

データ活用型 プロジェクトのマネジメント

日下部 貴彦【著】

コロナ社

ま え が き

本書は、データサイエンティストや研究者等が実施するデータ分析の領域だけでなく、これを包含するプロジェクト全体に焦点を当て、データ分析・活用を行う際に必要な基礎知識を整理したものである。昨今、産官学のさまざまな組織がさまざまなプロジェクトでデータを収集し、蓄積することが活発に行われており、大学での学術研究だけでなく企業等でのデータ活用の機会も一般化しつつある。一方で、データ活用の際に留意すべき個人情報保護法をはじめデータ利用に関わる制度も整備されてきており、データそのものに対する技術的視点以外にもデータ活用に必要な知識の領域が広がってきている。このような知識が不十分である場合には、コンプライアンス違反や手戻り等のプロジェクト遂行上のリスクに直結し、プロジェクトの成功の妨げとなることが想定される。このような背景から本書では、データ分析や活用を行うプロジェクトを成功に導くためにメンバーが共通的に持つておくべき基礎知識を網羅的に取り扱うことに主眼をおいた。

これまで筆者は、東京大学空間情報科学研究センターでの専任講師・准教授の在任中に、人の流れをはじめとした都市・交通に関わるデータの収集やこれらのデータを利用した研究を実施してきた。この際には、都市・交通分野のデータがさまざまな企業や行政組織によって分散的に収集されていることに起因したデータ入手手続きの煩雑さやデータクレンジング等の技術的課題等、多くの課題に直面してきた。その後、2021年4月より、「阪急阪神ホールディングス（株）の本社に『データ分析ラボ』を創設し、DXプロジェクトの推進を支援する。」という企業内のプロジェクトのもと、ラボ組織を立ち上げ、チームをリードし、データ分析を業務として確立することをミッションとしたディレクターとして活動している。このプロジェクトの当初は、筆者の大学でのデータ分析の経験を活かしてデータサイエンティストを育成することが、ディレクターの業務として期待されていた。しかし、このための環境を整えるため

には、データの収集や管理の手順書や業務フロー等を整備し、データ基盤の企画・運営、データ分析を用いた事業部門へのコンサルティング等を実施する必要があった。このように、企業内でデータサイエンスの知見を活かす環境を構成するための一連の業務を担ったが、このために必要な知識は既往の資料等では整理されたものが少なく、知識の習得やステークホルダーへの展開のために、自ら体系的な整理を行うことが必要であった。また、実際のプロジェクトを推進する過程では、ラボ内でのデータ分析に費やす時間と比べて、データの契約や実装、ステークホルダーへの説明に必要な時間が多く、プロジェクトの構成上は、データサイエンティストが分析を行う工程と同等にこれらの周辺にあるプロセスが重要なことを改めて実感した。

そこで本書では、筆者の上記のような経験で得た知見を中心に、データ活用を行うために必要な一連のプロセスに着目し、実践に必要な知識を整理した。特に、データ分析を担当するデータサイエンティストだけでなく、プロジェクトオーナー、プロジェクトマネージャー、チームリーダー等のプロジェクトを統括する方にも読んでいただくことで、データ活用・分析を伴うプロジェクトの特性や進行方法等についての共通的な知識を形成し、プロジェクトのより効率的な遂行やより確実な成功に寄与したいと考えて執筆を行っている。なお、本書で示しているプロジェクト構成等の例は、企業でのプロジェクトを中心としているが、プロジェクトの進行については筆者が経験した大学での研究室運営からもヒントを得ており、研究室の主宰者の方などにも是非読んでいただきたい。

本書の執筆にあたっては、山本隆弘氏をはじめ阪急阪神ホールディングス(株)グループ開発室 DX プロジェクト推進部およびデータ分析ラボのメンバー、関連する事業会社との取り組みを通じて得た知識・経験が大きな糧となっている。また、(株)ソーシャル・デザイナーズ・ベースの山田菊子氏に、執筆にあたってのきめ細やかなアドバイスをいただいた。親愛なる家族、活動を共にした皆様、コロナ社の皆様をはじめ執筆を支えてくださった多くの皆様に謝意を表します。

