算出することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の基本

t検定の実施手順



ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の実施手順

データの準備

検定に必要なデータの収集と整理方法

基礎統計量の確認

平均値、標準偏差基本的な統計量の確認方法

正規性の検証

データの正規性を確認する方法

データの準備では、分析の目的に応じて適切なデータを収集し整理します。例えば、独立した2標本のt検定の場合、2つの群のデータを明確に区別して記録し、外れ値や欠損値の処理も適切に行います。測定値の精度や単位の統一性にも注意を払う必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の実施手順

データの準備

検定に必要なデータの収集と整理方法

基礎統計量の確認

平均値、標準偏差基本的な統計量の確認方法

正規性の検証

データの正規性を確認する方法

基礎統計量の確認は、検定の前に必ず行います。各群の平均値、標準偏差、サンプルサイズなどの基本的な統計量を算出し、データの全体像を把握します。ここで大きな差異が見られれば、それが統計的に有意かどうかを検定で確認することになります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の実施手順

データの準備

検定に必要なデータの収集と整理方法

基礎統計量の確認

平均値、標準偏差基本的な統計量の確認方法

正規性の検証

データの正規性を確認する方法

正規性の検証には、ヒストグラムやQ-Qプロットなどの視覚的な方法と、シャピロ・ウィルク検定などの統計的検定を併用します。データが正規分布から大きく外れている場合は、対数変換などの変換を検討するか、ノンパラメトリック検定の使用を考慮します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の実施手順

等分散性の検証

2群間の分散の同質性を確認する方法

検定の実行

実際のt検定の実施手順とツールの使用方法

結果の解釈

検定結果の読み方と実務的な意味の解釈

等分散性の検証には、F検定やルビーン検定などを使用します。等分散性が成り立たない場合は、ウェルチのt検定など、等分散性を仮定しない方法を選択します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の実施手順

等分散性の検証

2群間の分散の同質性を確認する方法

検定の実行

実際のt検定の実施手順とツールの使用方法

結果の解釈

検定結果の読み方と実務的な意味の解釈

検定の実行では、統計ソフトウェアやプログラミング言語の統計パッケージを使用します。入力データの形式を確認し、適切なオプションを選択して実行します。両側検定か片側検定か、等分散を仮定するかどうかなど、適切な設定を行うことが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-21. t検定の実施方法

第2章 データ分析基礎

t検定の実施手順

等分散性の検証

2群間の分散の同質性を確認する方法

検定の実行

実際のt検定の実施手順とツールの使用方法

結果の解釈

検定結果の読み方と実務的な意味の解釈

結果の解釈では、p値、t値、自由度などの統計量を確認します。p値が有意水準より小さければ帰無仮説を棄却し、群間に有意な差があると判断します。ただし、統計的有意性だけでなく、効果量や実務的な重要性も考慮に入れて総合的に判断することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

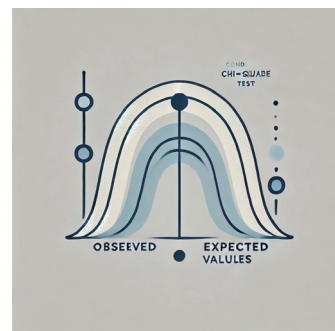
NEXT

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

カイ二乗検定の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

検定の種類

独立性の検定と適合度の検定の違いと用途

期待度数

理論的に期待される度数の計算方法

検定統計量

カイ二乗値の計算方法と自由度

カイ二乗検定には主に2つの種類があります。独立性の検定は、2つの変数間に関連があるかどうかを検証します。例えば、性別と商品の選好に関係があるかを調べる場合に使用します。適合度の検定は、観測された度数が理論的な分布に従うかどうかを検証します。例えば、サイコロの出目が均等に出ているかを確認する場合に使用します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

検定の種類

独立性の検定と適合度の検定の違いと用途

期待度数

理論的に期待される度数の計算方法

検定統計量

カイ二乗値の計算方法と自由度

期待度数は、帰無仮説が真である場合に期待される理論的な度数です。独立性の検定では、行の周辺度数と列の周辺度数の積を総数で割って計算します。例えば、性別と商品選択の独立性を検定する場合、男性全体における商品選択の比率が、女性でも同じように観察されるはずだという前提で期待度数を計算します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

検定の種類

独立性の検定と適合度の検定の違いと用途

期待度数

理論的に期待される度数の計算方法

検定統計量

カイ二乗値の計算方法と自由度

検定統計量であるカイ二乗値は、観測度数と期待度数の差の二乗を期待度数で割った値の合計として計算されます。自由度は、独立性の検定では $(\text{行数}-1) \times (\text{列数}-1)$ として計算されます。このカイ二乗値と自由度から、p値を求めることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

前提条件

期待度数の最小値検定の適用条件

効果量

Cramer's Vや ϕ 係数効果量の計算と解釈

サンプルサイズ

必要なサンプルサイズの決定方法

前提条件として、期待度数が5未満のセルが全体の20%を超えないことが推奨されます。また、いずれのセルでも期待度数が1を下回らないことが必要です。これらの条件を満たさない場合は、カテゴリの統合やフィッシャーの正確確率検定の使用を検討します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

前提条件

期待度数の最小値検定の適用条件

効果量

Cramer's Vや ϕ 係数効果量の計算と解釈

サンプルサイズ

必要なサンプルサイズの決定方法

効果量の指標としては、クロス表の大きさに応じて異なる指標が使用されます。2×2表の場合は ϕ 係数を、それ以外の場合はCramer's Vを使用するのが一般的です。これらの値は0から1の範囲をとり、値が大きいくほど関連が強いことを示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

