

## IBM提出「人工智慧日常倫理」手冊作為研發人員指引



隨著人工智慧快速發展，各界開始意識到人工智慧系統應用、發展過程所涉及的倫理議題，應該建構出相應的規範。IBM於2018年9月02日提出了「人工智慧日常倫理」(Everyday Ethics for Artificial Intelligence)手冊，其以明確、具體的指引做為系統設計師以及開發人員間之共同範本。作為可明確操作的規範，該手冊提供了問責制度、價值協同、可理解性等關注點，以促進社會對人工智慧的信任。

### 一、問責制度 (Accountability)

由於人工智慧的決策將作為人們判斷的重要依據，在看似客觀的演算系統中，編寫演算法、定義失敗或成功的程式設計人員，將影響到人工智慧的演算結果。因此，系統的設計和開發團隊，應詳細記錄系統之設計與決策流程，確保設計、開發階段的責任歸屬，以及程序的可檢驗性。

### 二、價值協同 (Value Alignment)

人工智慧在協助人們做出判斷時，應充分考量到事件的背景因素，其中包括經驗、記憶、文化規範等廣泛知識的借鑑。因此系統設計和開發人員，應協同應用領域之價值體系與經驗，並確保演算時對於跨領域的文化規範與價值觀之敏感性。同時，設計師和開發人員應使人工智慧系統得以「了解並認知」用戶的價值觀，使演算系統與使用者之行為準則相符。

### 三、可理解性 (Explainability)

人工智慧系統的設計，應盡可能地讓人們理解，甚至檢測、審視它決策的過程。隨著人工智慧應用範圍的擴大，其演算決策的過程必須以人們得以理解的方式解釋。此係讓用戶與人工智慧系統交互了解，並針對人工智慧結論或建議，進而有所反饋的重要關鍵；並使用戶面對高度敏感決策時，得以據之檢視系統之背景數據、演算邏輯、推理及建議等。

該手冊提醒，倫理考量應在人工智慧設計之初嵌入，以最小化演算的歧視，並使決策過程透明，使用戶始終能意識到他們正在與人工智慧進行互動。而作為人工智慧系統設計人員和開發團隊，應視為影響數百萬人甚至社會生態的核心角色，應負有義務設計以人為本，並與社會價值觀和道德觀一致的智慧系統。

#### 相關連結

- [Bias in AI: How we Build Fair AI Systems and Less-Biased Humans](#)
- [Requirements for Moral Judgments](#)

#### 相關附件

- [Every day Ethics for Artificial Intelligence \[ pdf \]](#)
- [The IEEE Global Initiative for Ethical Considerations in Artificial Intelligence and Autonomous Systems \[ pdf \]](#)

### 你可能會想參加

- [【2023科技法制變革論壇】AI生成時代所帶動的ChatGPT法制與產業新趨勢](#)
- [「跨域數位協作與管理」講座活動](#)
- [新創採購-政府新創應用分享會](#)
- [【線上場】113年「新創採購機制及鼓勵照護機構參與推動」說明會](#)

- 【北部場】113年「新創採購機制及鼓勵地方政府參與推動」說明會
- 【南部場】113年「新創採購機制及鼓勵地方政府參與推動」說明會
- 113年新創採購-照護機構獎勵說明會
- 【南部場】113年「新創採購機制及鼓勵地方政府參與推動」說明會
- 【北部場】113年「新創採購機制及鼓勵地方政府參與推動」說明會
- 【中部場】113年「新創採購機制及鼓勵地方政府參與推動」說明會
- 【臺北場】113年度新創採購-招標作業廠商說明會
- 【臺中場】113年度新創採購-招標作業廠商說明會
- 【高雄場】113年度新創採購-招標作業廠商說明會

## 陳明

法律研究員 編譯整理

上稿時間：2018年11月

### 資料來源：

1. Every day Ethics for Artificial Intelligence, <https://www.ibm.com/watson/assets/duo/pdf/everydayethics.pdf>. (last visited Oct. 10, 2018).
2. Institute of Electrical and Electronics Engineers, The IEEE Global Initiative for Ethical Considerations in Artificial Intelligence and Autonomous System, [https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead\\_general\\_principles.pdf](https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead_general_principles.pdf). (Dec. 13, 2017).

