

LQ *The Lab's Quarterly*

2018 / a. XX / n. 4 (ottobre-dicembre)



DIRETTORE

Andrea Borghini

COMITATO SCIENTIFICO

Albertini Françoise (Corte), Massimo Ampola (Pisa), Gabriele Balbi (Lugano), Matteo Bortolini (Padova), Massimo Cerulo (Perugia), Marco Chiappesi (Pisa), Franco Crespi (Perugia), Sabina Curti (Perugia), Gabriele De Angelis (Lisboa), Paolo De Nardis (Roma), Teresa Grande (Cosenza), Elena Gremigni (Pisa), Roberta Iannone (Roma), Anna Giulia Ingellis (València), Mariano Longo (Lecce), Domenico Maddaloni (Salerno), Stefan Müller-Doohm (Oldenburg), Gabriella Paolucci (Firenze), Massimo Pendenza (Salerno), Walter Privitera (Milano), Cirus Rinaldi (Palermo), Antonio Viedma Rojas (Madrid), Vincenzo Romania (Padova), Angelo Romeo (Perugia), Giovanni Travaglini (Kent).

COMITATO DI REDAZIONE

Luca Corchia (segretario), Roberta Bracciale, Massimo Cerulo, Cesar Crisosto, Elena Gremigni, Antonio Martella, Gerardo Pastore

CONTATTI

thelabs@sp.unipi.it

I saggi della rivista sono sottoposti a un processo di double blind peer-review.

La rivista adotta i criteri del processo di referaggio approvati dal Coordinamento delle Riviste di Sociologia (CRIS): cris.unipg.it

I componenti del Comitato scientifico sono revisori permanenti della rivista.

Le informazioni per i collaboratori sono disponibili sul sito della rivista:

<https://thelabs.sp.unipi.it>

ISSN 1724-451X



Quest'opera è distribuita con Licenza
Creative Commons Attribuzione 4.0 Internazionale

“The Lab’s Quarterly” è una rivista di Scienze Sociali fondata nel 1999 e riconosciuta come rivista scientifica dall’ANVUR per l’Area 14 delle Scienze politiche e Sociali. L’obiettivo della rivista è quello di contribuire al dibattito sociologico nazionale ed internazionale, analizzando i mutamenti della società contemporanea, a partire da un’idea di sociologia aperta, pubblica e democratica. In tal senso, la rivista intende favorire il dialogo con i molteplici campi disciplinari riconducibili alle scienze sociali, promuovendo proposte e special issues, provenienti anche da giovani studiosi, che riguardino riflessioni epistemologiche sullo statuto conoscitivo delle scienze sociali, sulle metodologie di ricerca sociale più avanzate e incoraggiando la pubblicazione di ricerche teoriche sulle trasformazioni sociali contemporanee.

2018 / a. XX / n. 4 (ottobre-dicembre)

Gli algoritmi come costruzione sociale

A cura di
Antonio Martella, Enrico Campo e Luca Ciccarese

Enrico Campo, Antonio Martella, Luca Ciccarese	<i>Gli algoritmi come costruzione sociale. Neutralità, potere e opacità</i>	7
SAGGI		
Massimo Airoidi, Daniele Gambetta	<i>Sul mito della neutralità algoritmica</i>	25
Chiara Visentin	<i>Il potere razionale degli algoritmi tra burocrazia e nuovi idealtipi</i>	47
Mattia Galeotti	<i>Discriminazione e algoritmi. Incontri e scontri tra diverse idee di fairness</i>	73
Biagio Aragona, Cristiano Felaco	<i>La costruzione socio-tecnica degli algoritmi. Una ricerca nelle infrastrutture di dati</i>	97
Aniello Lampo, Michele Mancarella, Angelo Piga	<i>La (non) neutralità della scienza e degli algoritmi. Il caso del machine learning tra fisica fondamentale e società</i>	117
Luca Serafini	<i>Oltre le bolle dei filtri e le tribù online. Come creare comunità "estetiche" informate attraverso gli algoritmi</i>	147
Costantino Carugno, Tommaso Radicioni	<i>Echo chambers e polarizzazione. Uno sguardo critico sulla diffusione dell'informazione nei social network</i>	173

LIBRI IN DISCUSSIONE

Irene Psaroudakis	Mario Tirino, Antonio Tramontana, <i>I riflessi di «Black Mirror»</i> . <i>Glossario su immaginari, culture e media della società digitale</i> , Roma, Rogas Edizioni, 2018, 280 pp.	203
Junio Aglioti Colombini	Daniele Gambetta, <i>Datacrazia. Politica, cultura algoritmica e conflitti al tempo dei big data</i> , Roma, D Editore, 2018, 360 pp.	209
Paola Imperatore	Safiya Umoja Noble, <i>Algorithms of Oppression: How Search Engines Reinforce Racism</i> , New York, New York University Press, 2018, 265 pp.	215
Davide Beraldo	Cathy O'Neil, <i>Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy</i> , New York, Broadway Books, 2016, 272 pp.	223
Letizia Chiappini	John Cheney-Lippold, <i>We Are Data: Algorithms and The Making of Our Digital Selves</i> , New York, New York University Press, 2017, 320 pp.	229



LA (NON) NEUTRALITÀ DELLA SCIENZA E DEGLI ALGORITMI

Il caso del machine learning tra fisica fondamentale e società

di *Aniello Lampo, Michele Mancarella, Angelo Piga**

Abstract

The impact of Machine Learning (ML) algorithms in the age of big data and platform capitalism has not spared scientific research in academia. In this work, we will analyze the use of ML in fundamental physics and its relationship to other cases that directly affect society. We will deal with different aspects of the issue, from a bibliometric analysis of the publications, to a detailed discussion of the literature, to an overview on the productive and working context inside and outside academia. The analysis will be conducted on the basis of three key elements: the non-neutrality of science, understood as the intrinsic relationship with history and society; the non-neutrality of the algorithms, in the sense of the presence of elements that depend on the choices of the programmer, which cannot be eliminated whatever the technological progress is; the problematic nature of a paradigm shift in favour of a science and a society dominated by data. The deconstruction of the presumed universality of scientific thought from the inside becomes in this perspective a necessary first step also for any social and political discussion. This is the subject of this work in the case study of ML.