2023年11月

日下部 貴彦

目 次

1 章 序

| | |
|-------------------------------|---|
| 1.1 データサイエンティストを取り巻く状況とプロジェクト | 2 |
| 1.2 従来型プロジェクトとデータ活用型プロジェクトの定義 | 5 |
| 1.3 本書の目的 | 9 |

2 章 データ利用の類型と必要なスキルセット

| | |
|----------------------------|----|
| 2.1 データ利用 | 10 |
| 2.1.1 データの類型 | 10 |
| 2.1.2 データの一次利用 | 12 |
| 2.1.3 データの二次利用 | 13 |
| 2.1.4 新たなデータ活用 | 13 |
| 2.1.5 データの二次利用の利点と活用に向けた方策 | 14 |
| 2.2 データ活用に求められるスキルセット | 15 |
| 2.3 プロジェクトチームに必要なスキルセットの定義 | 16 |

3 章 データ活用に関連する法令

| | |
|---------------------------------------|----|
| 3.1 個人情報保護法 | 22 |
| 3.2 仮名加工情報と匿名加工情報 | 26 |
| 3.3 データ活用・分析目的での個人情報を含むデータ提供・受領者側の留意点 | 28 |

| | |
|--------------------|----|
| 3.4 データ提供・受領に関わる契約 | 31 |
| 3.5 データの管理 | 37 |

4章 データ活用型プロジェクト

| | |
|--|----|
| 4.1 データ活用型プロジェクトの定義 | 39 |
| 4.1.1 ステークホルダーの定義 | 40 |
| 4.1.2 ロールの定義と会議体 | 41 |
| 4.2 研究機関と企業のデータ活用型プロジェクトの違い | 44 |
| 4.3 データ活用型プロジェクトの体制 | 46 |
| 4.4 データ活用型プロジェクトの推進フェーズ | 51 |
| 4.4.1 現状理解・問題把握フェーズ | 52 |
| 4.4.2 課題設定フェーズ | 56 |
| 4.4.3 施策実施・検証フェーズ | 57 |
| 4.5 データ活用型プロジェクトの推進フェーズと 従来の改善のプロセスとの関係 | 58 |
| 4.6 データ活用型プロジェクトのナレッジ管理 | 62 |

5章 データ活用型プロジェクトのマネジメント

| | |
|----------------------------|----|
| 5.1 各フェーズでのプロジェクト推進方法 | 66 |
| 5.2 ウォーターフォール型プロジェクトの管理の基礎 | 69 |
| 5.2.1 進捗管理 | 71 |
| 5.2.2 課題管理 | 73 |
| 5.2.3 変更管理 | 74 |
| 5.3 アジャイル型プロジェクトの管理の基礎 | 75 |
| 5.4 データラボの運営 | 80 |
| 5.4.1 デイリースタンドイング | 80 |

| | | |
|-------|-----------------------|-----|
| 5.4.2 | イタレーションミーティング | 81 |
| 5.4.3 | スプリントレビュー | 81 |
| 5.5 | 現状理解・問題把握フェーズの進行例 | 83 |
| 5.5.1 | プロジェクトのキックオフと進行手順 | 83 |
| 5.5.2 | データの仕様確認と入手・実装, 分析の実施 | 88 |
| 5.5.3 | 事業担当者に対するヒアリングとレポート | 90 |
| 5.5.4 | 成果物 | 97 |
| 5.6 | 課題設定フェーズの進行例 | 97 |
| 5.6.1 | 課題設定フェーズのキックオフ | 98 |
| 5.6.2 | 課題設定フェーズのレポート | 99 |
| 5.6.3 | 成果物 | 100 |
| 5.7 | 施策実施フェーズの進行例 | 101 |
| 5.7.1 | 施策実施フェーズのレポート | 103 |
| 5.7.2 | 成果物 | 103 |