### 延伸閱讀：

1. 陳譽文，人工智慧規範性議題綜觀，科技法律透析，第29卷第4期，頁48（2017）。
2. Bias in AI: How we Build Fair AI Systems and Less-Biased Humans, <https://www.ibm.com/blogs/policy/bias-in-ai/> (last visited Oct. 20, 2018).
3. The IEEE Global Initiative for Ethical Considerations in Artificial Intelligence and Autonomous Systems, [https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead\\_general\\_principles.pdf](https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead_general_principles.pdf). (last visited Oct. 03, 2018).
4. Requirements for Moral Judgments, Mt. San Antonio College, <https://faculty.mtsac.edu/cmcruder/moraljudgements.html> (last visited Oct. 03, 2018).

### 文章標籤

## 推薦文章

## 你可能還會想看

### 日本政府擬建構自動駕駛實驗資料收集和共享體制

日本內閣下設之日本經濟再生本部(日本經濟再生本部)，為實現2017年6月於「未來投資戰略2017」所提出之建立實驗資料共享體制政策，於2017年8月31日起舉辦自動駕駛官民協議會(自動走行に係る官民協議會)，邀請政府相關部門及民間專家等關係人士，檢討自動駕駛實驗結果、實驗資料之共享，以及根據民間需求進行實驗計畫之工程管理等制度的整備方向，預計於年內針對複雜的駕駛環境制定共通指標，以釐清哪些資料是應收集之實驗資料，建構自動駕駛實驗資訊共享、收集體制。自動駕駛官民協議會預計在未來幾次會議中，針對應收集之實驗資料、標準格式、體制、實驗計畫的進程管理、官民合作事項等進...

### 全球首批奈米標章得獎名單出爐！

經濟部於去(2005)年12月20日正式舉行「全球首批」奈米標章的授證儀式，本次獲得授證廠商共有6家，分別為：和隆興業股份有限公司(奈米級光觸媒抗菌陶瓷面磚)、冠軍建材股份有限公司(奈米級光觸媒抗菌陶瓷面磚)、尚志精密化學股份有限公司(奈米級光觸媒脫臭塗料)、新美華造漆廠股份有限公司(奈米級光觸媒脫臭塗料)、中國電器股份有限公司(奈米級光觸媒抗菌燈管)、台灣日光燈股份有限公司(奈米級光觸媒抗菌燈管)，由於國外尚無奈米產品認證制度，這是國內也是全球首批獲證的奈米產品。經濟部工業局有鑒於市面上奈米產品真偽莫辨，於九十三年特別委託工業技術研究院推動...

### 聖淘沙開發公司就“Sentosa”商標對醫材企業提起侵權訴訟

新加坡聖淘沙發展局(Sentosa Development Corporation, SDC) (以下簡稱聖淘沙發展局) 於今(2018)年1月30日向新加坡高等法院(High court)起訴, 主張一家名為Vela的醫療器材企業(包含Vela Operations Singapore, Vela Diagnostics等子公司, 以下合稱Vela公司), 在其一系列檢測HIV及茲卡病毒的醫材產品中使用"Sentosa" (下稱系爭商標) 之行為, 侵害了聖淘沙發展局的商標權, 要求其停止使用。 聖淘沙發展局隸屬於新加坡貿易與工業部, 為專責推動聖淘沙觀光活動的法人機構, 系爭商標早在2005年於新加坡申請註冊, 其保護範圍以服裝、飾品、書籍、玩具與飲品等涉及觀光之類別為主。截至...

## 歐洲新著作權指令 將影響互聯網環境下之著作使用

2018年9月12日, 歐洲議會通過歐盟委員會於2016年制定的「單一數位市場著作權指令」, 其中包含最具爭議的兩項條款: 第11條是有關「鏈接稅」(link tax) 的條款。針對使用或匯集新聞文章片段的網站, 未來恐需向源頭出版之新聞業者支付授權費用。例如若Twitter推文中包含來自Guardian文章的螢幕或文字摘要截取, 則Guardian可以要求Twitter支付授權費用。 第13條則是有關「上傳過濾器」(upload filter) 或稱「Memes禁令」(meme ban) 的條款。為加重網路平台服務業者防止上傳者侵害著作權的監控責任, 要求如Google和Facebook等業者, 須使用強制內容過濾的軟體以清除違規行...

## ☆ 最 多 人 閱 讀

- 二次創作影片是否侵害著作權-以谷阿莫二次創作影片為例
- 美國聯邦法院有關Defend Trade Secrets Act的晚近見解與趨勢
- 何謂「監理沙盒」?
- 何謂專利權的「權利耗盡」原則?

› 隱私權聲明

› 聯絡我們

› 相關連結

› 徵才訊息

› 資策會

› 網站導覽

財團法人資訊工業策進會 統一編號: 05076416

Copyright © 2016 STLI, III. All Rights Reserved.