Keywords

Neutrality, science, algorithms, machine learning, scientific paradigm, physics

* ANIELLO LAMPO lavora all'Internet Interdisciplinary Institute (IN3), Universitat Oberta de Catalunya (UOC), Barcelona (Spain). Email: nello.lampo@gmail.com

MICHELE MANCARELLA è Ricercatore indipendente. Email: mancarell@gmail.com

ANGELO PIGA lavora all'ICFO - Institut de Ciències Fotòniques di Barcelona. Email: angelo.piga@gmail.com

1. INTRODUCTION

«It's time to ask: what can science learn from Google?» si chiedeva provocatoriamente Chris Anderson, editor di Wired, in un celebre articolo del 2008 (Anderson, 2008). La domanda suggerisce di guardare al cuore del metodo scientifico. Costruire modelli teorici, falsificarli, eseguire dei test, è lo schema che ogni scienziato introietta ed utilizza, come risultato di secoli di epistemologia. L'idea di Anderson è che nell'era dei petabyte questo approccio risulterebbe obsoleto: nel XXI secolo, infatti, la mole di dati a disposizione sarebbe tale da poter estrarre previsioni da essi senza chiamare in causa modelli, dichiarando dunque “la fine della teoria”. Una “rivoluzione copernicana”, un cambio di paradigma, il cui principale perno tecnico-teorico risiederebbe nel *Machine Learning*, ovvero «la capacità di una macchina di imparare senza essere stata programmata per farlo» (Arthur, 1959). Quasi 10 anni dopo l'articolo di Anderson, Google ha investito nel solo 2016 tra i 20 e i 30 miliardi di dollari in ML ed in generale in Intelligenza Artificiale¹ (*Artificial Intelligence*, d'ora in poi AI), di cui il 90% in ricerca e sviluppo. Nello stesso periodo di tempo, assistiamo ad una crescita delle pubblicazioni scientifiche in questo campo².

La fisica non è esente da questa esplosione. Se l'interesse per tecniche avanzate di estrazione di informazione dai dati non dovrebbe sorprendere, sarebbe sbagliato provare a spiegare le ragioni del fenomeno appellandosi solo a criteri interni alla scienza, come ad esempio l'efficienza tecnica e la convenienza formale.

Al contrario, il boom del ML in fisica esemplifica il rapporto inestricabile tra il sapere scientifico ed il contesto storico nel quale esso viene prodotto. Una delle posizioni principali in accademia rispetto all'analisi di tale legame può essere condensata in un'unica parola: neutralità. Quest'espressione sottintende che l'attività dello scienziato, e dunque i suoi metodi ed i contenuti che produce, seguirebbero una progressione lineare, sconnessa dal contesto sociale e politico. Al contrario, in questo lavoro partiremo dal principio secondo cui la

¹ In letteratura, esistono diversi termini che si riferiscono alla possibilità di programmare una macchina per prendere decisioni autonomamente (cioè senza un modello predittivo). Le differenze sono sottili e spesso i termini sono usati come sinonimi, anche a sproposito. In particolare, il Machine Learning è una branca dell'Intelligenza Artificiale, la quale include altri settori come ad esempio la robotica. Nella sezione II introdurremo nel dettaglio altre differenze tecniche e terminologiche che ci serviranno nel corso dell'articolo. Quando non diversamente specificato ci riferiremo all'insieme di queste tecniche col nome Machine Learning.

² I dati verranno discussi in dettaglio nella sezione III.

scienza è a tutti gli effetti un'attività umana, ed in quanto tale viene esercitata in una data epoca storica, all'interno di un preciso sistema di pensiero, di produzione e, per arrivare alla quotidianità dello scienziato, nel quadro di un certo rapporto lavorativo. Di conseguenza la ricerca scientifica nel suo complesso riflette necessariamente questi aspetti del mondo circostante: in questo senso diremo che essa ha un carattere *non neutrale*.

La discussione intorno alla non-neutralità della scienza è stato uno degli argomenti principali del dibattito epistemologico tra fine ottocento e novecento ed oscilla fra l'idea estrema di una completa indipendenza della scienza dalla società, come ad esempio sostenuto dai convenzionalisti (Poincaré, 1905) e quella opposta dell'anarchismo epistemologico di Feyerabend (2002), passando per tutta una serie di posizioni intermedie rappresentate ad esempio da Popper (Popper, 1991), Lakatos (Lakatos, 1995) e Kuhn (Kuhn, 1973). In particolare, quest'ultimo propose un'analisi storica che evidenziava come la scienza si sviluppi secondo bruschi cambi di paradigma la cui affermazione dipende anche da fattori sociali.

L'obiettivo principale di questo lavoro è declinare tale dibattito nell'attualità assumendo come *case-study* l'utilizzo del ML nelle scienze pure. Proveremo a dimostrare che l'impiego di tali strumenti nell'ambito della fisica non può essere motivato prescindendo da una collocazione storica in quella che si sta configurando come "l'era dei dati e degli algoritmi".

Ciò richiede di quantificare se e in che misura un trend all'interno della ricerca in fisica si possa riscontrare; capirne le ragioni; discuterne i legami con il contesto produttivo e lavorativo. Parleremo in questo caso di "non-neutralità della scienza", nel senso chiarito poco sopra, mostrando che questo fenomeno non dipende da dinamiche interne alla scienza stessa, ma da fattori economici, sociali e politici.

In particolare, prenderemo in esame il contesto lavorativo accademico, pesantemente sovradeterminato dai tagli e le riforme imposte dopo la crisi del 2008, che portano a precarietà e riduzione delle posizioni permanenti. Questa situazione può condurre il ricercatore a scelte dei suoi programmi di ricerca dettate dalla necessità di acquisire competenze spendibili fuori dall'accademia piuttosto che dalle prospettive scientifiche. Vedremo come il boom del ML in accademia vada interpretato infatti alla luce della grande espansione del settore dei big data e dell'AI, che rappresenta oggi un settore industriale capace di attrarre forza lavoro con alti salari e investimenti.

Il rapporto della scienza "pura" con le dinamiche sociali e politiche

del tempo va letto però in un senso ulteriore oltre quello appena discusso. Ad esempio, ne *L'Ape e l'Architetto* (Ciccotti et al., 1976) si evidenzia «il ruolo sovrastrutturale che la produzione di scienza pura svolge in quanto forma specifica di cultura». Intendiamo riferirci in particolare al fatto che le scienze pure³ contribuiscono alla creazione di un sistema di linguaggio, di nozioni, di metodi, aspettative, legati alla possibilità di descrivere il mondo in termini di modelli matematici⁴ anche al di fuori di esse⁵.