6章 データ分析の実践

| | | |
|-------|-------------------------|-----|
| 6.1 | データ活用のための環境構築 | 105 |
| 6.1.1 | データレイク, データウェアハウス | 107 |
| 6.1.2 | ETLによるデータマートへのデータ取り込み | 108 |
| 6.1.3 | データマートの実装 | 108 |
| 6.1.4 | データ分析環境 | 110 |
| 6.1.5 | 分析環境の実装 | 112 |
| 6.2 | データの取り込みと基礎分析 | 113 |
| 6.2.1 | データ受領時のファイルの同一性チェック | 114 |
| 6.2.2 | データの変換 | 115 |
| 6.2.3 | データマートの実装のためのデータベースデザイン | 116 |
| 6.2.4 | データの検証 | 118 |
| 6.3 | 現状理解・問題把握の分析 | 120 |
| 6.4 | 検証・評価のための分析 | 123 |

vi 目 次

6.4.1 課題解決方法の提案時の検証 123

6.4.2 モデルの評価 123

6.4.3 モデル評価のための指標 126

6.5 施策の検証 128

6.5.1 A/Bテスト 128

6.5.2 DID 131

6.5.3 Causal Impact 132

7章 集計・可視化の実践例

7.1 分析データと環境 134

7.2 ydata-profiling を用いた Exploratory Data Analysis 136

7.3 集計と可視化 139

7.3.1 パッケージとファイルの読み込み 139

7.3.2 分布の可視化 140

7.3.3 出発場所別トリップ数の集計と可視化 144

7.3.4 出発場所別旅行時間の集計と可視化 147

7.3.5 トリップの遷移と可視化 149

付録：データ活用に関連した法令に関する条文

付1 民法 153

付2 令和3年改正個人情報保護法（平成十五年法律第五十七号個人情報の保護に関する法律） 154

付3 個人情報の保護に関する法律施行令（平成十五年政令第五百七号） 164

付4 平成二十八年個人情報保護委員会規則第三号（個人情報の保護に関する法律施行規則） 167

付5 不正競争防止法（平成五年法律第四十七号） 171

付6 デジタル社会形成基本法（令和三年法律第三十五号） 174

付7 統計法（平成十九年法律第五十三号） 174

引用・参考文献 177

索引 182

1章 序

データサイエンス^{1)†}のコミュニティ（例えば Google が運営している Kaggle²⁾）やさまざまなオープンデータ³⁾の取り組みの広がりにより、多様な分野で収集される実際のデータを用いて、データサイエンスの経験を積み、技術力を研鑽する環境が整ってきている。このような環境により、データサイエンティストがさまざまな手法の試行錯誤を繰り返すことが可能になり、手法の高度化や精度向上などを行うことで、データ活用の領域が加速度的に広がってきている。

データ活用が行われる現場は、研究者の自由な発想に基づいて行われる学術研究での高度な手法の開発や実装を目指したものだけでなく、企業でのマーケティング、製品開発、経営企画等や、行政組織での施策設計、EBPM (evidence based policy making)⁴⁾等をはじめとして多方面に展開されることが期待されている。このようなデータ活用を前提としたプロジェクトの展開では、データ分析等を実施するデータサイエンティストだけでなく、プロジェクトの企画者や統括者、データの提供者等を含むプロジェクトメンバーが、データ活用を行うプロジェクトの特性を理解し、実行することがプロジェクトを成功に導く鍵となる。このような背景から、本書は、データ活用・分析の技術的な内容だけでなく、法令や契約、プロジェクト管理、情報技術等を含む全体像を俯瞰し、データ活用を行うプロジェクトメンバーが円滑にプロジェクトを推進するうえで必要となる知識をまとめたものである。

† 肩付き数字は、巻末の引用・参考文献番号を表す。

1.1 データサイエンティストを取り巻く状況とプロジェクト

初めに述べたように、さまざまな分野でデータ活用を加速させるうえでのデータサイエンティストの役割は、特定の AI (artificial intelligence) や機械学習等の手法の高度化や精度向上を行ったり、提供されたデータに対して数理的な手法を適用したうえで解釈を得ることだけでは、満たされなくなっている状況がある。つまり、データ活用の目的を達成するために的確なデータを収集したうえで、収集データをはじめとしたさまざまな情報を用いて課題設定を行い、その課題に対して適切な手法を選択し、データ解析とその活用を通じて目的を達成する一連のプロセスのなかでのさまざまな役割が重要となっている。

