

# A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles

(Un approccio logico-metaontologico al problema della *(meta)data veracity* nei sistemi di estrazione automatica di metadati da articoli scientifico-giuridici)

Simone Cuconato\*

## Abstract

In an increasingly data-driven world, the question of data – or metadata – veracity is now a central issue not only in the world of information but also in the legal one. Data veracity describes a closeness to truth on a higher level than a measure such as accuracy does. High veracity data is data that can be relied upon when making decisions, thus reducing the risk of basing choices on untrue information. The article uses epistemic logic **T** to model structured metadata automatically extracted from legal papers, and the tools of metaontology to propose a definition of veracity as truthmaker.

**Keywords:** (meta)data veracity; truthmaker; applied epistemic logic; metadata modelling<sup>§</sup>

---

\* Department of Computer Engineering, Modelling, Electronics and Systems Engineering, University of Calabria and Istituto di Informatica e Telematica (IIT) - CNR, Via P.Bucci, 87036, Rende (CS), Italy; simone.cuconato@unical.it.

<sup>§</sup> Received on May 9th, 2022. Accepted on December 28th, 2022. Published on December 30th, 2022. DOI: 10.23756/sp.v10i1.784. ISSN 2282-7757; eISSN 2282-7765. ©Cuconato. This paper is published under the CC-BY licence agreement.

### Sunto

In un mondo sempre più guidato dai dati, la questione della veridicità dei dati – o metadati – è una questione ormai centrale non solo nel mondo dell’informazione, ma anche in quello giuridico. La veridicità dei dati descrive una vicinanza alla verità ad un livello più alto di una misura come l'accuratezza. I dati ad alta veridicità sono dati su cui si può fare affidamento quando si prendono decisioni, riducendo, in questo modo, il rischio di fondare le proprie scelte su informazioni non veritiere. L’articolo usa la logica epistemica **T** per modellare metadati strutturati estratti automaticamente da articoli scientifico-giuridici, e gli strumenti della metaontologia per proporre una definizione di *veracity* come *truthmaker*.

**Parole chiave:** (*meta*)data veracity; *truthmaker*; logica epistemica applicata; modellizzazione dei metadati

## 1. Introduzione

Sono passati più di dieci anni da quando Chris Anderson, allora caporedattore dell'influente rivista tecnologica *Wired*, pubblicò un articolo intitolato “*The End of Theory: The Data Deluge Makes the Scientific Method Obsolete*”<sup>1</sup>. L'articolo di Anderson è diventato rapidamente il manifesto ideologico dell’entusiasmo “datacentrico” ed è articolato lungo due punti chiave.

Primo: “fidatevi, è conveniente”. Google ci ha insegnato che non è importante capire perché una pagina web è "migliore" di un'altra, ma è sufficiente fidarsi dell'ordinamento prodotto dall'algoritmo PageRank. La comodità di ricevere una risposta molto semplice a una domanda potenzialmente molto complicata, senza dover necessariamente sviluppare alcuna analisi semantica o causale, è diventata presto la chiave del successo di Google.

Secondo: “i modelli scientifici sono obsoleti”. La disponibilità senza precedenti di dati prodotti più o meno consapevolmente da tutti noi ci permette di ripensare radicalmente la relazione tra i dati e i meccanismi che li generano.

---

<sup>1</sup> [2].

*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

Secondo Anderson possiamo smettere di cercare modelli: invece di procedere per "congetture" e "confutazioni" nello spiegare le osservazioni, il diluvio di dati ci permette di rinunciare al laborioso compito di costruire modelli per i fenomeni di interesse, in favore del compito molto più facile di analizzare le correlazioni individuate da sofisticati algoritmi statistici.

In questo lavoro, ci muoveremo in una direzione opposta a quella tracciata da Anderson.

Primo: "non fidarsi". Oggi più che mai è necessario porre l'accento sulla qualità delle informazioni e sulla veridicità dei dati: un'analisi semantica è necessaria.

Secondo: "i modelli scientifici sono fondamentali". È molto difficile pensare ai dati senza che essi rispondano a un'ipotesi di modellizzazione. L'idea semplicistica che *petabyte* di dati possano essere autosufficienti e che i dati possano essere visti come un sostituto della modellazione scientifica non è sostenibile.

In particolare, in questo contributo: *i*) proporrò un modello basato sulla logica epistemica per formalizzare i metadati estratti da articoli scientifico-giuridici tramite sistemi di estrazione automatica fondati sull'intelligenza artificiale; *ii*) porrò l'accento sul tema della qualità dei dati grazie alla definizione di un principio metaontologico di veridicità come *truthmaker*. Se fino a pochi anni fa il costo dell'informazione era l'aspetto più rilevante, al contrario, oggi la qualità delle informazioni è diventata più importante che mai. Per questo motivo, la veridicità dei dati – nel nostro caso metadati – è stata proposta come la quarta "V" – accanto a Volume, Varietà e Velocità – dei *big data*<sup>2</sup>.