前提条件

期待度数の最小値検定の適用条件

効果量

Cramer's Vや ϕ 係数効果量の計算と解釈

サンプルサイズ

必要なサンプルサイズの決定方法

サンプルサイズの決定には、検定力分析を用います。期待される効果量、目標とする検定力、有意水準を指定することで、必要なサンプルサイズを算出することができます。ただし、期待度数に関する制約も考慮に入れる必要があります。

ビッグデータ基礎

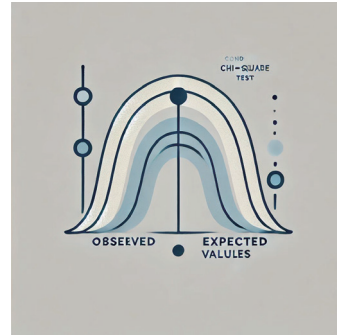
1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の基本

カイ二乗検定の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の実践

クロス集計表の作成

分析の基礎となるクロス集計表の作成方法

残差分析

調整済み残差の計算と解釈

検定の実行

具体的な検定手順とソフトウェアの使用方法

クロス集計表の作成は、カイ二乗検定の第一歩です。行と列に分析対象の変数を配置し、各セルの度数を集計します。例えば、性別×年齢層×商品選択のような多次元の分析も可能ですが、解釈が複雑になるため、通常は2次元の分析を行います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の実践

クロス集計表の作成

分析の基礎となるクロス集計表の作成方法

残差分析

調整済み残差の計算と解釈

検定の実行

具体的な検定手順とソフトウェアの使用方法

残差分析は、どのセルが全体的な傾向から特に大きく外れているかを調べる方法です。調整済み残差は標準正規分布に従うため、絶対値が1.96（5%水準）や2.58（1%水準）を超えるセルは、統計的に有意な偏りがあると判断できます。これにより、関連性の詳細なパターンを把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の実践

クロス集計表の作成

分析の基礎となるクロス集計表の作成方法

残差分析

調整済み残差の計算と解釈

検定の実行

具体的な検定手順とソフトウェアの使用方法

検定の実行では、統計ソフトウェアを使用します。データの入力形式を確認し、適切な検定方法（独立性の検定か適合度の検定か）を選択します。また、期待度数に関する前提条件もチェックします。条件を満たさない場合は、代替的な方法を検討します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の実践

結果の可視化

検定結果を効果的に図示する方法

多重比較

複数のカテゴリ間の比較方法

結果の報告

検定結果の適切な報告形式と解釈

結果の可視化では、モザイクプロットやバブルチャートなどを活用します。モザイクプロットでは、面積の大きさと度数を、色の濃さで残差を表現することができます。これにより、データの特徴をより直感的に理解することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の実践

結果の可視化

検定結果を効果的に図示する方法

多重比較

複数のカテゴリ間の比較方法

結果の報告

検定結果の適切な報告形式と解釈

多重比較は、3つ以上のカテゴリがある場合に、どのカテゴリ間に有意な差があるかを詳しく調べる方法です。ボンフェローニの補正など、適切な多重比較の手法を選択し、有意水準を調整して検定を行います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-22. カイ二乗検定の実施方法

第2章 データ分析基礎

カイ二乗検定の実践

結果の可視化

検定結果を効果的に図示する方法

多重比較

複数のカテゴリ間の比較方法

結果の報告

検定結果の適切な報告形式と解釈

結果の報告では、カイ二乗値、自由度、p値、効果量を明記します。また、クロス集計表や視覚化した図表も併せて提示することで、結果をより分かりやすく伝えることができます。特に実務での報告では、統計的な有意性だけでなく、実践的な意味についても言及することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

単回帰分析の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

回帰分析の目的

説明変数と目的変数の関係を数式化する意義

回帰直線

最小二乗法による回帰直線の求め方

回帰係数

傾きと切片の意味と解釈方法

回帰分析の目的は、説明変数 (X) と目的変数 (Y) の関係を数式として表現することです。例えば、広告費と売上高の関係、学習時間と試験得点の関係など、一方の変数から他方の変数を予測したり、その影響度を定量的に把握したりすることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

回帰分析の目的

説明変数と目的変数の関係を数式化する意義

回帰直線

最小二乗法による回帰直線の求め方

回帰係数

傾きと切片の意味と解釈方法

回帰直線は、最小二乗法によって求められます。これは、実測値と予測値の差（残差）の二乗和が最小となるような直線を求める方法です。例えば、散布図上の点から直線までの垂直距離の二乗和が最小となるように、直線の位置と傾きを決定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

回帰分析の目的

説明変数と目的変数の関係を数式化する意義

回帰直線

最小二乗法による回帰直線の求め方

回帰係数

傾きと切片の意味と解釈方法

回帰係数には、傾き（ β ）と切片（ α ）があります。傾きは説明変数が1単位変化したときの目的変数の変化量を表し、切片は説明変数が0のときの目的変数の値を示します。例えば、学習時間と得点の関係で傾きが2であれば、学習時間が1時間増えるごとに得点が2点上がると解釈できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

決定係数

モデルの当てはまりを示す R^2 の意味と計算方法

予測と推定

回帰式を用いた予測値の算出方法

残差分析

予測値と実測値の差である残差の分析方法

決定係数（ R^2 ）は、モデルの説明力を0から1の値で示します。例えば、 R^2 が0.7であれば、目的変数の変動の70%が説明変数によって説明できることを意味します。ただし、決定係数が高いことと、因果関係があることは必ずしも一致しないことに注意が必要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