Ciò costituisce una questione determinante al giorno d'oggi, in quanto spesso si considerano il “dato”, il modello matematico o l'algoritmo come elementi oggettivi, portatori di una “verità” imparziale, il cui impiego servirebbe a giustificare decisioni politiche come conseguenze logiche di un ordine naturale. Decostruire la pretesa universalità del “sistema scienza” *dall'interno delle scienze pure* diventa quindi un compito cardine per ogni discussione sociologica o politica che ruoti attorno all'utilizzo di strumenti scientifici. Un secondo obiettivo del lavoro è proprio l'inizio di tale decostruzione per quanto riguarda il ML e l'AI. Parleremo in questo caso, più specificamente, di “non-neutralità dell'algoritmo”, intendendo con questo concetto la presenza *all'interno dell'algoritmo stesso* di elementi che implicano scelte a discrezione del programmatore.

Tale discorso assume una rilevanza ancora maggiore accettando l'ipotesi di Anderson. Se rinunciamo alla teoria, cioè alla costruzione di *framework* predittivi che guidino le osservazioni, falsificabili o migliorabili, non resta infatti che ridursi all'idea secondo cui “correlation is enough”. Nell'uso di algoritmi per regolare la vita sociale e politica, questa affermazione va molto al di là del piano epistemologico - che pure è problematico di per sé. Elevare correlazione a strumento di previsione significa infatti amplificare, rafforzare e perpetuare lo *status quo*, in particolare i rapporti di forza e di potere all'interno della società.

Un ultimo obiettivo del lavoro è quindi problematizzare l'idea di un

³ Qui e in seguito useremo il termine “scienze pure” per indicare in particolare fisica e matematica. In altri casi, specificheremo la differenza.

⁴ Parliamo di “modelli” riferendoci a modelli predittivi, come nel metodo scientifico. In questo contesto, il claim di Anderson è volto a sostituire modelli con dati grezzi. Discuteremo in seguito la rilevanza di questo aspetto.

⁵ Un caso rilevante per la presente discussione viene dalla sociologia, di cui ad esempio Goldthorpe promuove una formulazione statistica in termini di teoria delle popolazioni (Goldthorpe, 2016), mentre Latour (Latour, 2010) denuncia come «Sociology has been obsessed by the goal of becoming a quantitative science» considerando irraggiungibile, se non inutile, questo obiettivo. Per una critica simile che prende direttamente piede dall'ipotesi di Anderson si veda anche (Crawford, 2011).

cambio di paradigma verso un'epoca completamente *data-driven*. Innanzitutto, l'analisi della letteratura nella fisica ci permetterà di valutare dall'interno se questo sia effettivamente il caso. D'altra parte, abbiamo appena visto che una prospettiva prettamente interna alla scienza non è sufficiente. Il concetto stesso di paradigma ci impone di mettere in relazione la scienza con il contesto storico. In altre parole, prenderemo *sul serio* l'affermazione di Anderson, mostrando come, sebbene non giustificata limitando lo sguardo alla sola fisica, essa acquisisce un senso diverso prendendo in considerazione "ciò che possiamo imparare da Google", cioè l'aura di universalità e legittimità che gli algoritmi (e il ML in particolare) acquisiscono nell'era dei big data. Concluderemo che questa prospettiva mette lo scienziato di fronte al dovere di assumere e comunicare con chiarezza come la scienza e i suoi strumenti siano attività umane e per questo criticabili e mai del tutto neutrali. La discussione del caso della fisica presentata in questo lavoro vuole essere un passo in questa direzione.

La struttura dell'articolo è la seguente. Nella sezione I, discuteremo il quadro alla base di questo studio: la non-neutralità della scienza. Nella sezione II, introdurremo il concetto di ML e alcuni elementi tecnici e storici che motivano l'esplosione del suo utilizzo negli ultimi 10 anni. Passeremo quindi a studiare l'utilizzo di algoritmi di ML in fisica, nella sezione III. Presenteremo innanzitutto una stima del fenomeno sulla base di un'analisi bibliometrica, per poi discutere nel dettaglio alcuni casi concreti: la fisica delle alte energie (sezione IIIa), l'astrofisica (sezione IIIb) e la fisica delle basse energie (sezione IIIc). Vedremo che l'analisi storica, bibliometrica e della letteratura conducono a considerare il ruolo delle grandi piattaforme e del mercato del lavoro. Questo sarà l'oggetto della sezione IV. Infine, nella sezione V saremo in grado di sviluppare un parallelo tra la fisica e altri campi, e di commentare i concetti di non-neutralità della scienza e degli algoritmi e il problema del cambio di paradigma.

2. LA NON-NEUTRALITÀ DELLA SCIENZA

Con non-neutralità della scienza intendiamo l'impossibilità di una demarcazione netta fra la struttura della scienza – ovvero le sue teorie, il cosiddetto metodo scientifico e la sua organizzazione interna – da un lato e il contesto storico, la struttura della società e i rapporti di forza interni ad essa, dall'altro.

Il lavoro di molti filosofi della scienza si concentra sulla possibilità di isolare e definire norme secondo cui lo scienziato dovrebbe

procedere: è il problema del “metodo”. Ancora oggi la pretesa *oggettività* della scienza, il suo valore e la sua funzione sociale sono spesso difese sulla base di una pretesa universalità del metodo scientifico. L’attuale visione scientifica è costruita sulla base di tale metodo, a partire da Galileo e Bacone, passando per Newton e poi Einstein, fino alle moderne misure che completano la conoscenza del modello standard delle particelle elementari.

Già Galileo era consapevole dell’importanza della presenza di una teoria speculativa che guidasse l’osservazione e nel *Dialogo* esprime la sua ammirazione per gli scienziati capaci di «anteporre quello che il discorso gli dettava, a quello che le sensate esperienze gli mostravano apertissimamente in contrario». Secoli dopo, passando attraverso una lunga storia che vede protagonisti filosofi come Descartes, Locke, Kant, i “positivisti” dell’Ottocento, i “convenzionalisti” di inizio Novecento, fino al neopositivismo, sono probabilmente Karl Popper e Thomas Kuhn ad aver affermato un lessico e un approccio divenuti standard per il lavoro degli scienziati.

In particolare, nel “razionalismo critico” popperiano (Popper, 1991) il procedere della scienza avviene per “prova ed errore”, avanzando ipotesi e teorie e testandone le previsioni. Tali test possono condurre unicamente a falsificazioni e mai a una verifica, dato che lo scienziato non può essere certo che non esistano altri fatti, ancora sconosciuti, che invalidano la teoria. Il falsificazionismo popperiano ha avuto un impatto talmente forte nel dibattito epistemologico da diventare il discriminante che tuttora definisce la scienza, e in particolare l’elemento che la eleva al di sopra di altre discipline che non portano in sé la possibilità di correggere i propri errori. In questo concetto oggi si cela la supposta oggettività della scienza, che si configura come prototipo esemplare dell’attività razionale umana, la strada più sicura, se non verso la certezza, almeno verso la verità secondo un progresso lineare. Nella visione di Popper le idee cambiano secondo criteri tutti interni alla scienza, il cui sviluppo risulta dunque cumulativo e svincolato dalla storia della società: in questo senso potremmo dire *neutrale*.