Ma cosa sono i metadati? Nel mondo dell'informazione, i metadati rappresentano la base informativa di "secondo livello", che descrive, struttura e gestisce i dati primari o le informazioni su cui vengono appoggiate le risorse

---

<sup>2</sup> Si veda [11].

informativa<sup>3</sup>. Attualmente i metadati sono necessari non solo per gestire, conservare e reperire gli oggetti informativi, ma rappresentano le pedine fondamentali nel *Semantic Web* avendo un ruolo chiave nell'indicizzazione e nell'identificazione, nella classificazione e nella catalogazione, nella conservazione, nella verifica dell'integrità e dell'affidabilità e nella gestione dei diritti, nonché nella distribuzione, nella ricerca e nel recupero delle risorse digitali. I metadati, siano essi descrittivi, strutturali, amministrativi o per la *long term digital preservation*, alla fine, sono accomunati da un unico obiettivo multifunzionale: quello di contribuire a una gestione e conservazione più chiara e modulare degli oggetti digitali. La metadattazione automatica permette di estrarre direttamente i metadati dalle fonti documentali. L'estrazione dei metadati come tecnologia fondamentale per il processo automatico dei documenti ha avuto un grande successo in numerose applicazioni e domini. Molte delle soluzioni proposte da tali sistemi sono basate tecniche sub-simboliche di intelligenza artificiale, come il *machine learning* (ML).

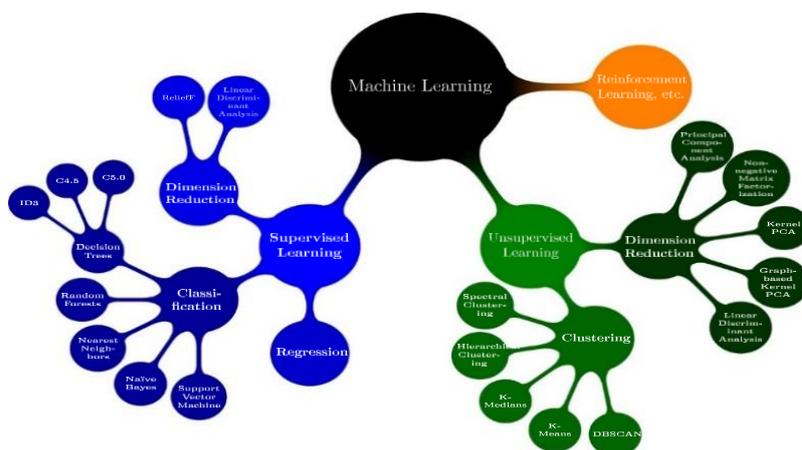
In generale, le modalità con cui il ML permette agli algoritmi di fare apprendimento con i dati sono classificate in cinque categorie: *i*) apprendimento supervisionato, nel quale vengono presentati al modello scelto gli esempi formati dagli input e relativi output desiderati con lo scopo di far apprendere una regola generale in grado di mappare gli input negli output; *ii*) apprendimento non supervisionato, nel quale vengono forniti al modello scelto solo gli esempi formati dagli input, senza alcun output atteso, con lo scopo di fargli apprendere in autonomia una qualche struttura nei dati d'ingresso; *iii*) apprendimento semi-supervisionato, nel quale vengono combinati i due approcci precedenti con una prima fase supervisionata sui dati aventi input e output associato, e una successiva fase non supervisionata su dati di cui non si conosce l'output associato; *iv*) apprendimento con rinforzo, con il quale si interagisce con un ambiente dinamico in cui raggiungere un certo obiettivo e a mano a mano che si esplora il dominio del problema vengono forniti dei *feedback* in termini di ricompense o punizioni secondo il comportamento eseguito; *v*) apprendimento con trasferimento, nel quale il modello scelto impara ad affrontare un certo

---

<sup>3</sup> Una delle migliori introduzioni ai metadati è [13].

*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

problema generico e, successivamente, si prende la conoscenza creata per usarla nell'affrontare un altro problema simile o più specifico. La figura 1 mostra una tassonomia dei metodi di apprendimento automatico.



**Figura 1** Una illustrazione schematica della tassonomia dei metodi di ML<sup>4</sup>

Sempre più sistemi di estrazione automatica di metadati basati su tecniche di apprendimento automatico sono diventati strumenti centrali nel mondo dell'informazione. In questo lavoro, useremo CERMINE per estrarre metadati da articoli scientifico-giuridici. CERMINE è un *framework open-source* per estrarre metadati strutturati da articoli scientifici nativi digitali. Il *framework* è basato su un *workflow* modulare e le implementazioni della maggior parte dei passi sono basate su tecniche di apprendimento automatico supervisionato e non supervisionato. Il *workflow* modulare, rappresentato in figura 2, consiste in tre percorsi (*ii* e *iii* eseguiti in parallelo): *i*) il percorso di estrazione della struttura di base richiede un file pdf come input e produce una struttura gerarchica in formato TrueViz. TrueViz è uno strumento in grado di classificare le entità di ogni pagina della struttura in quattro categorie: zone, linee, parole e caratteri. A sua volta, ogni zona è etichettata secondo altre quattro categorie: metadati, riferimenti, corpo e altro; *ii*) il percorso di estrazione dei metadati analizza le parti di metadati della struttura gerarchica, il risultato è un insieme di metadati del documento in formato XML; *iii*)

<sup>4</sup> Figura tratta da [1].

l'estrazione dei riferimenti estrae una lista di riferimenti bibliografici dal documento.

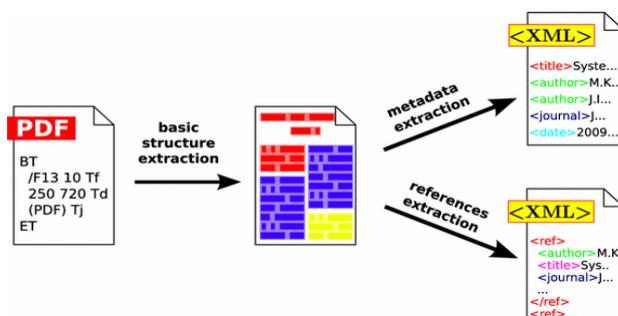


Figura 2 L'architettura del *workflow* di CERMINE<sup>5</sup>

Nel contesto degli articoli di ricerca, i metadati sono di solito di natura descrittiva e detengono una grande importanza poiché forniscono una breve panoramica su un articolo scientifico mostrando informazioni come il titolo, i suoi autori, la rivista, la bibliografia, ecc. Spesso, i ricercatori tendono a decidere la pertinenza dell'articolo con il loro dominio di interesse basandosi sulle informazioni dei metadati. Lo scopo di questo articolo è quello di applicare la logica epistemica ai sistemi di estrazione automatica di metadati da articoli scientifico-giuridici e di proporre una definizione di veridicità come *truthmaker*.