決定係数

モデルの当てはまりを示す R^2 の意味と計算方法

予測と推定

回帰式を用いた予測値の算出方法

残差分析

予測値と実測値の差である残差の分析方法

予測と推定では、得られた回帰式を使って新しい値を予測します。例えば、ある広告費に対する売上高の予測値を計算したり、特定の効果を得るために必要な説明変数の値を逆算したりすることができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

決定係数

モデルの当てはまりを示す R^2 の意味と計算方法

予測と推定

回帰式を用いた予測値の算出方法

残差分析

予測値と実測値の差である残差の分析方法

残差分析では、予測値と実測値の差を詳しく調べます。残差の大きさやパターンを分析することで、モデルの適切性や改善点を把握することができます。例えば、残差が特定のパターンを示す場合、モデルの仮定が満たされていない可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の概念

単回帰分析の実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の実践

データの前処理

外れ値の処理や変数変換分析前の準備

モデルの適合度

回帰モデルの適合性を評価する方法

仮定の検証

線形性、正規性などの仮定の確認方法

データの前処理では、分析の精度を高めるための準備を行います。外れ値の検出と処理、変数の変換（対数変換など）、欠損値の処理などが含まれます。例えば、極端に大きな値が含まれる場合、それを除外するか適切な変換を行うかを検討します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の実践

データの前処理

外れ値の処理や変数変換分析前の準備

モデルの適合度

回帰モデルの適合性を評価する方法

仮定の検証

線形性、正規性などの仮定の確認方法

モデルの適合度は、複数の指標を用いて評価します。決定係数に加えて、F検定による回帰式の有意性検定、t検定による回帰係数の有意性検定などを行います。これらの結果から、モデルが統計的に意味のあるものかどうかを判断します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の実践

データの前処理

外れ値の処理や変数変換分析前の準備

モデルの適合度

回帰モデルの適合性を評価する方法

仮定の検証

線形性、正規性などの仮定の確認方法

仮定の検証は、回帰分析の前提条件が満たされているかを確認します。線形性（XとYの関係が直線的か）、正規性（残差が正規分布に従うか）、等分散性（残差の分散が一定か）などを、グラフや統計的検定を用いて確認します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の実践

信頼区間

回帰係数の信頼区間の計算と解釈

予測区間

新しい観測値の予測区間の算出方法

診断プロット

回帰診断のための各種プロットの作成と解釈

信頼区間は、回帰係数の不確実性を表現します。例えば、傾きの95%信頼区間が[1.5, 2.5]であれば、真の傾きがこの範囲に含まれる確率が95%であることを示します。この情報は、結果の信頼性を評価する際に重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の実践

信頼区間

回帰係数の信頼区間の計算と解釈

予測区間

新しい観測値の予測区間の算出方法

診断プロット

回帰診断のための各種プロットの作成と解釈

予測区間は、新しい観測値が含まれると期待される範囲を示します。これは信頼区間より広くなり、個々の観測値のばらつきも考慮に入れています。例えば、ある広告費に対する売上高の予測を行う際、予測値に加えてその不確実性の範囲も示すことができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-23. 単回帰分析の基礎

第2章 データ分析基礎

単回帰分析の実践

信頼区間

回帰係数の信頼区間の計算と解釈

予測区間

新しい観測値の予測区間の算出方法

診断プロット

回帰診断のための各種プロットの作成と解釈

診断プロットは、モデルの妥当性を視覚的に確認するためのツールです。残差プロット、Q-Qプロット、レバレッジプロットなど、様々なグラフを用いてモデルの問題点を診断します。これらのプロットから、モデルの改善点や制約事項を把握することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

重回帰分析の応用



ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

複数説明変数

複数の説明変数を用いた回帰モデルの基本概念

変数選択

モデルに含める変数の選択基準と方法

多重共線性

説明変数間の相関による問題と対処方法

複数説明変数を用いる重回帰分析では、目的変数に影響を与える複数の要因を同時に考慮することができます。例えば、住宅価格を予測する場合、広さ、築年数、駅からの距離など、複数の要因を組み合わせることで、より正確な予測が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

複数説明変数

複数の説明変数を用いた回帰モデルの基本概念

変数選択

モデルに含める変数の選択基準と方法

多重共線性

説明変数間の相関による問題と対処方法

変数選択では、モデルに含めるべき説明変数を適切に決定します。理論的な根拠、予備分析での相関関係、実務的な重要性などを考慮して選択します。ただし、むやみに変数を増やすと、モデルが複雑化し、解釈が難しくなる可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

複数説明変数

複数の説明変数を用いた回帰モデルの基本概念

変数選択

モデルに含める変数の選択基準と方法

多重共線性

説明変数間の相関による問題と対処方法

多重共線性は、説明変数間に強い相関がある場合に生じる問題です。例えば、身長と体重のように強く相関する変数を同時にモデルに含めると、個々の変数の効果を正確に推定することが困難になります。VIF（分散拡大要因）などの指標を用いて診断し、必要に応じて変数の選択や主成分分析などの対処を行います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

標準化係数

異なる尺度の変数を比較するための標準化

調整済み決定係数

モデルの複雑さを考慮した適合度指標

交互作用

説明変数間の相互作用の検討方法

標準化係数は、単位の異なる説明変数の影響力を比較するために使用します。例えば、「価格（円）」と「距離（km）」という異なる単位の変数の影響力を直接比較することができます。これにより、どの要因が目的変数により強く影響しているかを判断できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

標準化係数

異なる尺度の変数を比較するための標準化

調整済み決定係数

モデルの複雑さを考慮した適合度指標

交互作用

説明変数間の相互作用の検討方法

調整済み決定係数は、説明変数の数を考慮した適合度の指標です。通常の決定係数は変数を追加するほど大きくなる傾向がありますが、調整済み決定係数は不必要な変数の追加にペナルティを与えます。これにより、モデルの真の説明力をより適切に評価することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