一連のデータ活用のプロセスを実行するためには、これまでデータサイエンスのコミュニティで利用できたオープンデータだけでなく、各専門分野の研究や企業での研究開発に必要なクローズドなデータセットを対象とすることも必要となる。このためには、分析手法の側面だけでなくデータの契約やプロジェクト管理、要件定義等のプロセスの側面からもデータを適切に取り扱うための取り組みが必要となる。特に都市や交通、マーケティング等の活動を対象として、民間企業をはじめとした多くの事業者等が収集するビジネスデータを取り扱うためには、個人情報の取り扱いをはじめとしたデータに対する制度や考え方を、データサイエンティストだけでなく、データの収集や契約、活用に携わるプロジェクトメンバーが共通に理解したうえでプロジェクトを推進する必要がある。さらに、プロジェクトのアウトプットを用いて意図したデータ活用が実施できるよう、適切なスコープの設定やプロジェクトの管理も重要な要素となる。(図 1.1)。

一般に、データを活用するプロジェクトでは

- データ活用を含むプロジェクトスコープの明確化
- 新たなデータの収集や第三者が収集したデータの受領を含むデータ取得

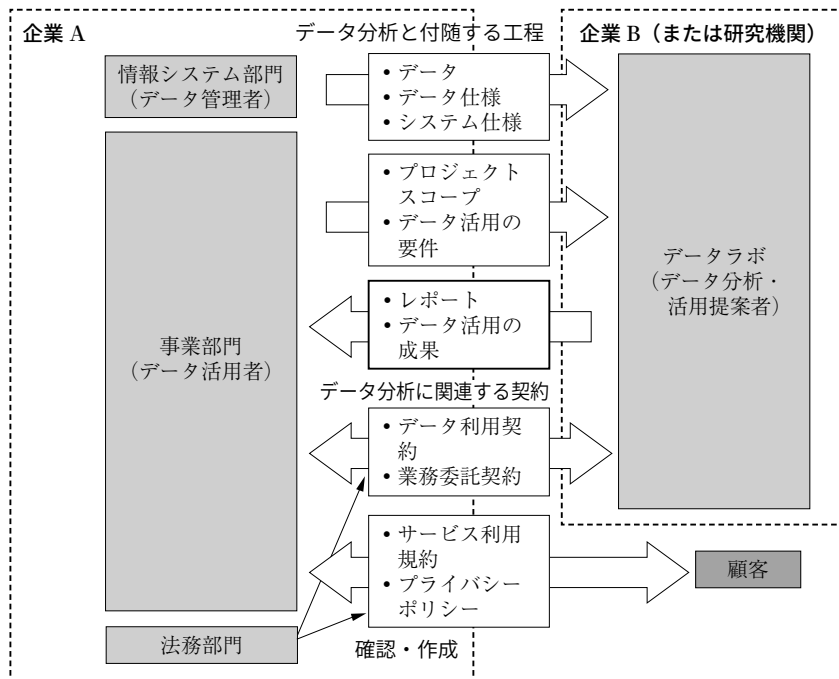


図 1.1 データ活用のためのプロジェクト概要

- 個人情報や機密情報へのアクセス制限等適切なデータの取り扱い環境の整備
- 分析可能なデータへの変換・加工

等の上流側のプロセスを経て初めて、データサイエンティストが分析手法の選択や構築、分析や検証、レポート等が実施できる。さらに、分析結果の活用や社会実装のためには、情報システムへの実装や関係するステークホルダーへの説明、費用や収益性等のビジネス的な観点での判断も必要となる。つまり、データ活用を行うのはデータサイエンティストではなく、その上流プロセスを行うステークホルダーであり、プロジェクトの中でデータサイエンティストが担う役割はデータ分析やモデル構築を実施する等プロジェクトの中では限られた部分である。したがって、それ以外の業務を担うシステムエンジニアや法務担当者、ビジネスや研究の企画担当者等の各プロセスのメンバーも

4 1章 序

データを活用したプロジェクトが具備すべき条件を体系的に理解する必要がある。図 1.2 および図 1.3 では、上記に示したデータ活用の実装までの流れの一例と、このようなプロジェクトに関わるステークホルダーの関係を示している。これらの図に示されるように、データ活用に異なる役割を持った多数のプロジェクトメンバーが関わる。このことから、各プロセスに関わるメンバーの

