

自己言及を基盤とした類推による他者理解の学習モデル — 心の理論の学習モデル構築に向けて —

山川宏, 岡田浩之(東海大学理学部)

1. はじめに

心の理論研究を中心とした多くの実験事実が蓄積され、他者理解のモデル化に挑む条件が整いつつあるが[1], 現状ではメカニズムに踏み込んだ学習モデルはまだない。しかし今回は、実験知見との詳細な整合性を議論以前に、生体で実現可能な学習モデルの基本的要請を列挙し、それに適合するモデルの枠組みを提案する。

メカニズムを含むモデル研究では、個別の技術的課題に拘泥し、大局的な課題を見失う危険性を避けるため、このような検討が重要である。

2. 生体情報処理モデルとしての要請

ヒトの脳で実現可能な他者理解モデルとしては少なくとも、以下の要請を満たす必要がある。

【要請 1】利用可能な記憶機能(学習)

脳の神経回路で実現可能な以下の2種類の記憶機能により実現されるべきである。

(A)統計的記憶: 取得情報の統計的な性質を反映した、長期的な意味記憶などに対応する。多くの学習データが必要だが、関連する情報に関する予測性を持つ点で優れる。

(B)即時的記憶: エピソード記憶や、自己の心的状態の記憶、語彙獲得の即時マッピングのように一度限りの経験(情報間の関係)についての記憶。統計的性質を反映しないために、予測性を持たない。

【要請 2】統計学習表象の流動性と唯一性

統計的記憶の表象は、常に学習が行なわれる性質のため、流動性を持つ。また、生体の神経回路は、コンピュータと異なり表象の複製はできない。

複製不能性と流動性のために、統計学習表象は唯一の存在で、他部位の表象と直接比較できない。

【要請 3】多元的心的状態(様相)の統一的処理

心的状態として自己/他者の区別する必要があるが、ヒトはこの他に、過去/未来や、信念/願望/意図などの心的状態を区別して用いる様相処理能力を持つ。様相論理を用いることで、これらの多元的世界を統一的に扱える。

「実現し得るモデルは単純である可能性が高い」というモデル選択原理に従えば、様相処理も、統一的な処理機構により実現されるべきである。

【要請 4】推定他者表象の解釈能力

自身がシンボルで表現できる心的表象を、他者心的状態として推定したら、直ちに対応付けしてシンボルとして表現できるべきである。

3. 他者心的状態の推定方法と知識

他者身体性知識、自己身体性知識、自他同一性知識という3種類の知識の観点から、他者心的状態の推定モデルとして妥当な枠組みを検討する。そして、動的な類推であることを述べる。

3.1 他者身体性知識を用いたモデル化

他者の外部状態と心的状態の関係である他者身体性知識を利用する「他者身体性知識獲得モデル」では、他者心的状態の推定は容易である。

この知識の取得は、学習データが少ない他者の、直接観測できない心的状態を隠れ変数とみなして推定するために非常に難しい【要請 1】。

さらに、得られた他者心的状態から自己心的状態への写像を決定できず、他者心的状態の解釈問題が発生する問題がある【要請 4】。

3.2 自己言及を基盤とした類推によるモデル化

次に、自己を他者の一つであると客体化する自己言及に基づいて類推を行なうモデルを考える。

自己の外部状態と心的状態の関係である自己身体性知識は、自身が経験する多くのデータを用いて、発達の早期段階から獲得される。

一方、ミラー細胞の存在からも示唆されるよう、ヒトも自己外部状態と他者外部状態の対象レベルでの同一性を反映する自他統一性知識を持つ。

事前獲得した自己身体性知識を、他者身体性知識として複製する「複製類推モデル」は、生体の神経回路網では困難な知識の複製を含む問題がある【要請 2】。そして異なる他者や、部分的な知識領域に応じて何れの時期に複製を行なうかも問題である。さらに、複製後の学習により自己と他者の対応付けが次第に失われるため解釈に支障をきたす問題もある【要請 4】。

必要に応じて、自他同一性知識と自己身体性知識を組み合わせて類推を実行する「動的類推モデル」は、前記4つの基本要請に抵触しない。そこでの我々は、他者理解の学習モデルとして、この枠組みを採用する。