Un elemento importante consiste appunto nello sfondare la barriera tra storia e scienza, tra società, politica e progresso lineare del pensiero scientifico. È soprattutto a Lakatos e Kuhn che si deve una svolta in questa direzione. Ogni standard o modello che ha dominato la cultura scientifica nel corso della storia è in realtà il prodotto di una “visione del mondo” che trascende la sola scienza ed è dettato dall’epoca storica e da un insieme di credenze e di sistemi di pensiero. È questo che Kuhn definisce un *paradigma* (Kuhn, 1972). La conoscenza scientifica segue

un andamento alternato tra lunghi periodi di “scienza normale”, fatti di accumulazione di dati seguendo criteri interni alla scienza e a un paradigma vigente, e “rivoluzioni” che portano alla sostituzione di un paradigma con un altro, in cui la frontiera tra il “fuori” e il “dentro” la scienza cade. Diversi paradigmi sono incommensurabili e non è possibile scartarne uno in base alle prescrizioni del falsificazionismo.

C’è di più: i dati sperimentali candidati alla verifica logica di una teoria sono sempre contaminati dalle assunzioni insite in un paradigma, e i dati stessi sono sempre *theory-laden* (carichi di teoria).

È Paul Feyerabend a spingersi oltre (Feyerabend, 2002). In effetti, nessuna teoria scientifica è mai in accordo con *tutti* i fatti noti del suo tempo, e i canoni di “razionalità scientifica” non esistono se svincolati dal contesto storico. L’unica soluzione è allora scrivere “contro il metodo”, che significa riconoscere che gli standard sono «una guida che a sua volta è parte di un’attività guidata ed è mutata da questa»: cioè l’osmosi tra scienza, epoca, scelte politiche è continua, totale e inevitabile. È in questi termini che si deve parlare di scienza.

In questo quadro, il risultato principale è che la scienza va considerata come un’attività umana, il cui sviluppo non si può scindere dal percorso storico regolato da fattori sociali e politici. Nel resto del testo specializzeremo questo punto di vista al lavoro del fisico teorico degli ultimi anni, ed in particolare analizzeremo il caso dell’exploit nell’uso del ML.

3. DAL MACHINE LEARNING ALLE DEEP NEURAL NETWORKS

In questa sezione, introdurremo alcuni concetti e definizioni che ci saranno utili per inquadrare il problema e per la discussione seguente.

Come già accennato nell’introduzione, il termine *Machine Learning* (ML) si riferisce alla «capacità di una macchina di apprendere da sola dai dati, senza essere stata esplicitamente programmata per farlo»⁶. I primi lavori di ML risalgono agli anni 50, con semplici modelli ad output binario (attivo/inattivo) per spiegare il funzionamento dei neuroni. Uno di questi modelli è il *Perceptron* (Rosenblatt, 1958). Venne rapidamente realizzato che diversi Perceptrons (anche chiamati “neuroni”) possono essere combinati in un *layer* per fornire un output a più classi – il “layer di output”. Tale struttura è l’esempio più semplice di *rete neurale* o *Neural Network* (NN). Nel caso delle NN, i layers

⁶A S. Arthur (Arthur, 1959) viene attribuito il conio del termine Machine Learning. Per un’introduzione storica, si veda l’introduzione in (Goodfellow 2016). Per un’introduzione formale e rigorosa del problema, si veda ad es. (Mostafà, 2012).

possono inoltre essere combinati in successione per raggiungere livelli più alti di precisione e astrazione (si parla di “*hidden layers*”, “*layer nascosti*”, per i layers più interni). Per quanto quello delle NN sia un esempio di ML, esistono molti algoritmi di ML che non si basano su reti neurali. In seguito, ci riferiremo con ML al secondo tipo di algoritmi e con NN al primo.

È solo negli anni 2000 che la possibilità di sviluppare reti con molti layers viene sfruttata fino in fondo, per tre motivi. Innanzitutto, la disponibilità di dati diventa esponenzialmente più grande e molto più varia⁷. Il secondo motivo è l'aumento, anche questo di diversi ordini di grandezza, della potenza di calcolo, in particolare grazie all'utilizzo delle cosiddette *Graphic Processing Unit* (GPU) (Raina 2009). Questi due progressi hanno inoltre portato – terzo motivo – alla necessità di affinare alcuni algoritmi noti per renderli più performanti (Hinton, 2012; Maas, 2013).

Grazie a questi tre elementi chiave (volume di dati-potenza di calcolo-nuovi algoritmi), si inizia a parlare di *Deep Learning* (DL), o di *Deep Neural Networks* (DNN, reti neurali profonde), riferendosi proprio alla possibilità di costruire reti a molti layer.

Nella prossima sezione, studieremo innanzitutto la presenza dei tre termini sopra introdotti nella letteratura scientifica, e commenteremo in dettaglio dei casi specifici.

4. IL MACHINE LEARNING NELLA FISICA FONDAMENTALE

In questa sezione discuteremo l'utilizzo del ML nella ricerca. Cercheremo di dare una stima della portata del fenomeno, e mostreremo che si riscontra un'impennata di pubblicazioni a riguardo in tutti i settori dopo il 2010 circa. Ci concentreremo poi sulla fisica e su alcune branche particolarmente significative, analizzando in dettaglio dei casi specifici.