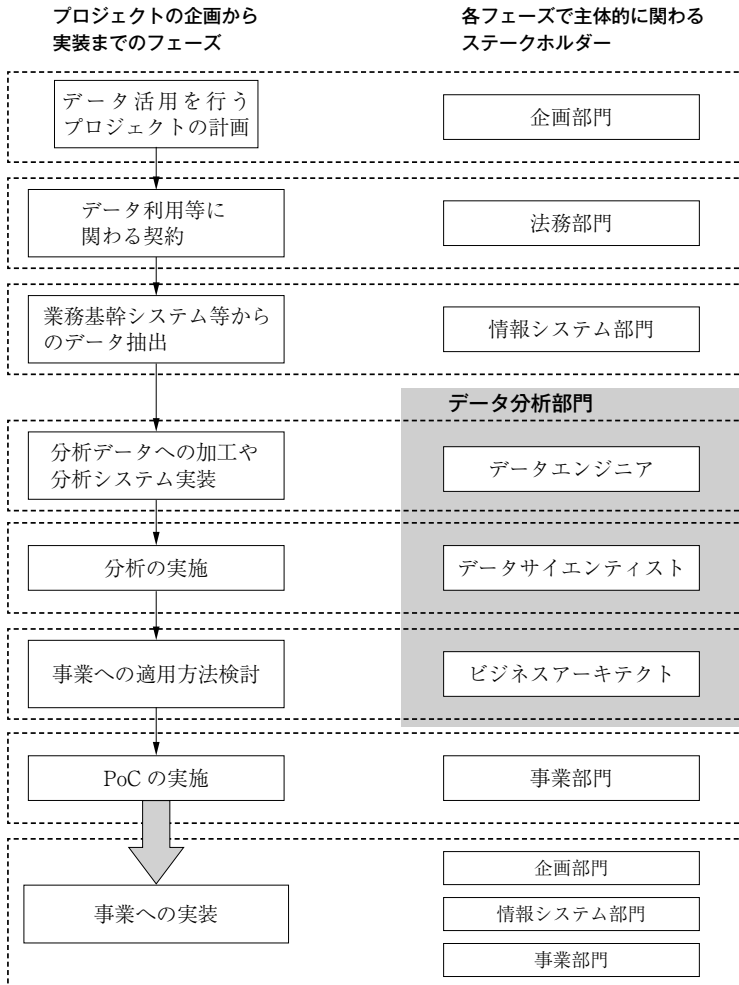


図 1.2 データ活用の実装までの流れとステークホルダーの例

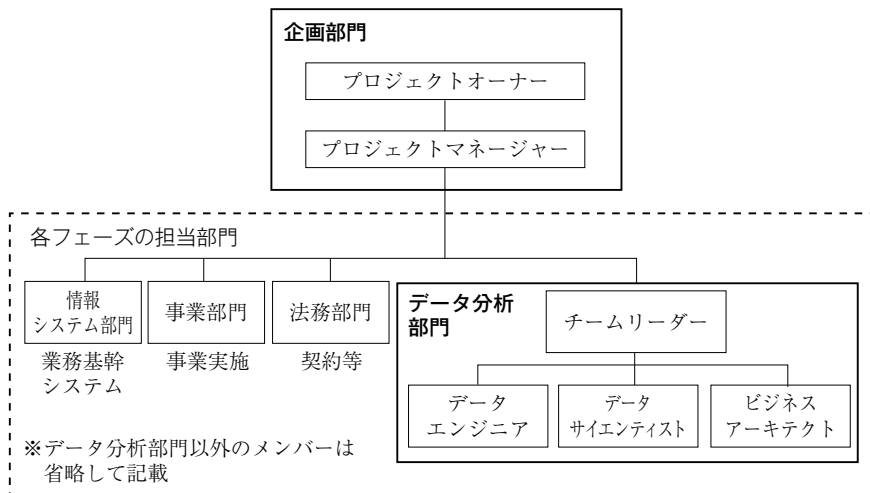


図 1.3 データ活用に必要なプロジェクトチームの例

データの活用に対する理解が十分でない場合、データが分析要件を満たさないなどの技術的な問題が発生したり、契約上の要件が不十分で不適切な個人情報の利用や使用条件の違反などのコンプライアンス上の問題が発生し、プロジェクトの失敗につながる恐れもある。