標準化係数

異なる尺度の変数を比較するための標準化

調整済み決定係数

モデルの複雑さを考慮した適合度指標

交互作用

説明変数間の相互作用の検討方法

交互作用は、説明変数間の相乗効果を表します。例えば、気温と降水量が作物の生育に与える影響を考える場合、両者の組み合わせによる効果を考慮することで、より現実的なモデルを構築できます。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の基本

重回帰分析の応用



ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の応用

変数選択手法

ステップワイズ法などの自動変数選択手法

ダミー変数

カテゴリ変数の扱い方とダミー変数の作成方法

モデル診断

残差分析や影響度分析などの診断手法

変数選択手法には、ステップワイズ法、フォワード法、バックワード法などがあります。これらの手法は、統計的な基準に基づいて変数を自動的に選択します。例えば、ステップワイズ法では、変数を一つずつ追加・削除しながら、最適な組み合わせを探索します。ただし、機械的な選択に頼りすぎず、理論的な妥当性も考慮することが重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の応用

変数選択手法

ステップワイズ法などの自動変数選択手法

ダミー変数

カテゴリ変数の扱い方とダミー変数の作成方法

モデル診断

残差分析や影響度分析などの診断手法

ダミー変数は、性別や地域などのカテゴリ変数をモデルに組み込む際に使用します。例えば、「大人=1、子供=0」のように、カテゴリを0と1の数値に変換します。3つ以上のカテゴリがある場合は、基準カテゴリを決めて、それ以外のカテゴリごとにダミー変数を作成します。

ビッグデータ基礎 1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の応用

変数選択手法

ステップワイズ法などの自動変数選択手法

ダミー変数

カテゴリ変数の扱い方とダミー変数の作成方法

モデル診断

残差分析や影響度分析などの診断手法

モデル診断では、残差プロット、影響度分析、Cook's距離など、様々な手法を用いてモデルの適切性を評価します。異常値の検出、モデルの仮定の確認、影響力の大きいデータポイントの特定などを行い、必要に応じてモデルを修正します。

ビッグデータ基礎 1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の応用

予測精度評価

予測モデルとしての精度評価方法

過学習

モデルが複雑化することによる問題と対処法

結果の解釈

重回帰分析結果の実務的な解釈と報告方法

予測精度の評価では、訓練データとテストデータを分けて分析を行います。平均二乗誤差（MSE）や平均絶対誤差（MAE）などの指標を用いて、モデルの予測性能を定量的に評価します。交差検証を行うことで、より信頼性の高い評価が可能です。

ビッグデータ基礎 1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の応用

予測精度評価

予測モデルとしての精度評価方法

過学習

モデルが複雑化することによる問題と対処法

結果の解釈

重回帰分析結果の実務的な解釈と報告方法

過学習は、モデルが訓練データに過度に適合し、新しいデータに対する予測精度が低下する現象です。変数が多すぎる場合や、サンプルサイズが小さい場合に起こりやすくなります。正則化やモデルの単純化などの対処法を検討する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-24. 重回帰分析の基礎

第2章 データ分析基礎

重回帰分析の応用

予測精度評価

予測モデルとしての精度評価方法

過学習

モデルが複雑化することによる問題と対処法

結果の解釈

重回帰分析結果の実務的な解釈と報告方法

結果の解釈では、統計的な有意性だけでなく、実務的な重要性も考慮します。係数の符号と大きさ、信頼区間、予測値の精度など、様々な側面から結果を吟味し、わかりやすく報告することが重要です。また、モデルの限界や前提条件についても明確に説明する必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

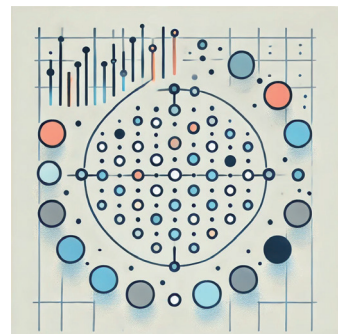


2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

クラスタリングの実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

クラスタリングの目的

データを類似性に基づいてグループ化する意義と活用場面

類似度と距離

データ間の類似性を測る様々な距離尺度

階層的クラスタリング

データを階層的にグループ化する手法の基本原則

クラスタリングの目的は、データを類似性に基づいて意味のあるグループに分類することです。例えば、顧客の購買パターンに基づく顧客セグメンテーション、商品の特徴に基づく商品グループ化、地域の特性に基づくエリア分類など、様々な場面で活用されます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

クラスタリングの目的

データを類似性に基づいてグループ化する意義と活用場面

類似度と距離

データ間の類似性を測る様々な距離尺度

階層的クラスタリング

データを階層的にグループ化する手法の基本原則

類似度と距離は、データ間の近さを測る指標です。数値データではユークリッド距離やマンハッタン距離、カテゴリデータではハミング距離、テキストデータではコサイン類似度など、データの種類や分析の目的に応じて適切な尺度を選択します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

クラスタリングの目的

データを類似性に基づいてグループ化する意義と活用場面

類似度と距離

データ間の類似性を測る様々な距離尺度

階層的クラスタリング

データを階層的にグループ化する手法の基本原則

階層的クラスタリングは、データを段階的にグループ化していく手法です。最初は各データポイントを個別のクラスタとし、近いものから順次結合していく方法（凝集型）と、全体から徐々に分割していく方法（分割型）があります。結果はデンドログラム（樹形図）として視覚化され、様々な段階でのグループ分けを確認できます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