4.1. *L'impennata delle pubblicazioni: i dati*

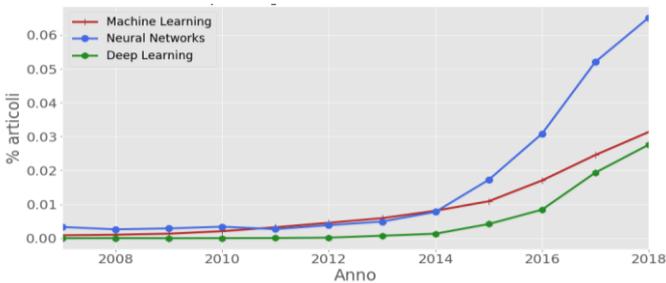
Nella letteratura scientifica, l'archivio open-access online

⁷ Consideriamo alcuni dataset di riferimento: “Modified National Institute of Standards and Technology” (MNIST) è il più largo degli anni 90 e conta 60 mila cifre scritte a mano. Oggi lo Street View House Numbers (SVHN) – una raccolta di immagini di numeri messa a disposizione da Google Street View – è di un ordine di grandezza più largo. Il dataset ImageNet conta oltre 14 milioni di immagini, contenenti gli oggetti più diversi. Infine, il dataset rilasciato dal “Workshop of Machine Translation” (WMT) per la traduzione automatica arriva a 1 miliardo di dati.

www.arxiv.org è lo spazio più largamente utilizzato per condividere articoli o preprint. Il sito conta, nel 2018, più di 10mila nuovi articoli al mese⁸. Dato il suo largo utilizzo, una ricerca tra gli articoli contenuti nel database rappresenta un mezzo utile per avere una stima della presenza di un trend. Abbiamo effettuato quindi una ricerca degli articoli che contengano le espressioni “Machine Learning”, “Neural Networks” o “Deep Learning” nel titolo o nell’abstract. I dati presentati nelle figure che seguono sono stati raccolti il 14 Dicembre 2018.

Il grafico in Figura 1 mostra la percentuale di articoli contenenti ML (linea rossa), DL (linea verde) e NN (linea blu) nell’abstract o nel titolo. È immediato notare una netta impennata a partire dal 2014, preceduta da un lieve incremento a partire dal 2010 circa. Il numero di articoli a tema NN è passato da 113 nel 2008 a 4575 nel 2017, pari ad un incremento totale di 40 volte e relativo di poco più di 30. Nel 2018 il numero di articoli supera gli 8000. Simili stime valgono per gli altri due termini, con un incremento relativo pari a circa 55 volte per ML. Il termine DL è addirittura assente prima del 2008, mentre conta più di 1500 articoli nel 2017. Come si capisce dal grafico, tale cambiamento di tendenza nel numero assoluto di articoli non è spiegabile con il generale incremento degli articoli inviati ad Arxiv, che è invece aumentato di solamente il doppio, da 4970 articoli inviati nel mese di gennaio 2008 a 10332 articoli nel mese di dicembre 2017.

Fig. 1. Percentuale di articoli Arxiv con “Machine Learning”, “Deep Learning” o “Neural Networks”) in titolo o abstract



È interessante chiedersi anche quali campi di ricerca risentono maggiormente di questa tendenza e in seguito concentrarsi sulla fisica. Arxiv divide gli articoli in 8 categorie: Fisica, Informatica (Computer

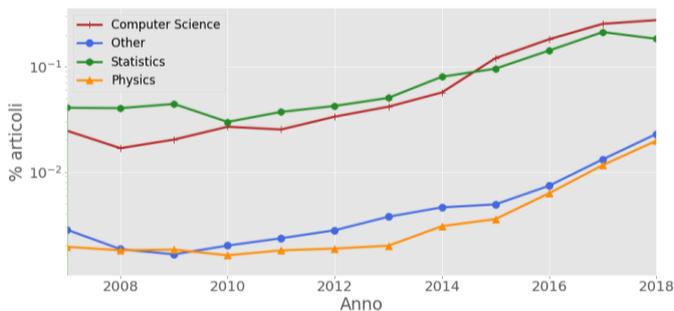
⁸ https://arxiv.org/stats/monthly_submissions

Science), Matematica, Finanza Quantitativa, Statistica, Biologia Quantitativa, Economia, Energia Elettrica e Scienza dei Sistemi.

In Figura 2, mostriamo il numero di articoli su Arxiv con una qualsiasi delle espressioni ML, DL e NN nel titolo o nell'abstract, divise per categorie. I dati sono normalizzati rispetto al numero totale di articoli usciti nello stesso anno nella stessa categoria e sono visualizzati in scala logaritmica. È chiara quindi una crescita esponenziale delle ricerche con queste parole chiave. Com'era prevedibile, la "Computer Science" (linea rossa), domina il trend nei valori assoluti degli articoli pubblicati, con una significativa impennata dopo il 2014. Tra gli altri campi, seguono la "Statistica" (linea verde), e proprio la "Fisica" (linea arancione), per cui l'incremento è ancor più significativo (circa 10 volte tra il 2008 e il 2018). La linea blu tratteggiata rappresenta la somma di tutte le altre sezioni tematiche dell'archivio ed anche in questo caso si nota una forte crescita che dimostra la trasversalità nell'utilizzo di queste nuove tecniche pressoché in tutti gli ambiti della fisica.

Il periodo dell'impennata, in particolare in Informatica, segue un momento di svolta nei progressi nel campo del riconoscimento di immagini. Nel 2012, una DNN ha vinto una celebre competizione, la *ImageNet Large Scale Visual Recognition Challenge*, con un margine di oltre il 10% in più di precisione (Krizhevsky, 2012) rispetto al secondo classificato. Questo avvenimento è oggi considerato come il vero e proprio inizio della "rivoluzione del Deep Learning" ed ha convinto definitivamente le maggiori piattaforme (Google, Facebook, Microsoft, etc.) a investire massicciamente in progetti di ricerca e sviluppo sulle reti neurali. L'ipotesi, che confermeremo in seguito con esempi specifici, è che anche la fisica abbia risentito di questa svolta.

Fig. 2. Percentuale di articoli Arxiv con "Machine Learning", "Deep Learning" o "Neural Networks") in titolo o abstract, divisi per categoria

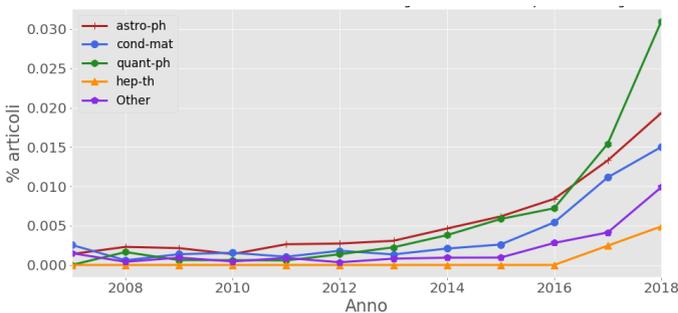


Restringiamoci ora alla sola categoria “Fisica”. Questa è a sua volta ripartita in 32 diverse sottocategorie⁹. In Figura 3, riportiamo il grafico delle percentuali di articoli nella categoria “Fisica”, contenenti una delle tre espressioni ML, DL e NN in titolo o abstract, e divisi per sottocategoria.