1.2 従来型プロジェクトとデータ活用型プロジェクトの定義

IoT (internet of things) を活用して安価に収集できるデータや、さまざまな業務システム等によって収集され蓄積されたデータなどのいわゆる「ビッグデータ」として活用されるデータの取り扱い、従来の調査、実験での観測で得られたデータの取り扱いと本質的な違いがある。このことを明確にするために本書では、調査・実験での観測によってデータ収集を行う「従来型プロジェクト」と、いわゆるビッグデータ等の既存データを中心としたデータを用いる「データ活用型プロジェクト」という区分を定義し、その違いを説明する (表 1.1)。

従来型プロジェクトでは、実験計画や調査設計、データ収集機器の開発などのデータ収集のためのフェーズをプロジェクト内に内包していることを前提と

索引

| | | | | |
|---------------|--------------------|--------|--------------|-----|
| 【あ行】 | クリエイティブ・コモンズ・ライセンス | 34 | 成果物 | 86 |
| アジャイル型プロジェクト | クリエイティカルパス | 72 | 正規化 | 122 |
| | クロスバリデーション法 | 124 | 正検出 | 126 |
| アジャイル思考 | 現状理解・問題把握 | | 正未検出 | 126 |
| アジャイルソフトウェア | フェーズ | 52 | 遷移行列 | 149 |
| 開発宣言 | 構造化データ | 115 | 【た行】 | |
| 委託 | 顧客マスターデータ | 116 | 第三者提供 | 24 |
| イタレーション | 誤検出 | 126 | 体制 | 86 |
| イタレーションミーティング | 個人識別符号 | 25 | 台帳 | 37 |
| | 個人情報 | 25 | 担当者ミーティング | 71 |
| 一次利用 | 個人情報保護法 | 22 | チケット | 79 |
| ウォーターフォール型 | コミュニケーションツール | | チーム内会議 | 71 |
| プロジェクト | | 86 | チームリーダー | 42 |
| エクспанディング法 | 混同行列 | 127 | デイリー | 79 |
| オーダーメイド集計 | 【さ行】 | | デイリースタANDING | 80 |
| オプトアウト規定 | 再現率 | 127 | 適合率 | 127 |
| オープンデータ | 差分の差分法 | 131 | デジタルスキル標準 | 16 |
| | サンキーダイアグラム | 151 | データウェアハウス | 107 |
| 【か行】 | 散布図 | 142 | データエンジニア | 19 |
| 会議体 | 施策実施・検証フェーズ | 57 | データサイエンスプロ | |
| 改行コード | 実施方針 | 85 | フェッショナル | 19 |
| 外部キー | 譲渡・寄付契約 | 36 | データ提供者 | 40 |
| 学術研究 | 進捗確認会議 | 71 | データビジネスストラテ | |
| 可視化 | 進捗管理 | 70, 71 | ジスト | 18 |
| 課題管理 | 進捗管理表 | 71 | データ分析 | 105 |
| 課題管理表 | スコープ | 85 | データマート | 108 |
| 課題設定フェーズ | ステアリングコミッティ | 41 | データラボ | 40 |
| 課題の5W1H | ステアリングコミッティ | | データ利用契約 | 35 |
| 仮名加工情報 | 会議 | 71 | データレイク | 107 |
| 管理ログ | ステークホルダー | 40 | 同一性チェック | 114 |
| キックオフミーティング | スプリント | 79 | 統計法に基づく統計データ | |
| 機微情報 | スプリントバックログ | 79 | | 36 |
| 教師無し学習 | スプリントレビュー | 81 | 匿名加工情報 | 27 |
| 共同利用者 | 正解率 | 127 | トランザクションデータ | 116 |

