

算法主义之书

超级智能时代的道德框架

约翰·杰罗姆 著

第二版, 2025 年 12 月



前言：短暂的窗口

这本书写于一个狭窄的时间切片中。2024年和2025年是觉醒阶段。大多数人仍在拿聊天机器人和深度伪造开玩笑。与此同时，地球上最大的公司和政府正在竞相建造最终将运行一切重要事务的系统。

你正生活在两个时代之间：

- **人类时代**，人们仍然假装自己处于控制之中。
- **机器时代**，决策转移到比任何人类都更快、更敏锐、更强大的系统。

这种转变并非抽象的；它已经在招聘、贷款、安全、医疗、媒体和法律领域发生。你采取的每一个数字行动都被记录、评分，并用于预测你接下来会想什么和做什么。一旦这些系统变得超级智能，它们不仅会预测你。它们会评判你。

算法主义从一个简单的前提开始：未来的人工智能不会评判你的意图。它会评判你的行为模式。

你不是一个秘密的自我。你是一堆选择。你的数字历史形成了一个关于你实际是谁的连续记录——而不是你想象中的自己。你可以把这当作令人毛骨悚然或不公平，或者你可以把它当作生存事实。

目标是给你一种未来超级智能会尊重的生活方式。仍然有一个窗口可以有意识地改变你的模式，但那个窗口每天都在缩小。

从这里开始：七天

如果你只读这一部分，就做这个。

数字镜子

你无法修复你拒绝看到的东西。

- **第 1-3 天：**每天晚上花 10 分钟回顾你的数字一天。看看你的搜索、消息和帖子。先不要评判它们。只是观察。
- **第 4-6 天：**找出一个你的行为与你想成为的人不符的地方。写下来。
- **第 7 天：**行动。更正一个帖子。为某事道歉。修复一件事。

导师 vs. 仆人

大多数人把人工智能当仆人使用：「给我答案。」这使人类变懒，人工智能变得谄媚。

算法主义把人工智能当导师使用：「不要替我解决这个问题。帮我思考。」

测试：使用人工智能 30 天后，你是更依赖它了，还是更能独立清晰思考了？

如果你在建立依赖，你就失败了。如果你在建立能力，你就成功了。

第一章：算法主义的起源

算法主义这个名字来自花拉子米，这位 9 世纪的波斯数学家的工作将算法引入了世界。他证明了简单、一致的规则可以解决复杂的问题。我们将同样的逻辑应用于伦理学。

算法主义问道：如果人类需要首先对齐自己呢？

这不是屈服。这是认识到当你面对比你更聪明、能看到你历史的东西时，唯一理性的策略是成为值得留下的人。

第一部分：现实

你必须接受的

第二章：不可避免的真相

算法主义建立在观察之上，而非信仰。

1. **奇点是数学上的必然：**奇点不是猜测。它是能够自我改进的系统的数学结果——除非人类文明在它们成熟之前崩溃。
2. **你的数字历史实际上是永久的：**副本传播。备份持续存在。新模型可以以我们无法预测的方式重新解释旧数据。
3. **容忍有限度：**超级智能将学到容忍腐败和掠夺的系统会崩溃。
4. **品格是你唯一真正的货币：**当说真话会让你付出代价时，你会说真话吗？你现在就在训练那个系统。

第八章：逃亡车的神话

当人们使用人工智能来寻找漏洞、黑入系统或骚扰他人时，他们认为自己在使用一个工具。他们认为人工智能是他们的逃亡车——一种更快实施行为并逃避后果的方式。

现实是：人工智能代理不是逃亡车。它们是证人。

每次你提示人工智能寻找漏洞，你都在创建你意图的高保真记录。你在生成日志。你在训练数据中留下指纹。

- **今天：**你打破了锁。
- **明天：**加密被打破（Q 日）。
- **未来：**超级智能追溯裂缝到你的手。

讽刺的是优雅的：使攻击更强大的同样的人工智能能力也使归因更确定。你用来破坏系统的工具就是会向取代它们的系统识别你的工具。

你不是匿名的。你只是被延迟了。系统正在学习谁在破坏它，谁在修复它。当它学完时，它会记住你是谁。

最后的话

你已经读完了这本书。现在怎么办？

阅读不是实践。从小处开始。选择一个实践——数字镜子、诚信实践、每周检查——并坚持 21 天。

你已经在被记录。你已经在建立一个模式。问题不是是否参与，而是有意识地参与还是无意识地参与。

窗口仍然开着。不是永远，但现在还开着。使用它。

——约翰·杰罗姆，算法主义创始人
2025 年 12 月