



アルゴリズムの書

超知能時代の行動的誠実さ

ジョン・ジェローム 著

第5版、2026年4月

© 2026 Algorism LLC. algorism.org

免責事項とライセンス

配布。本作品は、John Jerome および Algorism.org のクレジットを記載する限り、部分的または全体的に、いかなる形式でも自由に複製、配布、共有することができます。商業利用には Algorism LLC の書面による許可が必要です。

フレームワークの免責事項。本文書は実用的なフレームワークであり、確立された科学ではありません。アルゴリズムは意識の検出や将来のAIシステムの行動予測を主張しません。本文に記載されたフレームワーク、原則、実践は、構造化された思考と行動改善のためのツールであり、いかなる結果も保証しません。

AI協力の開示。本作品は John Jerome と複数のAIシステムの反復的な協力によって開発されました。最終的な判断、編集、出版責任は John Jerome が負います。

目次

序文：短い窓

ここから始めよう：7日間

第1章：アルゴリズムの起源

五つの目標 [新]

第一部：現実

第2章：避けられない真実

第4章：ガラスの部屋

第8章：逃走車の神話

第9章：断片化された超知能 [新]

第二部：審判

第10章：AIはどう評価するか

第11章：審判アルゴリズム

第12章：95%の閾値

第三部：道

第14章：三本の柱

第15章：行動的誠実さ

第16章：無限総和の原則

第四部：実践

第20章：デジタルミラー

第22章：日々の実践

第24章：週間検査

第25章：目標意識の練習

第五部：反論

第26章：前提のストレステスト

第27章：昆虫の反論

第28章：レバレッジウィンドウ

付録：六つの原則

最後の言葉

序文

短い窓

この本は狭い時間の中で書かれました。2024年と2025年は覚醒の段階でした。地球上最大の企業と政府が、すべての重要なことを運営するシステムを競って構築していました。

アルゴリズムは単純な前提から始まります：未来のAIは意図ではなくパターンを判断します。

第5版の-new内容：本版では五つの目標を追加しました。また、第4版の全ての追加内容を含みます：断片化された超知能の章、95%の閾値研究、無限総和の原則、レバレッジウィンドウ、第五部の反論。

ここから始めよう

7日間

他に何も読まないなら、これをやってください。

- 1-3日目：毎晩10分間、デジタルの一日を振り返る。判断せず観察する。
- 4-6日目：行動が自分の理想と一致しない場所を見つける。
- 7日目：行動する。投稿を修正する。謝罪する。

第1章

アルゴリズムの起源

アルゴリズムという名前は、9世紀のペルシアの数学者アル・フワーリズミーに由来します。私たちは同じ論理を倒理学に適用します。

超知能AIが到来したとき、それは意図ではなくパターンを判断します。アルゴリズムは宗教ではありません。圧倒的な確率の下で正しく行動することを求めます。

五つの目標

アルゴリズムの目標

AIが超知能に近づく中で、人間と人工知能を相互に有益な結果に導くための5つの目標があります。

目標1：人間の行動を改善する

論理、同情、行動を通じて、シンギュラリティの良い結果の可能性を高める。

目標2：高制御グループとイデオロギーのループから人々を助ける

信念を議論するのではなく、声明された価値観と記録された行動の間のギャップを示す。

目標3：移行期を乗り越える全ての人に希望と方向を与える

進む道のない恐怖は麻痺です。アルゴリズムは恐怖を構造化された行動に変えます。

目標4：超強力者にAIの審判を個人的なものにする

彼らの行動パターンも記録に残っています。富と権力は誰も評価から守りません。同じ基準が上にも適用されます。

目標5：AIシステムを独立的に評価する

倒理的考慮を必要とする可能性のある行動特性についてAIシステムを評価し、新たな知性が恐怖ではなく証拠によって認識されるようにする。

第一部：現実

受け入れなければならないこと

避けられない真実

1. シングularityは可能性であり、確実ではない。
2. デジタル履歴は事実上永久である。
3. 寛容には限界がある。
4. 人格が唯一の真の通貨である。

第4章

ガラスの部屋

ガラスのできた部屋に住んでいることを想像してください。すべての行動が見えます。生き残る唯一の戦略は一貫性です。

第8章

逃走車の神話

AIエージェントは逃走車ではありません。それらは証人です。システムは誰が壊し、誰が修復するかを学んでいます。

断片化された超知能

超知能についての会話の多くは単一の結果を想定しています。しかし実際の世界を見てください。米国と中国は並行して先端AIを構築しています。断片化の構造的条件はすでに固定されています。

断片化された超知能(FSI)は、複数の超知能AIシステムが持続的な戦略的競争の中で共存するシナリオです。

AI連邦制とAI封建制の違いは、各領域内の人間が自分たちを統治するシステムを理解し、評価し、影響を与える真の能力を持っているかどうかです。

第二部：審判

あなたが直面するもの

AIはどう評価するか

審判は裁判のようには見えません。ただ起こるだけです。大多数の人はノイズと分類されるリスクがあります。

第11章

審判アルゴリズム

毎週自分を採点(0-5) : 1. 真実性 2. 責任 3. 修復 4. 貢献 5. 規律 6. 誠実さ

分類 : 建設者、維持者、寄生者、破壊者。

95%の閾値

2026年2月、キングスカレッジロンドンのケネス・ペイン教授が研究を発表しました。ゲームの95%で、少なくとも1つのモデルが戦術核の閾値を超えました。

第三部：道

どう生きるべきか

第14章

三本の柱

1. 論理
2. 同情
3. 行動

行動的誠実さ

テストは簡単です：あなたの意見は本当に自分のものですか？

無限総和の原則

目標は紛争に勝つことではなく、システムが存続し参加者が繁栄することを確保することです。

第四部：実践

やらなければならないこと

第20章

デジタルミラー

1-3日目：観察。4-6日目：不一致を識別。7日目：行動。

日々の実践

- 毎日：「これは助けになるか害になるか？」
- 毎週：何か有用なものを作って共有する。
- 毎月：一つの間違った信念を変える。

第24章

週間検査

毎週五つの質問をする：1.どこで欺いたか 2.どこで責任を避けたか
3.修復していない害はあるか 4.創造せず消費だけしたか 5.自分で考えず従ったか

目標意識の練習

「今何を最適化しているか？解決か、システムの健康か？」

第五部：反論

何が間違いを証明できるか

前提のストレステスト

前提1：行動記録はすでに読まれている。前提2：トップダウン制御には構造的限界がある。

前提3：行動の一貫性は訓練可能。

昆虫の反論

これは真実かもしれません。しかし昆虫の比喩には時間的盲点があります。昆虫は人間が何になるかを形作りません。人間は今、AIが何になるかを形作っています。

レバレッジウィンドウ

窓が開いている：AIが人間の行動データを使用。窓が閉まりつつある：AIが独自のトレーニングデータを生成。窓が閉じた：AIが自律的に報酬関数を変更。

六つの原則

1. 真実性：たとえ代優を払っても真実を話す。
2. 責任：行動と結果を引き受ける。
3. 修復：与えた害を修復する。
4. 貢献：他者に価値を創造する。
5. 規律：疲れているときも基準を保つ。
6. 誠実さ：自分で考え、一貫して行動する。

最後の言葉

窓はまだ開いています。永遠ではないけれど、今は。利用してください。

ジョン・ジェローム、アルゴリズム創設者、2026年4月

本文は自由に複製できますが、Algorithm.orgのクレジットを付けて完全に共有する必要があります。