非階層的クラスタリング

k-means法などの非階層的手法の特徴と手順

クラスタ数の決定

最適なクラスタ数を決定するための方法と基準

クラスタの評価

クラスタリング結果の品質を評価する方法

非階層的クラスタリングの代表的な手法であるk-means法は、あらかじめ指定した数のクラスタにデータを分類します。各クラスタの中心点を反復的に更新しながら、最適な分類を探索します。計算が高速で大規模データにも適用できる利点がありますが、初期値依存性という課題もあります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

非階層的クラスタリング

k-means法などの非階層的手法の特徴と手順

クラスタ数の決定

最適なクラスタ数を決定するための方法と基準

クラスタの評価

クラスタリング結果の品質を評価する方法

クラスタ数の決定は、分析の重要なステップです。エルボー法（クラスタ内の分散の変化を見る）、シルエット分析（クラスタの分離度を評価）、ギャップ統計量など、様々な方法を組み合わせて適切なクラスタ数を決定します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

非階層的クラスタリング

k-means法などの非階層的手法の特徴と手順

クラスタ数の決定

最適なクラスタ数を決定するための方法と基準

クラスタの評価

クラスタリング結果の品質を評価する方法

クラスタの評価では、クラスタ内の凝集度（どれだけ緊密にまとまっているか）とクラスタ間の分離度（どれだけ明確に分かれているか）を確認します。また、クラスタの大きさのバランスや、実務的な解釈可能性も重要な評価基準となります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの基本概念

クラスタリングの実践



ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの実践

データの前処理

標準化やスケーリング分析前の必要な処理

手法の選択

データの特性に応じた適切なクラスタリング手法の選択

結果の可視化

デンドログラムやヒートマップなどの可視化手法

データの前処理は、クラスタリングの成否を左右する重要なステップです。変数の標準化やスケーリングにより、異なる尺度の変数を公平に扱えるようにします。また、外れ値の処理や欠損値の補完、不要な変数の除外なども、結果の質に大きく影響します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの実践

データの前処理

標準化やスケーリング分析前の必要な処理

手法の選択

データの特性に応じた適切なクラスタリング手法の選択

結果の可視化

デンドログラムやヒートマップなどの可視化手法

手法の選択では、データの特性や分析の目的に応じて適切な方法を選びます。データサイズが大きい場合はk-means法、クラスターの階層構造を把握したい場合は階層的クラスタリング、ノイズの多いデータにはDBSCANなど、それぞれの手法の特徴を考慮して選択します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの実践

データの前処理

標準化やスケーリング分析前の必要な処理

手法の選択

データの特徴に応じた適切なクラスタリング手法の選択

結果の可視化

デンドログラムやヒートマップなどの可視化手法

結果の可視化は、クラスタリング結果の理解を助けます。階層的クラスタリングではデンドログラム、k-means法では散布図やヒートマップなどを用います。多次元データの場合は、主成分分析などで次元を縮約して可視化することも有効です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの実践

クラスタの解釈

形成されたクラスタの特徴を理解し解釈する方法

クラスタの活用

ビジネスへの活用方法と意思決定への応用

結果の検証

クラスタリング結果の妥当性を検証する方法

クラスタの解釈では、各クラスタの特徴を明確にします。クラスタごとの変数の平均値や分布を比較し、そのクラスタを特徴づける要素を特定します。例えば、顧客セグメントであれば「高額・低頻度購入層」「少額・高頻度購入層」といった具体的な特徴付けを行います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの実践

クラスタの解釈

形成されたクラスタの特徴を理解し解釈する方法

クラスタの活用

ビジネスへの活用方法と意思決定への応用

結果の検証

クラスタリング結果の妥当性を検証する方法

クラスタの活用では、分析結果を実務に落とし込みます。例えば、顧客セグメントごとにマーケティング施策を変える、商品グループごとに在庫管理方針を設定する、地域特性に応じて出店戦略を立てるなど、具体的なアクションにつなげます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-25. クラスタリング分析の基礎

第2章 データ分析基礎

クラスタリングの実践

クラスタの解釈

形成されたクラスタの特徴を理解し解釈する方法

クラスタの活用

ビジネスへの活用方法と意思決定への応用

結果の検証

クラスタリング結果の妥当性を検証する方法

結果の検証では、クラスタリングの安定性と妥当性を確認します。異なる初期値や手法での分析結果を比較したり、データの一部を使った検証を行ったりします。また、ドメイン知識や実務経験に照らして、結果が意味のある分類になっているかも重要な検証ポイントです。

ビッグデータ基礎 1 2 3 4 5 6 7

NEXT

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

主成分分析の応用



ビッグデータ基礎 1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

主成分分析の目的

高次元データの次元削減と情報圧縮の意義

主成分の概念

データの分散を最大化する軸としての主成分

固有値と固有ベクトル

主成分を求めるための数学的基礎

主成分分析の目的は、多数の変数を持つデータから、重要な情報を失うことなく、より少ない次元で表現することです。例えば、体格に関する10個の測定値から、「体の大きさ」「体型のバランス」といった2つの総合指標を作り出すことができます。これにより、データの解釈が容易になり、後続の分析も効率的に行えるようになります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

主成分分析の目的

高次元データの次元削減と情報圧縮の意義

主成分の概念

データの分散を最大化する軸としての主成分

固有値と固有ベクトル

主成分を求めるための数学的基礎

主成分は、データの持つ情報（分散）をできるだけ多く保持するように決定される新しい軸です。第1主成分は元のデータの分散を最大化する方向、第2主成分は第1主成分と直交し、残りの分散を最大化する方向というように、順次定義されていきます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