## 2. Logica epistemica standard

La logica epistemica è un'estensione della logica classica che ha come oggetto di studio gli enunciati di credenza e di sapere<sup>6</sup>. Nell'epistemologia contemporanea è ampiamente condivisa l'idea secondo cui la verità è una condizione necessaria della conoscenza. Per tale motivo: *i*) si dice che la conoscenza è *fattiva*, ossia si presuppone la verità della proposizione conosciuta; *ii*) perché si abbia conoscenza è necessario intrattenere una credenza; *iii*) la credenza deve essere giustificata. Per lungo tempo la verità, la credenza e la

<sup>5</sup> Figura tratta da [16].

<sup>6</sup> Per un'introduzione italiana alle logiche intensionali e modali si veda [8], [9], [17], mentre per uno studio più mirato alla logica epistemica si veda [18].

giustificazione sono state considerate condizioni congiuntamente sufficienti perché si abbia conoscenza. Dagli anni '60 in poi, grazie ai lavori di Gettier<sup>7</sup>, gli epistemologi contemporanei hanno sostenuto che, oltre alle tre suddette condizioni, ne occorrono altre<sup>8</sup>. Tuttavia, per quanto i logici siano particolarmente interessati al complesso dibattito che si è sviluppato tra gli epistemologi riguardo alla strategia da adottare per caratterizzare esaustivamente la conoscenza, nelle logiche epistemiche la conoscenza è generalmente caratterizzata come semplice credenza vera. In questo modo, i logici trattano le attribuzioni di conoscenza e credenza come formule contenenti operatori modali. Semanticamente questo significa che, nel valutare il valore di verità di una formula associata a un operatore epistemico, si prenda in esame un insieme di circostanze alternative a quelle attuali. Tali circostanze alternative prendono in letteratura il nome di *mondi possibili*. Poniamo per esempio che un soggetto creda che Mario Draghi sia il presidente del Consiglio italiano e che Barack Obama sia il Presidente degli Stati Uniti. I mondi compatibili con le sue credenze saranno tutti e soli i mondi in cui è vero che Mario Draghi è il presidente del Consiglio italiano e che Barack Obama è il Presidente degli Stati Uniti. Ma in base alla semantica dei mondi possibili il possesso o meno della conoscenza dipende da come stanno le cose nel mondo attuale: il nostro soggetto non può sapere che Barack Obama è il Presidente degli Stati Uniti, dato che è falso.

Sintatticamente, il linguaggio della logica epistemica proposizionale è il linguaggio della logica proposizionale classica con l'aggiunta di uno specifico operatore epistemico unario tale che

$K_a\varphi$  si legge "l'agente  $a$  sa che  $\varphi$ "

In generale, un agente può essere una persona reale, un giocatore in un gioco, un robot, una macchina, un "processo" o, nel nostro caso, un *framework* di estrazione automatica di metadati. Hintikka ha fornito una prima pionieristica formalizzazione delle attribuzioni di credenza in un linguaggio modale

---

<sup>7</sup> [10].

<sup>8</sup> Approfondire tale tema ci porterebbe lontani dagli scopi di questo lavoro. Rimandiamo a chi fosse interessato a un recente approccio formale al problema di Gettier a [20].

sfruttando delle strutture semantiche dette *model set*. Tuttavia, sarà solo dopo le pubblicazioni dei lavori di Kripke che Hittinka elaborerà un'interpretazione semantica degli operatori epistemici che possiamo presentare in termini di semantica standard dei mondi possibili secondo le seguenti linee:

$K_a\varphi$ : è vera in un mondo possibile  $w$  a condizione che  $\varphi$  sia vera in tutti i mondi compatibili con le credenze intrattenute dal soggetto epistemico in  $w$ .

Pertanto l'idea intuitiva delle logiche modali epistemiche è associare a un dato soggetto epistemico un insieme di mondi, che corrispondono a tutte le situazioni compatibili con le credenze del soggetto stesso. Vediamo ora come catturare queste intuizioni in termini formali.

**Definizione 1** [Sintassi di  $\mathcal{L}_K$ ] Dato un insieme  $\mathcal{P}$  di variabili proposizionali ed un insieme finito di agenti  $\mathcal{A}$ , definiamo il linguaggio epistemico  $\mathcal{L}_K$  come segue:

$$\varphi := p \mid \neg\varphi \mid \varphi \wedge \varphi \mid K_a\varphi$$

Dove  $p \in \mathcal{P}$  e  $a \in \mathcal{A}$ .

**Definizione 2** [Modello epistemico] Dato  $\mathcal{P}$  ed  $\mathcal{A}$  un modello epistemico  $M: \langle W, R^{\mathcal{A}}, V^{\mathcal{P}} \rangle$  è una tripla dove

- $W \neq \emptyset$  è un insieme di mondi possibili  $w_i$ , a volte chiamato il dominio di  $M$ , e denotato  $\mathcal{D}(M)$ ;
- $R^{\mathcal{A}}$  è una funzione, che produce una relazione di accessibilità  $R_a \subseteq W \times W$  per ogni agente  $a \in \mathcal{A}$ ;
- $V^{\mathcal{P}}: W \rightarrow (\mathcal{P} \rightarrow \{\text{vero}, \text{falso}\})$  è una funzione tale che per ogni  $p \in \mathcal{P}$  e per ogni  $w_i \in W$ , determina quale sia il valore di verità  $V^{\mathcal{P}}(w_i)(p)$  di  $p$  nel mondo possibile  $w_i$ .