I dati mostrano che le categorie dominanti sono “Astrofisica” (astro-ph), “Fisica della Materia Condensata” (cond-mat), “Fisica Quantistica” (quant-ph) e “Fisica teorica delle Alte Energie” (hep-th), con un significativo incremento post-2014¹⁰. Discuteremo questi tre casi in dettaglio in seguito.

Anche se non analizzeremo tutti i casi in dettaglio in questo lavoro, è significativo notare l’incremento nei campi più disparati. Nella figura riassumiamo tutti i campi minori in un’unica curva indicata come “other”, che sembrano segnalare, come anticipato in precedenza, un interesse crescente per l’argomento che sembra partire da ambiti dove la sua applicazione è più immediata (come ad esempio l’astrofisica, come spiegheremo), ad ambiti in cui l’applicazione è più pionieristica (come la fisica teorica o la diagnosi di tumori in fisica medica). In effetti, tutti i settori ad eccezione dei primi due hanno un numero di pubblicazioni praticamente nullo fino a dopo il 2010. I dati di Arxiv presentati in questa sezione confermano, come ipotizzato, l’esistenza di un trend che eccede notevolmente la crescita del numero totale di articoli.

Fig. 3. Percentuale di articoli Arxiv con “Machine Learning”, “Deep Learning” o “Neural Networks”) in titolo o abstract nella categoria “Fisica”, divisi per sottocategoria



⁹ Una lista completa può essere consultata a <https://arxiv.org/archive/physics>.

¹⁰ Per rendere il grafico leggibile ed escludere informazione poco significativa, abbiamo scelto di riportare nel grafico solo le sottocategorie che hanno avuto un incremento di almeno un fattore 5 tra il 2008 e il 2018, e che abbiano avuto almeno 10 articoli nel 2018.

4.2. *La fisica delle alte energie*

Tecniche statistiche e algoritmi di previsione avanzati sono fondamentali nell'analisi dati. La fisica sperimentale non fa eccezione. In questa sezione discuteremo il caso più noto di una grande collaborazione in Fisica delle Alte Energie: l'acceleratore di particelle *Large Hadron Collider* (LHC) del CERN di Ginevra.

Una review apparsa su *Nature* (Radovic, 2018) spiega come il ML sia usato da lungo tempo nella collaborazione, ma segnala una svolta importante a partire dal 2014, con l'investimento massiccio nello studio di tecniche di DL. Tra le ragioni si citano esplicitamente il progresso nel riconoscimento immagini da parte di Google e la competizione *ImageNet* del 2012.

Del 2014 è anche il lancio di una competizione sulla piattaforma *Kaggle.com*, la *Higgs Boson Machine Learning Challenge* (Adam-Bourdarios, 2014), volta allo sviluppo di algoritmi per rivelare il bosone di Higgs (la famosa particella alla base del meccanismo che dà massa alla materia, al centro del premio nobel del 2013), vinta da una combinazione di 70 reti neurali, ciascuna con 3 layers di 600 neuroni¹¹.

Ancora del 2014 è un articolo che introduce a tutti gli effetti il DL in fisica delle alte energie, analizzandone alcune differenze sostanziali con le tecniche di ML già utilizzate (Baldi, *et al.*, 2014). Gli autori mettono l'accento sulla profonda trasformazione concettuale che il DL apporta. Gli algoritmi di ML infatti funzionano efficacemente solo con dati già preparati dai ricercatori sulla base della conoscenza del processo fisico in gioco. Questa procedura è nota come *feature engineering*. Ad esempio, l'impulso delle particelle rivelato a LHC permette di ricostruire la "massa invariante" del processo, una grandezza fisica ad alto potere discriminante (*high-level feature*), che è usata come input per il ML. Questa fase può essere estremamente complicata o computazionalmente costosa. Con una rete neurale profonda, si possono invece usare direttamente i dati "grezzi" del rivelatore (*low-level features*), raggiungendo una precisione equivalente: è la rete stessa che indipendentemente recupera l'informazione che prima veniva ottenuta tramite la conoscenza dei processi fisici in gioco.

Tuttavia, il *reverse engineering* di una NN, cioè la ricostruzione del processo decisionale che ha condotto all'output, è ad oggi impossibile: nel secondo caso, non possiamo quindi estrarre l'equivalente delle *high-level features*, e rinunciamo alla conoscenza del processo fisico in favore del potere predittivo. La transizione qui esemplificata, verso

¹¹ <https://github.com/melisgl/higgsmachinelearning/blob/master/doc/model.md>

un'analisi "cieca", è la cifra della rivoluzione insita nel DL e la ragione per cui il concetto di "cambiamento di paradigma" viene spesso evocato recentemente anche in fisica. In realtà il problema è più complesso e un modello falsificabile rimane comunque in gioco. Una rete neurale ha infatti bisogno di una fase cosiddetta di *training* (il suo "allenamento"), in cui l'algoritmo impara da una porzione di dati (*training set*) in cui è noto (in questo caso) se un segnale sia presente o no, per poi essere in grado di fornire previsioni. Il *training* richiede centinaia di migliaia di dati per essere efficace: per allenare la rete, il segnale cercato viene quindi calcolato sulla base del modello teorico da testare (nel caso di LHC, può essere ad esempio il Modello Standard – ad ora il più importante e preciso modello che prevede e cataloga le particelle fondamentali della natura – come sue estensioni), e i dati di allenamento generati attraverso simulazioni numeriche. Alla base del processo rimangono dunque un modello predittivo e la sua capacità di fornire previsioni falsificabili.

Un altro caso rilevante è l'esperimento NOvA per la rivelazione di neutrini (Ayes, 2007). In particolare, il tipo di rete utilizzata è basata sul modello *GoogLeNet* (Szegedy, 2014), vincitore della competizione *ImageNet* 2015, come esplicitamente spiegato dagli autori.

Per finire, si può citare lo studio delle proprietà di un'altra particella elementare, il quark beauty, a LHCb, in cui si è iniziato ad utilizzare modelli basati sull'algoritmo usato da Google per il suo traduttore (CERN, 2017).

4.3. *L'astrofisica*

In Figura 3 si vede che astro-ph è una delle sottocategorie con il maggior incremento in numero di pubblicazioni.