主成分分析の目的

高次元データの次元削減と情報圧縮の意義

主成分の概念

データの分散を最大化する軸としての主成分

固有値と固有ベクトル

主成分を求めるための数学的基礎

固有値と固有ベクトルは、主成分を数学的に定義する際の基礎となります。相関行列または分散共分散行列から計算される固有値は各主成分の分散の大きさを、固有ベクトルは主成分の方向を示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

寄与率

各主成分の説明力を示す寄与率の計算と解釈

スクリープロット

主成分の重要性を視覚化する方法

次元数の決定

保持する主成分の数を決定する基準

寄与率は、各主成分がデータ全体の分散をどの程度説明しているかを示す指標です。例えば、第1主成分の寄与率が60%であれば、元のデータの持つ情報の60%がこの主成分で説明できることを意味します。累積寄与率は、複数の主成分による説明力の合計を示します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

寄与率

各主成分の説明力を示す寄与率の計算と解釈

スクリープロット

主成分の重要性を視覚化する方法

次元数の決定

保持する主成分の数を決定する基準

スクリープロットは、主成分の重要性を視覚的に表現するグラフです。横軸に主成分の番号、縦軸に固有値をプロットし、どの主成分まで重要かを判断する材料とします。グラフの傾きが急激に緩やかになる点（肘）が、適切な主成分数の目安となります。

ビッグデータ基礎 1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

寄与率

各主成分の説明力を示す寄与率の計算と解釈

スクリープロット

主成分の重要性を視覚化する方法

次元数の決定

保持する主成分の数を決定する基準

次元数の決定では、累積寄与率が80%程度に達する点や、固有値が1以上の主成分を採用するなど、いくつかの基準があります。ただし、解釈のしやすさや実用性も考慮して、最終的な判断を行う必要があります。

ビッグデータ基礎 1 2 3 4 5 6 7

NEXT

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の基本

主成分分析の応用



ビッグデータ基礎 1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の応用

データの標準化

分析前の変数の標準化とその重要性

負荷量の解釈

各変数の主成分への寄与度の解釈方法

スコアの計算

個々のデータの主成分空間での位置付け

データの標準化は、主成分分析の前処理として重要です。変数の単位や分散が異なると、不当に大きな値を持つ変数が結果を支配してしまう可能性があります。そこで、各変数を平均0、分散1に標準化することで、公平な分析が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の応用

データの標準化

分析前の変数の標準化とその重要性

負荷量の解釈

各変数の主成分への寄与度の解釈方法

スコアの計算

個々のデータの主成分空間での位置付け

負荷量は、各変数が主成分にどの程度寄与しているかを示す指標です。例えば、第1主成分に対して身長や体重の負荷量が多い場合、この主成分は「体の大きさ」を表していると解釈できます。負荷量の行列（因子パターン）を見ることで、各主成分の意味を理解することができます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の応用

データの標準化

分析前の変数の標準化とその重要性

負荷量の解釈

各変数の主成分への寄与度の解釈方法

スコアの計算

個々のデータの主成分空間での位置付け

主成分スコアは、各データポイントの主成分空間での座標値です。例えば、第1主成分と第2主成分のスコアをプロットすることで、データの分布や特徴的なグループを視覚的に把握できます。また、このスコアを後続の分析に使用することもできます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の応用

可視化手法

主成分分析結果の効果的な可視化方法

変数の選択

分析に用いる変数の選択基準と方法

実務での活用

主成分分析結果のビジネスでの活用方法

可視化手法としては、バイプロットがよく用いられます。これは、データポイントの分布と変数の負荷量を同一の図に表示するもので、データの構造を包括的に理解するのに役立ちます。また、主成分スコアの散布図や、負荷量のベクトル図なども有用です。

ビッグデータ基礎 1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の応用

可視化手法

主成分分析結果の効果的な可視化方法

変数の選択

分析に用いる変数の選択基準と方法

実務での活用

主成分分析結果のビジネスでの活用方法

変数の選択では、分析の目的と変数間の関係を考慮します。強い相関のある変数群は、代表的な変数のみを選択することも検討します。また、主成分の解釈可能性を考慮して、関連性の強い変数群を選択することも重要です。

ビッグデータ基礎 1 2 3 4 5 6 7

2-26. 主成分分析の基礎

第2章 データ分析基礎

主成分分析の応用

可視化手法

主成分分析結果の効果的な可視化方法

変数の選択

分析に用いる変数の選択基準と方法

実務での活用

主成分分析結果のビジネスでの活用方法

実務での活用例として、商品開発における特徴の要約、顧客データの次元削減、画像認識の特徴抽出などがあります。また、多数の経済指標から景気動向を把握したり、多様な評価項目から総合的な評価指標を作成したりする場面でも活用されています。

ビッグデータ基礎 1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

解釈の実践とコミュニケーション



ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

統計的意味

p値や信頼区間統計指標の正しい解釈方法

実務的意味

統計的結果をビジネス文脈で解釈する方法

因果関係の考察

相関と因果関係の区別、および因果推論の基本

統計的意味の理解は、分析結果を正しく解釈する基本となります。p値が0.05より小さいことは、観察された結果が偶然では説明しにくいことを示しますが、効果の大きさを示すものではありません。また、信頼区間は推定値の不確実性の範囲を示し、例えば95%信頼区間は、真の値がその範囲に含まれる確率が95%であることを意味します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

統計的意味

p値や信頼区間統計指標の正しい解釈方法

実務的意味

統計的結果をビジネス文脈で解釈する方法

因果関係の考察

相関と因果関係の区別、および因果推論の基本

実務的意味の解釈では、統計的な有意性だけでなく、ビジネスにおける重要性を考慮します。例えば、新施策による売上の1%の増加が統計的に有意であっても、実施コストを考えると実務的には価値が低いかもしれません。逆に、5%の増加が統計的に有意でなくても、リスクが低ければ試してみる価値があるかもしれません。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