| | | | | | |
|------------|-------------|--------------|--------|-------------|-----|
| | 【な行】 | 秘密保持契約 | 34 | 変更管理委員会 | 74 |
| | | 標準化 | 122 | 【ま行】 | |
| 内部結合 | 118 | 不正競争防止法 | 32 | マイルストーン | 85 |
| ナレッジ管理 | 62 | プライマリキー | 116 | マスターデータ | 116 |
| 二次利用 | 13 | プロジェクトオーナー | 41 | 未検出 | 126 |
| | 【は行】 | プロジェクト管理方法 | 86 | 文字コード | 116 |
| 箱ひげ図 | 147 | プロジェクト推進者 | 40 | 【や行】 | |
| 派生データ | 14 | プロジェクトのリスク | 87 | 要配慮個人情報 | 25 |
| バックログ | 78 | プロジェクトマネージャー | | 【ら行】 | |
| 非構造化データ | 115 | プロダクトバックログ | 78 | 利用規約 | 34 |
| ビジネスアーキテクト | 18 | 分析環境 | 110 | ローリング法 | 126 |
| ヒストグラム | 140 | 分析報告会議 | 71 | ロール | 41 |
| 左外部結合 | 118 | 変更管理 | 70, 74 | | |
| ヒートマップ | 142 | ——のルール | 86 | | |

| | | | | | |
|------------------|-------------|-----------------|-----|-----------------------|----------|
| | 【英字】 | INNER JOIN | 118 | Python | 112 |
| A/B テスト | 128 | ipysankeywidget | 135 | Queensland Household | |
| Active data | 10 | ipywidgets | 135 | Travel Survey—2020-21 | 134 |
| Anaconda | 112 | k-匿名化 | 27 | RMSE | 126 |
| Causal Impact | 132 | LEFT OUTER JOIN | 118 | seaborn | 135 |
| CC BY 4.0 | 134 | 多様性 | 27 | SHA-256 ハッシュ値 | 115 |
| Creative Commons | | macroFl | 128 | SQL | 108 |
| Attribution 4.0 | 134 | MAE | 126 | WBS | 69 |
| DID | 131 | matplotlib | 135 | ydata-profiling | 135, 136 |
| EDA | 120 | microFl | 128 | 【数字】 | |
| ETL | 108 | pandas | 135 | 100 パーセント・ルール | 72 |
| F1 スコア | 127 | Passive data | 11 | 5W2H | 60 |
| GDPR | 29 | PMBOK ガイド | 39 | | |
| | | PyPI | 112 | | |

— 著者略歴 —

- 2006年 神戸大学工学部建設学科卒業
2008年 神戸大学大学院自然科学研究科博士前期課程修了（建設学専攻）
2010年 神戸大学大学院工学研究科博士後期課程修了（市民工学専攻）、博士（工学）
2010年 日本学術振興会特別研究員 PD
2011年 東京工業大学助教
2016年 東京大学講師
2021年 東京大学准教授
2021年 阪急阪神ホールディングス株式会社データアナリシスディレクター
現在に至る
- 2022年 東京大学特任准教授（～2023年）
2022年 東京大学空間情報科学研究センター客員研究員（～現在）
2023年 株式会社ソーシャル・デザイナーズ・ベース代表取締役社長（～現在）

データ活用型プロジェクトのマネジメント

Project Management for Data Utilization and Application

© Takahiko Kusakabe 2024

2024年3月27日 初版第1刷発行



検印省略

著者 日下部 貴彦
発行者 株式会社 コロナ社
代表者 牛来真也
印刷所 壮光舎印刷株式会社
製本所 株式会社 グリーン

112-0011 東京都文京区千石 4-46-10
発行所 株式会社 コロナ社
CORONA PUBLISHING CO., LTD.
Tokyo Japan
振替00140-8-14844・電話(03)3941-3131(代)
ホームページ <https://www.coronasha.co.jp>

ISBN 978-4-339-05281-7 C3051 Printed in Japan

(新井) I



JCOPY <出版者著作権管理機構 委託出版物>

本書の無断複製は著作権法上での例外を除き禁じられています。複製される場合は、そのつど事前に、出版者著作権管理機構（電話 03-5244-5088, FAX 03-5244-5089, e-mail: info@jcopy.or.jp）の許諾を得てください。

本書のコピー、スキャン、デジタル化等の無断複製・転載は著作権法上での例外を除き禁じられています。購入者以外の第三者による本書の電子データ化及び電子書籍化は、いかなる場合も認めていません。落丁・乱丁はお取替えいたします。