Questo dato non è sorprendente, visto che un telescopio fornisce immagini ottiche, e il riconoscimento di immagini è cruciale per estrarre efficacemente informazione. Come caso esplicito, consideriamo la rivelazione di *lenti gravitazionali forti*¹². Ad oggi l'analisi accurata di una sola lente gravitazionale può richiedere anche settimane di lavoro e una conoscenza molto approfondita dei processi fisici che intervengono nella sua osservazione. Un recente articolo apparso su Nature (Hezaveh,

¹² Le lenti gravitazionali forti sono un fenomeno previsto dalla teoria della relatività generale, che consiste nella formazione in un telescopio di più immagini della stessa sorgente luminosa, a causa della deflessione della luce da parte di strutture gravitanti lungo il percorso che essa compie per arrivare all'osservatore.

2017) mostra che usando delle NN¹³, il processo di rivelazione è circa 10 milioni di volte più veloce che nel caso standard ed ha una precisione paragonabile. Soprattutto, però, in modo analogo all'esempio di LHC, questo approccio non richiede alcuna conoscenza specifica dei processi fisici in gioco nella misura: i dati delle immagini del telescopio passano direttamente all'algoritmo. Precedentemente, la misura richiedeva l'inclusione di centinaia di migliaia di parametri supplementari per modellare il segnale e complesse tecniche statistiche per estrarre il risultato finale.

Un altro esempio è la rivelazione di *onde gravitazionali* (Abbott, 2016), in cui studi recenti (Gabbard, 2018) mostrano come sia possibile rivelare un segnale a partire dai dati grezzi, raggiungendo la stessa performance (o addirittura migliorandola) della lunga e complessa tecnica standard di analisi dati.

Altri casi di applicazioni alla frontiera dell'astrofisica moderna sono lo studio di *esopianeti* (Lam, 2018; Shallue, 2018) e gli esperimenti cosmologici di prossima generazione (Gillet, 2018).

Gli esempi discussi, come nel caso delle alte energie, vanno nella direzione di una scienza sempre più data-driven, rispetto a una fase precedente in cui la modellizzazione di tutte le componenti di un esperimento era determinante.

Tuttavia, anche in questo caso l'allenamento della rete richiede dei dataset di allenamento provenienti da simulazioni numeriche risultato di complessi codici che forniscono previsioni sulla base dei modelli che sono da testare nell'esperimento¹⁴. In ogni caso, tale processo è determinante per l'utilizzo delle NN che vengono così ridimensionate a efficaci *tool* per l'analisi dei dati.

4.4. *La meccanica quantistica e la fisica della materia condensata*

Gli ultimi due ambiti di ricerca su cui ci soffermiamo sono la fisica della materia condensata (cond-mat) e la meccanica quantistica (quant-ph).

A differenza della fisica delle particelle e dell'astrofisica, in cui i casi che abbiamo analizzato si riferiscono prevalentemente all'uso di tecniche di DL per l'analisi dati, la meccanica quantistica e la materia

¹³ Ancora una volta, una di queste reti è AlexNet, vincitrice della competizione ImageNet del 2012.

¹⁴ Nel caso delle onde gravitazionali, ad esempio, occorre risolvere le equazioni della relatività generale per un sistema di oggetti compatti orbitanti, il che è risultato di tecniche numeriche combinate con accurate approssimazioni analitiche. Nel caso del lensing gravitazionale, esistono codici di simulazione di immagini di galassie che combinano modelli a immagini già esistenti.

condensata sono il campo di sperimentazione di nuove frontiere. Elenchiamo le più importanti¹⁵.

Una prima menzione va fatta per il cosiddetto *quantum machine learning* (QML), cioè l'uso di algoritmi che sfruttano le enormi potenzialità del calcolo quantistico per risolvere problemi classici di ML in tempi e costi computazionali esponenzialmente minori. Il QML è un campo nato di fatto nel 2013 da una collaborazione tra MIT e Google¹⁶. Nel 2013 Google (in collaborazione con la NASA) acquista anche un esemplare del prototipo di computer quantistico D-Wave, giustificando l'investimento di 10 milioni di dollari proprio per le possibili applicazioni in ML¹⁷. Va notato tuttavia come il computer quantistico sia lontano dall'essere realizzato e il QML rimanga un ambito prettamente teorico. D'altra parte, tecniche di ML stanno venendo utilizzate proprio per migliorare alcuni processi sperimentali utili per la costruzione del computer quantistico. In questo senso, il ML è un utile strumento più che un cambio di paradigma (Biamonte, 2017).

Concludiamo la sottosezione con uno sguardo alle importanti applicazioni allo studio delle *transizioni di fase quantistiche*. Una delle categorie di problemi tradizionali è studiare le proprietà di atomi o gas di atomi interagenti a una temperatura ambiente spesso vicina allo zero assoluto. Le equazioni da risolvere per avere una completa conoscenza delle sue fasi e dell'evoluzione dinamica hanno immensi costi computazionali e metodi numerici o approssimazioni sono necessari. Si cerca quindi di usare il ML per classificare le fasi del sistema. L'approccio è recentissimo, e risale solo agli ultimi tre anni. Nel caso di algoritmi di *supervised learning* (Carrasquilla, 2017) – cioè in cui sia disponibile un training set in cui le fasi siano note – il ricercatore sa già quali sono le fasi del sistema e cerca di ottenere l'esatto punto della transizione di fase (ad esempio una temperatura).

Il metodo cosiddetto di *unsupervised learning* (Van Nieuwenburg, 2017) è più interessante nel senso di un possibile cambio di paradigma, perché non richiede la conoscenza pregressa delle fasi. In questo caso l'algoritmo è in grado di dividere automaticamente i dati in gruppi di affinità. Tuttavia, è solitamente necessaria una *feature engineering* altamente sofisticata che da sola potrebbe condurre al risultato voluto, senza necessariamente ricorrere al ML (Huembeli, 2018). Un altro

¹⁵ Per una rassegna tecnica si veda la recente review (Dunjko, 2018)

¹⁶ Vedi (Lloyd, 2013). Gli autori sono Seth Lloyd, guru mondiale dell'informazione quantistica e del computer quantistico, Patrick Rebentrost (un suo giovane collaboratore), entrambi affiliati MIT, e Masoud Mohseni, senior research del team Google AI.

¹⁷ Google and NASA Launch Quantum Computing AI Lab: <https://www.technology-review.com/s/514846/google-and-nasa-launch-quantum-computing-ai-lab/>

problema risiede nel fatto che spesso risulta necessario sapere in anticipo quante sono le fasi del sistema. Questo impedisce quindi di trovare fasi nuove e rende comunque necessaria una conoscenza qualitativa, ma profonda, del sistema stesso.