**Definizione 3** [Verità nel modello] Dato un modello  $M: \langle W, R^{\mathcal{A}}, V^{\mathcal{P}} \rangle$  definiamo la verità di una formula  $\varphi$   $M, w_i \models \varphi$  come segue:

$$\begin{array}{lll} M, w_1 \models p & \text{sse} & V(w_1)(p) = \text{vero} \text{ con } p \in \mathcal{P} \\ M, w_1 \models \varphi \wedge \psi & \text{sse} & M, w_1 \models \varphi \text{ e } M, w_1 \models \psi \end{array}$$

*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

$$\begin{array}{ll}
 M, w_1 \models \neg\varphi & \text{sse} \quad \text{non } M, w_1 \models \varphi \text{ (spesso scritto} \\
 & M, w_1 \not\models \varphi) \\
 M, w_1 \models K_a\varphi & \text{sse} \quad M, w_2 \models \varphi \text{ per ogni } w_2 \text{ tale che} \\
 & w_1 R_a w_2
 \end{array}$$

**Definizione 4** [Assiomi e regole di inferenza] Il sistema di prova della logica epistemica che useremo è assiomatizzato utilizzando gli assiomi di **T** e la regola del *modus ponens* e della necessitazione come riportato in tabella 1:

Sistema	Regole	Assiomi	Proprietà di R
<b>T</b>	<b>MP</b> e <b>Nec</b>	$K_a(\varphi \rightarrow \psi) \rightarrow (K_a\varphi \rightarrow K_a\psi)$  $K_a\varphi \rightarrow \varphi$	R riflessiva

**Tabella 1** Logica epistemica **T**

La riflessività di  $R$  garantisce che il principio

$$\mathbf{T} \quad K_a\varphi \rightarrow \varphi$$

sia valido.

### 3. Logica e metaontologia dei metadati

Vediamo ora come adattare la logica epistemica standard alla modellizzazione dei metadati<sup>9</sup>. A livello sintattico nel nostro modello avremo solo una particolare tipologia di proposizioni  $p_\varepsilon$

$$p_\varepsilon =_{def} \mathcal{E}_{m_i}^{d_i}$$

<sup>9</sup> Per un'applicazione della logica epistemica ai sistemi di estrazione automatica di metadati da articoli scientifici sul Covid-19 si veda [6].

dove  $\mathcal{E}_{m_i}^{d_i}$  si legge “estrae il metadato  $m_i$  dal documento  $d_i$ ”.

**Definizione 5** [Sintassi di  $\mathcal{L}_{K_\varepsilon}$ ] Dato un insieme  $\mathcal{P}_\varepsilon$  di variabili proposizionali ed un insieme finito di *framework*  $\mathcal{F}$ , definiamo il linguaggio epistemico  $\mathcal{L}_{K_\varepsilon}$  come segue

$$\varphi := p_\varepsilon | \neg\varphi | \varphi \wedge \varphi | K_a\varphi$$

Dove  $p_\varepsilon \in \mathcal{P}_\varepsilon$  e  $a \in \mathcal{F}$ .

A livello semantico, invece, sostituiremo il concetto di mondo possibile con quello di *estrazione possibile*. L’idea intuitiva alla base dell’applicazione della logica epistemica alla modellizzazione dei metadati è associare a un dato *framework* un insieme di possibili estrazioni, che corrispondono a tutte le situazioni compatibili con le credenze del *framework* stesso.

**Definizione 6** [Modello epistemico per metadati] Dato  $\mathcal{P}_\varepsilon$  e  $\mathcal{F}$  un modello epistemico per metadati  $M: \langle E, R^\mathcal{F}, V^{\mathcal{P}_\varepsilon} \rangle$  è una tripla dove

- $E \neq \emptyset$  è un insieme di estrazioni possibili  $e_i$ ;
- $R^\mathcal{F}$  è una funzione, che produce una relazione di accessibilità  $R_a \subseteq E \times E$  per ogni agente  $a \in \mathcal{F}$ ;
- $V^{\mathcal{P}_\varepsilon}: E \rightarrow (\mathcal{P}_\varepsilon \rightarrow \{\text{vero}, \text{falso}\})$  è una funzione tale che per ogni  $p_\varepsilon \in \mathcal{P}_\varepsilon$  e per ogni  $e_i \in E$ , determina quale sia il valore di verità  $V^{\mathcal{P}_\varepsilon}(e_i)(p_\varepsilon)$  di  $p_\varepsilon$  nell’estrazione possibile  $e_i$ .