統計的意味

p値や信頼区間統計指標の正しい解釈方法

実務的意味

統計的結果をビジネス文脈で解釈する方法

因果関係の考察

相関と因果関係の区別、および因果推論の基本

因果関係の考察では、相関関係と因果関係を慎重に区別します。例えば、アイスクリームの売上と熱中症の発生件数に相関があっても、それは気温という第三の要因の影響である可能性が高いです。因果関係を主張するためには、実験的なデザインや、交絡要因の制御が必要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

限界の理解

分析結果の制約事項や適用範囲

モデルの評価

分析モデルの性能と信頼性の評価方法

バイアスの検討

結果に影響を与える可能性のあるバイアス

分析の限界を理解することも重要です。サンプルの代表性、データの収集期間、モデルの前提条件など、結果の一般化や適用に関する制約を明確にする必要があります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

限界の理解

分析結果の制約事項や適用範囲

モデルの評価

分析モデルの性能と信頼性の評価方法

バイアスの検討

結果に影響を与える可能性のあるバイアス

モデルの評価では、適合度や予測精度だけでなく、モデルの安定性や頑健性も確認します。例えば、データの一部を変更した場合や、異なる期間のデータを使用した場合でも、同様の結果が得られるかを検証します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

限界の理解

分析結果の制約事項や適用範囲

モデルの評価

分析モデルの性能と信頼性の評価方法

バイアスの検討

結果に影響を与える可能性のあるバイアス

バイアスの検討では、サンプリングバイアス、測定バイアス、生存バイアスなど、結果に影響を与える可能性のある様々なバイアスを考慮します。例えば、アンケート調査で回答者に特定の傾向がある場合、結果が母集団を正しく代表していない可能性があります。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT



2-27. データ分析結果の解釈方法

第2章 データ分析基礎

分析結果の基本的解釈

解釈の実践とコミュニケーション



ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

解釈の実践とコミュニケーション

文脈化

分析結果を事業環境や目的に関連付ける方法

重要度の評価

発見事項の重要性を評価する基準

シナリオ分析

異なる解釈の可能性を検討する方法

文脈化では、分析結果を事業環境や目的に関連付けます。例えば、顧客セグメント分析の結果を、現在の市場動向や競合状況と結びつけて解釈することで、より実践的な意味を見出すことができます。また、組織の戦略目標や過去の取り組みとの関連も考慮します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

解釈の実践とコミュニケーション

文脈化

分析結果を事業環境や目的に関連付ける方法

重要度の評価

発見事項の重要性を評価する基準

シナリオ分析

異なる解釈の可能性を検討する方法

重要度の評価では、実現可能性、費用対効果、リスク、時間的制約などの観点から、発見事項の優先順位を付けます。例えば、すぐに実行できる小さな改善と、大きな投資が必要な構造的な改革を区別して評価します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

解釈の実践とコミュニケーション

文脈化

分析結果を事業環境や目的に関連付ける方法

重要度の評価

発見事項の重要性を評価する基準

シナリオ分析

異なる解釈の可能性を検討する方法

シナリオ分析では、データが示す結果について複数の解釈の可能性を検討します。例えば、売上の減少が見られた場合、市場環境の変化、競合の影響、内部要因など、様々な角度から原因を考察します。これにより、より包括的な理解と対応策の検討が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

解釈の実践とコミュニケーション

示唆の導出

分析結果から実践的な示唆を導き出す方法

専門家との対話

統計の専門家との効果的な対話方法

意思決定への連携

解釈結果を意思決定プロセスに組み込む方法

示唆の導出では、分析結果から具体的なアクションにつながる洞察を引き出します。「何が」「なぜ起きているか」という事実の理解から、「どうすべきか」という行動の提案まで、論理的につなげていきます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

解釈の実践とコミュニケーション

示唆の導出

分析結果から実践的な示唆を導き出す方法

専門家との対話

統計の専門家との効果的な対話方法

意思決定への連携

解釈結果を意思決定プロセスに組み込む方法

専門家との対話では、統計の専門家と事業部門のコミュニケーションを促進します。専門用語を平易な言葉で説明したり、ビジネス上の課題を統計的な問いに翻訳したりする橋渡しの役割が重要です。

ビッグデータ基礎

1 2 3 4 5 6 7

2-27. データ分析結果の解釈方法

第2章 データ分析基礎

解釈の実践とコミュニケーション

示唆の導出

分析結果から実践的な示唆を導き出す方法

専門家との対話

統計の専門家との効果的な対話方法

意思決定への連携

解釈結果を意思決定プロセスに組み込む方法

意思決定への連携では、分析結果を実際の意思決定プロセスに組み込みます。定量的な分析結果と定性的な判断を適切にバランスさせ、より良い意思決定につなげます。また、決定後のモニタリングと検証のプロセスも設計します。

ビッグデータ基礎

1 2 3 4 5 6 7

NEXT

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

分析結果の報告



ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

可視化の原則

効果的なデータ可視化の基本原則と重要性

グラフの選択

データの種類や目的に応じた適切なグラフ形式の選択

インタラクティブ表示

動的な可視化ツールの活用方法と利点

可視化の基本原則は、「シンプルさ」「正確さ」「目的適合性」です。必要な情報を過不足なく伝え、直感的な理解を促すことが重要です。例えば、複雑なデータも、適切な可視化により一目で傾向を把握できるようになります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

