

データ分析手法の原理の理解に対する自信の分析

Analysis of confidence in understanding principles of data analysis methods

安田 健人[†]

島川 博光[†]

Kento Yasuda Hiromitsu Shimakawa

1. はじめに

近年、IT 業界だけでなく、製造や物流、医療などの幅広い業界においてデータサイエンティストの需要が高まってきている。しかし、データサイエンスの技術習得には、確率統計、最適化、プログラミングなど多様な知識、スキルが求められ、くじける者が多い。その要因として、理解に自信がないことが考えられる。

本研究ではデータ分析手法の理解度を推定する。くじける人の多くが、知識の整理、順序立てた説明ができないことに着目する。これらの人の多くが、ある箇所で思い違い、勘違いして自信を無くしていることが考えられるため、思い違い、勘違いしているかどうかの評価から自信の揺らぎを分析する。データ分析手法の原理の理解に対する自信を分析することで、データサイエンス教育のくじけを防止する。

2. データサイエンスに対する理解と自信

2.1 データサイエンス教育

データサイエンティストを育成するためには、形だけのリテラシーレベルの教育では成果は得られない。データサイエンティストには、数学的知識だけでなく、現場特有の用語や概念を理解する能力も求められる。データサイエンティストは、分析手法に基づく数学的根拠を理解し、その理解のもと、多様な分析手法から最適なものを選んで、眼前の実データに運用できる能力を持たなければならない。このようにデータサイエンティストには、広い分野の知識と、それを実問題に適用する問題解決能力が求められる。

2.2 自信と理解度の関係

自信がある人は対象分野の概念や手続きをよく理解しており、一方で自信がない人はこれらに対する理解が乏しいと考えられる。したがって、自信のある人は知識の整理、順序立てた説明ができる人が多く、自信のない人は何かの分野、概念について思い違い、勘違いしているため、知識の整理、順序立てた説明ができないと考えられる。

本研究では、どの分野、概念が理解できていないかを、記述問題に対する回答をテキスト分析す

ることで判別する。理解できていない分野、概念を同定することで、くじけそうな学習者の支援を提案することができると考えられる。

3. 理解と自信の推定

3.1 全体の流れ

データを分析する課題を与える。被験者の解答からテキストデータを収集する。次に類似度の推定と理解度の分類をするために、正解例の単語の頻度で、理解しているかどうかを判別する。その判別には AnGAN を用いる。判別後は、重要単語の特定をするために、正解例で出てくる単語を頻度の高い順に列挙する。最後に、勘違い、思い違いしているかどうかを評価することで自信の揺らぎを分析する。

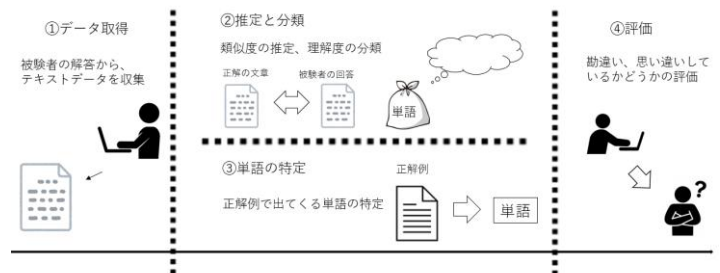


図 1 : 手法概要図

3.2 データ取得

データを分析する課題において、穴埋め課題、考察課題を出題する。本研究ではとくに後者に着目する。穴埋め課題では、データを分析するための Python コードに複数の空欄を設けて置き、この空欄を学習者が正しく補充できるかどうかを調べる。この課題から、学習者が分析のためのコードについての知識を正しく理解しているかを推定する。

考察課題では、上記の Python コードをデータに適用した結果、出力される数値が何を意味しているかを学習者に記述してもらう。考察課題から、学習者が、与えられた課題を分析するうえで指定された手法を使用する意義を理解しているか、かつ、その手法の適用結果を正しく読み解けていかを推定する。考察課題の回答文をテキストデータとして収集する。

3.3 類似度の推定と理解度の分類

学習者の回答と正解例との類似度は、`dog2vec` を使用して推定する。回答と正解例のそれぞれの単

[†] 立命館大学 Ritsumeikan University