**Definizione 7** [Verità del modello epistemico per metadati] Dato un modello epistemico per metadati  $M: \langle E, R^\mathcal{F}, V^{\mathcal{P}_\varepsilon} \rangle$  definiamo la verità di una formula  $\varphi$   $M, e_i \models \varphi$  come segue:

$$\begin{array}{lll} M, e_1 \models p_\varepsilon & \text{sse} & V(e_1)(p_\varepsilon) = \text{vero} \text{ con } p_\varepsilon \in \mathcal{P}_\varepsilon \\ M, e_1 & \text{sse} & M, e_1 \models \varphi \text{ e } M, e_1 \models \psi \\ \models \varphi \wedge \psi & & \\ M, e_1 \models \neg\varphi & \text{sse} & \text{non } M, e_1 \models \varphi \\ M, e_1 \models K_a\varphi & \text{sse} & M, e_2 \models \varphi \text{ per ogni } e_2 \text{ tale che} \\ & & e_1 R_a e_2 \end{array}$$

**Definizione 8** [Assiomi e regole di inferenza] Il sistema di prova della logica epistemica che useremo è assiomatizzato utilizzando gli assiomi di **T** e la regola del *modus ponens* e della necessitazione come riportato in tabella 2:

Sistema	Regole	Assiomi	Proprietà di R
<b>T</b>	<b>MP</b> e <b>Nec</b>	$K_a(\varphi \rightarrow \psi) \rightarrow (K_a\varphi \rightarrow K_a\psi)$  $K_a\varphi \rightarrow \varphi$	R riflessiva

**Tabella 2** Logica epistemica **T**

**Definizione 9** [Struttura  $\mathcal{S}$ ] Una struttura  $\mathcal{S}$  è della forma  $\mathcal{S} = \langle \mathcal{F}, E, \mathcal{P}_E, M, D \rangle$ , dove

$\mathcal{F} = \{a, b, c, \dots\}$  è un insieme non vuoto di *framework* di estrazione automatica di metadati,

$E = \{e_1, \dots, e_m\}$  è un insieme non vuoto di possibili estrazioni ( $|E| = m \in \mathbb{N}$ ),

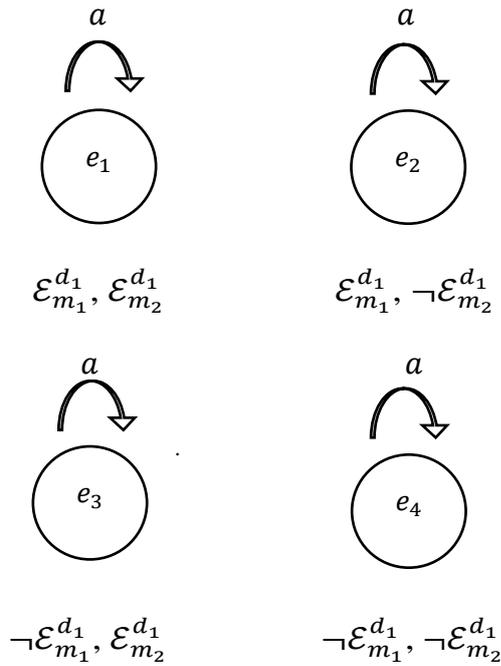
$\mathcal{P}_E = \{p_{e_1}, \dots, p_{e_m}\}$  è un insieme non vuoto di proposizioni ( $|\mathcal{P}_E| = m \in \mathbb{N}$ ),

$M = \{m_1, \dots, m_m\}$  è un insieme non vuoto di metadati ( $|M| = m \in \mathbb{N}$ ),

$D = \{d_1, \dots, d_m\}$  è un insieme non vuoto di documenti ( $|D| = m \in \mathbb{N}$ ).

$\mathcal{S}$  è una struttura nella quale occorrono possibili estrazioni  $E$ .  $\mathcal{F}$  è l'insieme dei *framework* di estrazione automatica di metadati, mentre  $\mathcal{P}_E$  è l'insieme delle proposizioni. Infine,  $M$  è l'insieme dei metadati e  $D$  l'insieme dei documenti (nel nostro caso di articoli scientifico-giuridici). All'interno della struttura possiamo rappresentare un modello relazionale usando un grafo in cui le

possibili estrazioni sono nodi e la relazione epistemica è indicata tramite frecce come illustrato in figura:



In questo grafo abbiamo una situazione nella quale dato un documento in entrata e due metadati, un agente di estrazione sa che si possono verificare quattro possibili estrazioni: l'estrazione in cui entrambi i metadati vengono estratti correttamente, l'estrazione in cui il metadato uno viene estratto correttamente mentre il due no, l'estrazione in cui il metadato due è estratto correttamente mentre l'uno no, ed infine l'estrazione in cui entrambi i metadati non sono riportati correttamente. Analizziamo, ora, più in dettaglio cosa vuol dire che in una estrazione una proposizione è vera o falsa. Come già sappiamo, la verità di una formula proposizionale dipende "dalla situazione del mondo", o nel caso di una proposizione epistemica "è vera in  $w$  a condizione che sia vera in tutti i mondi accessibili da  $w$ ". Le situazioni sono formalizzate usando valutazioni e in  $\mathcal{S}$  sappiamo che una proposizione  $p_{\mathcal{E}}$  "è vera in  $e$  a condizione che sia vera in tutte le possibili estrazioni accessibili da  $e$ "

*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

$$V^{\mathcal{P}\varepsilon}: E \rightarrow (\mathcal{P}_\varepsilon \rightarrow \{\text{vero}, \text{falso}\})$$

Inoltre, poiché sappiamo che  $p_\varepsilon$  ha la forma  $\mathcal{E}_{m_i}^{d_i}$  scriveremo che è vero (V) o falso (F) che “nell’estrazione  $e_i$  un framework estrae il metadato  $m_i$  dal documento  $d_i$ ” nel seguente modo

$$\underbrace{\mathcal{E}_{m_i}^{d_i}}_{e_i} = V/F$$

Ma cosa vuol dire che in una estrazione un metadato estratto è vero? Detto altrimenti, che cosa vuol dire che un *framework* estrae correttamente un metadato da un articolo scientifico-giuridico? Per rispondere a queste domande occorre presentare la teoria dei *truthmakers* e definire la veridicità come *truthmaker*. La teoria dei *truthmakers* è una interessante teoria metaontologica proveniente dal mondo della filosofia analitica che esplora la relazione tra ciò che è vero e ciò che esiste<sup>10</sup>. La teoria ha radici profonde nel pensiero occidentale e, da un lato, veicola una nostra intuizione emergente: se, ad esempio, è vero che il cane è sullo zerbino è perché il cane è “di fatto” sullo zerbino; dall’altro, rappresenta l’idea alla base di una celebre teoria della verità, ossia il corrispondentismo:

(C) Dire la verità è dire come “stanno le cose nel mondo”

La teoria dei *truthmakers* a cui ci rifaremo in questo lavoro è quella sviluppata dal filosofo australiano David Malet Armstrong in *Truth and Truthmakers*<sup>11</sup>:

(T) Per ogni verità,  $p$ , esiste un ente,  $T$ , tale che  $T$  rende vero  $p$  se e solo se non è possibile che  $T$  esista e  $p$  sia falso

---

<sup>10</sup> Il termine metaontologia – usato per la prima volta in [19]– indica l’indagine che mira a determinare quale sia il modo di caratterizzare la nozione di ontologia.

<sup>11</sup> [3]. Per una traduzione italiana delle principali opere di Armstrong si veda [7], mentre per un’introduzione al pensiero del filosofo australiano si veda [4]. Per un confronto tra la metaontologia del *Tractatus logico-philosophicus* di Wittgenstein e la teoria dei *truthmakers* di Armstrong rimandiamo a [5]. Infine, è importante specificare che: *i*) la relazione di *truthmaking* non è una relazione univoca: una proposizione può avere molti *truthmakers* e un oggetto può rendere vere molte proposizioni; e *ii*) ai fini di questo articolo non è necessario impegnarsi in una particolare ontologia dei *truthmakers* (come fatti, sostanze, proposizioni vere, ecc.).

Nel nostro dominio possiamo riformulare il principio  $\mathcal{T}$  armstronghiano come segue:

( $\mathcal{V}$ ) Per ogni proposizione vera,  $p_{\mathcal{E}}$ , esiste un documento,  $d$ , tale che  $d$  rende vero  $p_{\mathcal{E}}$  se e solo se non è possibile che  $d$  esista e  $p_{\mathcal{E}}$  sia falso

In questo modo dovremmo ormai essere in posizione di apprezzare pienamente il significato di questo principio. Consideriamo nuovamente il nostro schema

$$(1) \underbrace{\mathcal{E}_{m_i}^{d_i}}_{e_i} = V/F$$

Solleghiamo, ora, la tipica *truthmaking question*: in virtù di cosa (1) è vero? Ebbene, in base a  $\mathcal{V}$  dire che un *framework*  $a$  ha estratto correttamente un metadato  $m$  vuol dire che esiste un documento  $d$  che “rende vera” l’estrazione  $e$ .

## 4. Esempio di $\mathcal{S}$

Consideriamo ora in che modo possiamo modellare i metadati estratti da due differenti articoli scientifico-giuridici usando il framework CERMINE<sup>12</sup>. I metadati che terremo in considerazione negli esempi sono:  $m_1$  titolo,  $m_2$  autore e  $m_3$  rivista. Il primo documento  $d_1$  riguarda l’utilizzo di modelli bayesiani nell’ambito dell’argomentazione giuridica, invece, il secondo documento  $d_2$  analizza il percorso e le ragioni che hanno portato l’Unione Europea ad entrare in una nuova fase del costituzionalismo moderno (ossia il costituzionalismo digitale).

Posto un *framework*  $a$ , tre metadati  $m_1, m_2$  e  $m_3$  e due articoli scientifico-giuridici  $d_1$  e  $d_2$  avremo la seguente struttura  $\mathcal{S} = \langle \mathcal{F}, E, \mathcal{P}_{\mathcal{E}}, M, D \rangle$ :

$$\mathcal{F} = \{a\}$$

$$E = \{e_1, \dots, e_m\};$$

---

<sup>12</sup> Per l’estrazione dei metadati è stata utilizzata la risorsa gratuita online <http://cermine.ceon.pl/index.html>.

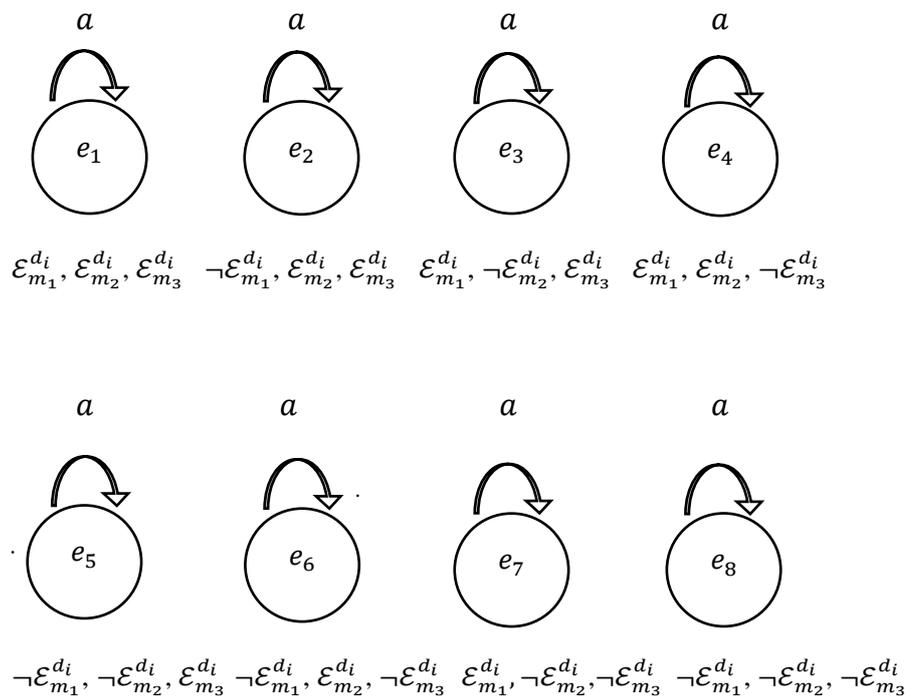
*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