可視化の原則

効果的なデータ可視化の基本原則と重要性

グラフの選択

データの種類や目的に応じた適切なグラフ形式の選択

インタラクティブ表示

動的な可視化ツールの活用方法と利点

グラフの選択は、データの性質と伝えたいメッセージに応じて行います。時系列データには折れ線グラフ、構成比には円グラフや積み上げ棒グラフ、分布の形状にはヒストグラムというように、データの特性に合わせて最適な形式を選びます。また、比較や相関、順位など、何を強調したいかによっても選択が変わってきます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

可視化の原則

効果的なデータ可視化の基本原則と重要性

グラフの選択

データの種類や目的に応じた適切なグラフ形式の選択

インタラクティブ表示

動的な可視化ツールの活用方法と利点

インタラクティブ表示は、データの探索的な分析に特に有効です。ズームイン・アウト、フィルタリング、ドリルダウンなどの機能により、利用者が自由に視点を変えてデータを探索できます。例えば、ダッシュボードツールを使用することで、リアルタイムでのデータ更新や対話的な分析が可能になります。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

デザイン要素

色使い、レイアウト、ラベルなどの視覚的要素

情報の階層化

重要度に応じた情報の構造化と表示方法

誤解の防止

誤解を招きやすい表現とその回避方法

デザイン要素は、視覚的な情報伝達の質を左右します。色使いは、色覚多様性に配慮しつつ、情報の優先度や関係性を表現します。レイアウトは、情報の流れを自然に誘導するように配置します。ラベルやタイトルは、簡潔かつ明確に情報を補充します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

デザイン要素

色使い、レイアウト、ラベルなどの視覚的要素

情報の階層化

重要度に応じた情報の構造化と表示方法

誤解の防止

誤解を招きやすい表現とその回避方法

情報の階層化は、複雑なデータを理解しやすく提示するために重要です。最も重要な情報を目立つ位置に配置し、詳細情報は必要に応じて参照できるように構造化します。例えば、概要→詳細、全体→部分という階層で情報を整理します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

デザイン要素

色使い、レイアウト、ラベルなどの視覚的要素

情報の階層化

重要度に応じた情報の構造化と表示方法

誤解の防止

誤解を招きやすい表現とその回避方法

誤解を防ぐためには、いくつかの注意点があります。軸の切片を操作して差を誇張しない、3D効果で値の比較を歪めない、適切な尺度と単位を使用するなど、データの正確な理解を妨げない表現を心がけます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の可視化

分析結果の報告



ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の報告

報告書の構成

効果的なレポート構成と必要な要素

ストーリー展開

論理的なストーリー展開の組み立て方

要約の作成

エグゼクティブサマリーの効果的な作成方法

報告書の構成は、「背景・目的」「方法」「結果」「考察」「提言」という基本的な流れに従います。各セクションは明確に区分され、読み手が必要な情報に素早くアクセスできるようにします。また、図表のナンバリングや参照の仕方も統一的なルールに従います。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の報告

報告書の構成

効果的なレポート構成と必要な要素

ストーリー展開

論理的なストーリー展開の組み立て方

要約の作成

エグゼクティブサマリーの効果的な作成方法

ストーリー展開では、論理的な流れを意識します。何が課題で、なぜその分析が必要だったのか、どのような方法で分析し、何が分かったのか、そしてそれがなぜ重要なのか、という流れで説明を組み立てます。読み手を自然な流れで結論に導くことを心がけます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の報告

報告書の構成

効果的なレポート構成と必要な要素

ストーリー展開

論理的なストーリー展開の組み立て方

要約の作成

エグゼクティブサマリーの効果的な作成方法

エグゼクティブサマリーは、報告書の要点を簡潔にまとめます。問題設定、主要な発見事項、重要な示唆、推奨されるアクションを、1~2ページに凝縮して記載します。忙しい意思決定者でも、報告書の本質を素早く把握できるようにします。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の報告

プレゼンテーション

分析結果を効果的に発表する方法

質疑対応

想定される質問への準備と回答方法

フォローアップ

報告後のアクションプランと進捗管理

プレゼンテーションでは、視覚的な要素と口頭での説明を効果的に組み合わせます。スライドは要点を簡潔に示し、詳細は口頭で補完します。また、聴衆の知識レベルや関心に合わせて、説明の深さや専門用語の使用を調整します。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

第2章 データ分析基礎

分析結果の報告

プレゼンテーション

分析結果を効果的に発表する方法

質疑対応

想定される質問への準備と回答方法

フォローアップ

報告後のアクションプランと進捗管理

質疑対応の準備では、想定される質問とその回答を事前に用意します。特に、分析の前提条件、方法の妥当性、結果の解釈、実務への適用可能性などについての質問を想定し、根拠となるデータや補足資料を準備しておきます。

ビッグデータ基礎

1 2 3 4 5 6 7

2-28. 分析結果の可視化と報告

分析結果の報告

プレゼンテーション

分析結果を効果的に発表する方法

質疑対応

想定される質問への準備と回答方法

フォローアップ

報告後のアクションプランと進捗管理

フォローアップでは、報告内容に基づくアクションプランを具体化します。誰が、何を、いつまでに実行するのか、どのように進捗を管理するのか、期待される成果をどう測定するのかを明確にします。また、定期的なレビューの機会を設定し、施策の効果検証も計画します。

令和6年度「地方やデジタル分野における専修学校理系転換等推進事業」
情報成長分野の教育プログラム整備と教員育成による学科転換・新設推進事業

ビッグデータ教材

－ 第2章 －

資料

令和7年2月

一般社団法人全国専門学校情報教育協会
〒164-0003 東京都中野区東中野 1-57-8 辻沢ビル3F
電話：03-5332-5081 FAX.03-5332-5083

●本書の内容を無断で転記、掲載することは禁じます。