$$\mathcal{P}_{\mathcal{E}} = \{p_{\mathcal{E}_1}, \dots, p_{\mathcal{E}_m}\}$$

$$M = \{m_1, m_2, m_3\}$$

$$D = \{d_1, d_2\}$$

In  $\mathcal{S}$  avremo il seguente universo di possibili estrazioni:



Con il primo documento  $d_1$  il *framework*  $a$  estrae correttamente tutti i metadati. Nella figura 3 il metadato “rivista” è evidenziato in giallo, il metadato “titolo” in verde e il metadato “autore” in rosso, mentre nella figura 4 è riportato l’XML dell’estrazione operata dal *framework*  $a$  con i relativi metadati evidenziati con gli stessi colori



Figura 3 Documento  $d_1$

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <article xmlns:xlink="http://www.w3.org/1999/xlink">
3   <front>
4     <journal-meta>
5       <journal-title-group>
6         <journal-title>Artificial Intelligence and Law</journal-title>
7       </journal-title-group>
8     </journal-meta>
9     <article-meta>
10      <title-group>
11        <article-title>Modelling competing legal arguments using Bayesian model compa
12      </title-group>
13      <contrib-group>
14        <contrib contrib-type="author">
15          <string-name>Martin Neil</string-name>
16          <xref ref-type="aff" rid="aff0"></xref>
17          <xref ref-type="aff" rid="aff1"></xref>
18          <xref ref-type="aff" rid="aff2"></xref>
19          <xref ref-type="aff" rid="aff3"></xref>
20        </contrib>
21        <contrib contrib-type="author">
22          <string-name>Norman Fenton</string-name>
23          <xref ref-type="aff" rid="aff0"></xref>
24          <xref ref-type="aff" rid="aff1"></xref>
25          <xref ref-type="aff" rid="aff2"></xref>
26          <xref ref-type="aff" rid="aff3"></xref>
27        </contrib>
28        <contrib contrib-type="author">
29          <string-name>David Lagnado</string-name>
30          <xref ref-type="aff" rid="aff0"></xref>
31          <xref ref-type="aff" rid="aff1"></xref>
32          <xref ref-type="aff" rid="aff2"></xref>
33          <xref ref-type="aff" rid="aff3"></xref>
34        </contrib>
35        <contrib contrib-type="author">
36          <string-name>Richard David Gill</string-name>
37          <xref ref-type="aff" rid="aff0"></xref>
38          <xref ref-type="aff" rid="aff1"></xref>
39          <xref ref-type="aff" rid="aff2"></xref>
40          <xref ref-type="aff" rid="aff3"></xref>
41        </contrib>

```

Figura 4 XML estrazione documento  $d_1$ <sup>13</sup>

<sup>13</sup><http://cermine.ceon.pl/cermine/task.html?jsessionid=E3B413EA00B92B9A1080DD1943910FD5?task=6507409268409556231>,  
<http://cermine.ceon.pl/cermine/download.html?type=nlm&task=6507409268409556231>.

*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

Poiché *a* estrae correttamente tutti i metadati si verifica l'estrazione possibile  $e_1$ :

- $\underbrace{\mathcal{E}_{m_1}^{d_1}}_{e_1} = V$
- $\underbrace{\mathcal{E}_{m_2}^{d_1}}_{e_1} = V$
- $\underbrace{\mathcal{E}_{m_3}^{d_1}}_{e_1} = V$

Invece, con il secondo documento  $d_2$  il *framework a* estrae correttamente due metadati su tre. Anche in questo caso, nella figura 5 il metadato “rivista” è evidenziato in giallo, il metadato “titolo” in verde e il metadato “autore” in rosso, mentre nella figura 6 è riportato l'XML dell'estrazione operata dal *framework a* con i relativi metadati evidenziati con gli stessi colori. Si noti come il metadato  $m_3$  “rivista” non sia estratto correttamente



**Figura 5** Documento  $d_2$

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <article xmlns:xlink="http://www.w3.org/1999/xlink">
3   <front>
4     <journal-meta>
5       <journal-title-group>
6         <journal-title>I</journal-title>
7       </journal-title-group>
8     </journal-meta>
9     <article-meta>
10      <article-id pub-id-type="doi">10.1093/icon/mbab001</article-id>
11      <title-group>
12        <article-title>The rise of digital constitutionalism in the European Union</a
13      </title-group>
14      <contrib-group>
15        <contrib contrib-type="author">
16          <string-name>Giovanni De Gregorio</string-name>
17          <xref ref-type="aff" rid="aff1"></xref>
18        </contrib>
19        <aff id="aff0">
20          <label></label>
21          <institution>John Danaher, The Threat of Algocracy: Reality</institution>
22          <addr-line>Resistance and Accommodation, 29(3) P</addr-line>
23        </aff>
24        <aff id="aff1">
25          <label></label>
26          <institution>PhD Candidate, University of Milano-Bicocca</institution>
27          <addr-line>Milan</addr-line>
28        </aff>
29        <country country="IT">Italy</country>
30        <institution>Academic Fellow, Bocconi University</institution>
31        <addr-line>Bocconic</addr-line>
32      </contrib-group>
33      <pub-date>2021</year>

```

Figura 6 XML estrazione documento  $d_2$ <sup>14</sup>

Poiché  $a$  estrae correttamente i metadati  $m_1$  e  $m_2$  ma non il metadato  $m_3$  si verifica l'estrazione possibile  $e_4$ :

- $\underbrace{\mathcal{E}_{m_1}^{d_2}}_{e_4} = V$
- $\underbrace{\mathcal{E}_{m_2}^{d_2}}_{e_4} = V$
- $\underbrace{\mathcal{E}_{m_3}^{d_2}}_{e_4} = F$

## 5. Conclusioni

In questo articolo, a partire dai metadati estratti da articoli scientifico-giuridici tramite sistemi di estrazione automatica di metadati basati sull'intelligenza artificiale, abbiamo usato la logica epistemica e la metaontologia per descrivere (modellare) formalmente i metadati estratti e definire un principio di veridicità come *truthmaker*. In particolare, in ambito giuridico un potenziale agente

<sup>14</sup><http://cermine.ceon.pl/cermine/task.html?jsessionid=99D5F9BDA2ED2DD6DCC3F9FF8EF5F9B1?task=3230711466372700550>,  
<http://cermine.ceon.pl/cermine/download.html?type=nlm&task=3230711466372700550>.

*A logical-metaontological approach to the problem of (meta)data veracity in systems for automatic extraction of metadata from scientific-legal articles*

automatico deve essere in grado di operare legittimamente, e di produrre decisioni e atti giuridicamente validi. Tuttavia, la capacità di istruire un agente automatico non può prescindere dalla costruzione di specifici modelli formali che consentano di utilizzare metadati ad alta veridicità e, conseguentemente, di ridurre il rischio di fondare le decisioni di un agente su informazioni non veritiere. In un mondo sempre più guidato dai dati, la modellizzazione logico-metaontologica dei metadati permette la costruzione di soluzioni per la sistematizzazione delle informazioni estratte e l'intersezione tra il diritto, la Data Science e l'intelligenza artificiale.

## **Riferimenti bibliografici**

- [1] B. M. Abdel-Karim, N. Pfeuffer e O. Hinz. Machine learning in information systems – a bibliographic review and open research issues. *Electronic Markets*, 2021,1-28. 2021.
- [2] C. Anderson. The End of Theory: The Data Deluge Makes the Scientific Method Obsolete. *Wired*. 2008. <http://www.wired.com/2008/06/pb-theory/>.
- [3] D.M. Armstrong. *Truth and Truthmakers*. Cambridge University Press, Cambridge. 2004.
- [4] F.F. Calemi. *Le radici dell'essere. Metafisica e metaontologia in David Malet Armstrong*. Armando Editore. Roma. 2013.
- [5] S. Cuconato. Mondi di Wittgenstein. Metaontologia del 'Tractatus' e teoria dei 'truthmakers' di Armstrong. *Rivista Italiana di Filosofia Analitica junior*, 5:2, Special Issue: Metaphysics, 53-65, 2014.
- [6] S. Cuconato. Epistemic logic for metadata modelling from scientific papers on COVID-19. *Science & Philosophy – Journal of Epistemology, Science and Philosophy*. 9 (2),83-96. 2021.
- [7] A. d'Atri. *Ritorno alla Metafisica*. Bompiani, Milano, 2012.
- [8] M. Frixione, S. Iaquinto e M. Vignolo. *Introduzione alle logiche modali*. Laterza, Roma-Bari, 2016.

- [9] S. Galvan. *Logiche intensionali. Sistemi proposizionali di logica modale, deontica, epistemica*. Franco Angeli, Milano, 1991.
- [10] E. Gettier. Is Justified True Belief Knowledge?. *Analysis*, 23,121-123, 1963.
- [11] M. G. Lozano, J. Brynielsson, U. Franke, M. Rosell, E. Tjörnhammara, S. Varga e V. Vlassov. Veracity assessment of online data. *Decision Support Systems*, 2020.
- [12] T. Lukoianova e V.L. Rubin. Veracity Roadmap: Is Big Data Objective, Truthful and Credible?, *Advances in Classification Research Online*, 4-15, 2014.
- [13] J. Pomerantz. *Metadata*, MIT Press, Cambridge, MA, 2015.
- [14] D. Snow. Adding a 4th V to BIG Data – Veracity, <http://dsnowondb2.blogspot.se/2012/07/adding-4th-v-tobig-data-veracity.html>, 2012.
- [15] D. Tkaczyk, P. Szostek, P. Jan Dendek, M. Fedoryszak e Ł. Bolikowski. CERMINE — automatic extraction of metadata and references from scientific literature. Conference: 2014 11th IAPR International Workshop on Document Analysis Systems, 2014.
- [16] D. Tkaczyk, P. Szostek, P. Jan Dendek, M. Fedoryszak e Ł. Bolikowski. CERMINE — automatic extraction of metadata and references from scientific literature. *International Journal on Document Analysis and Recognition (IJDAR)*, Springer, 2015.
- [17] G. Turbanti. *Logica e mondi possibili*. Pisa university press, Pisa, 2020.
- [18] H. van Ditmarsch, J. Halpern, W. van Der Hoek e B. Kooi. *Handbook of Epistemic Logic*. College Publications, 2015.
- [19] P. van Inwagen. *Meta-Ontology*. *Erkenntnis*, 48, 223–250, 1998.
- [20] T. Williamson. Gettier Cases in Epistemic Logic. *Inquiry: An Interdisciplinary Journal of Philosophy*, 56, 1-14, 2013.