5. PIATTAFORME E ACCADEMIA

In più di un caso discusso nelle sezioni precedenti è emerso esplicitamente un legame tra la ricerca sviluppata da grandi piattaforme come Google e il suo successivo impiego in fisica. Questo è indicativo della necessità di guardare al di fuori della fisica per comprendere fino in fondo il fenomeno. In questa sezione daremo una misura quantitativa di quale sia il trend degli investimenti privati nell'ambito dei big data, mostrando un exploit che fa da attrattore nel mercato del lavoro. Successivamente mostreremo come nell'ambito accademico si assista da diversi anni ad una costante diminuzione di fondi e possibilità di accesso a posizioni permanenti: una situazione che porta naturalmente i giovani ricercatori a spostarsi spontaneamente su ambiti che in termini di competenze acquisite hanno maggiori prospettive lavorative in settori privati esterni all'accademia e alla ricerca di base.

Le grandi piattaforme, tra i maggiori attori del capitalismo mondiale, basano oggi i loro profitti quasi esclusivamente su dati, algoritmi e risorse di calcolo: ad esempio, oltre il 95% dei profitti di Facebook proviene dalla vendita di spazi pubblicitari "adattati" ai gusti dell'utente¹⁸; l'80% degli streaming di Netflix proviene da raccomandazioni piuttosto che da ricerche dirette (Gomez-Uribe, 2016). Al di là delle maggiori piattaforme, questo trend si estende ai settori più disparati. Secondo un rapporto di J.P. Morgan's¹⁹, tra quelli più rilevanti in cui il ML si rivela essenziale ci sono banche e sicurezza digitale, finanza²⁰, usi governativi, comunicazione e servizi, assicurazioni, trasporti, sanità, educazione.

Una stima della dimensione di questo mercato – fornita sempre da J.P. Morgan's – parla di 58 miliardi di dollari entro il 2021, a partire da circa 12 miliardi nel 2017; questo rappresenta uno dei settori tecnologici in più forte crescita, con un tasso di crescita annua composto (CAGR) di quasi 50%. Il 44% delle imprese oggetto dello studio indica il settore di

¹⁸ <https://www.statista.com/statistics/267031/facebooks-annual-revenue-by-segment/>

¹⁹ <https://flamingo.ai/wp-content/uploads/2017/11/JPMorganAnInvestorsGuideToArtificialIntelligencev2.pdf>

²⁰ Nel maggio 2017, J.P. Morgan ha pubblicato il più grande report esistente su Big Data e Machine Learning in finanza (J.P. Morgan, 2017)

ML/AI come la tecnologia che avrà il maggior impatto sull'impresa nella prossima decade, mentre un'altra indagine condotta da O'Reilly Media su un diverso campione²¹ dà lo stesso settore al 61%.

Da questo quadro si capisce l'entità dello studio e dell'applicazione del ML negli ultimi 5-10 anni al di fuori dell'accademia. Di tale tendenza si trova una chiara traccia anche nei dati che riguardano il mercato del lavoro.

Secondo il rapporto *LinkedIn's 2017 US Emerging Jobs Report* pubblicato dal sito LinkedIn del Dicembre 2017²², le due occupazioni che hanno visto il più alto incremento dal 2012 sono *Machine Learning Engineer* e *Data Scientist*, con crescite di 9.8 e 6.5 volte rispettivamente. Tra le prime tre posizioni nei lavori svolti precedentemente da chi ha assunto le posizioni di *Machine Learning Engineer* e *Data Scientist* troviamo *Research Assistant* e *Teaching Assistant*, a testimonianza di un flusso che proviene dall'accademia.

Ancora, il rapporto 2018 dell'azienda IBM *The quant crunch: how the demand of data science skills is disrupting the job market*²³ parla di un incremento del 40% dei posti di *Data Scientist* nel 2016 e di una proiezione di crescita uguale di qui al 2020. Tra questi, circa il 40% richiede un titolo di studio dal master in su, mentre il 78% almeno tre anni di esperienza.

Questi dati suggeriscono un forte bisogno di manodopera altamente formata in discipline scientifiche, che sembrano in molti casi provenire dall'accademia. Guardiamo dunque alle condizioni di lavoro proprio nell'accademia.

L'ipotesi che prendiamo in considerazione è che molto dell'interesse dei giovani ricercatori verso il ML risiede nella possibilità di acquisire tecniche spendibili in un mercato del lavoro che, abbiamo visto, è in enorme espansione, al contrario del panorama nel campo dell'università: poche prospettive di posto fisso, alta competizione, precarietà e stress poco sostenibili a fronte di bassi salari.

Negli USA e in Europa, almeno fino al 2008, le università non sono in grado di assorbire i dottori e questo trend è in costante peggioramento

²¹ 2018 Outlook: Machine Learning and Artificial Intelligence, A Survey of 1,600+ Data Professionals (14 pp., PDF, no opt-in).

²²<https://economicgraph.linkedin.com/research/LinkedIns-2017-US-Emerging-Jobs-Report>

²³<https://public.dhe.ibm.com/common/ssi/ecm/im/en/iml14576usen/analytics-analytics-platform-im-analyst-paper-or-report-iml14576usen-20171229.pdf>

in tutti i paesi avanzati²⁴ (Cyranoski, 2011). K. Powell (Powell, 2015) fa notare come negli USA il numero di postdoc nelle scienze è cresciuto del 150% tra il 2000 e il 2012, senza però nessun parallelo incremento delle posizioni fisse e delle *tenure-track*. Come risultato, sarà il settore privato (o pubblico, ma governativo, cioè non accademico) a impiegare i dottori.

Sempre negli USA, l'ultimo report del *National Science Foundation*, il *Science and Engineering indicator 2018*²⁵ mostra come la percentuale di dottori di ricerca con speranza di ottenere una posizione di tenure o *tenure-track* è solo del 20% circa. Dal 2000 ad oggi si riscontra una netta diminuzione di posizioni da *full-time faculty* e contratti a termine. Restringendosi alle scienze fisiche, quelle che più ci riguardano, se nel 1973 il 21% del totale tra le posizioni in accademia erano occupate da fisici, oggi la percentuale è scesa al 14%.

Se la scarsità di posizioni stabili obbliga probabilmente molti ricercatori ad uscire dall'accademia contro la propria vocazione, i più alti stipendi in ambito privato hanno comunque una forte attrattiva: in media un postdoc americano guadagna 43.000 dollari nell'università e ben 73.000 nel privato come stipendio iniziale (Powell, 2015).

Da più parti (Gould, 2015; Powell, 2015; Cyranoski, 2011) si lancia l'allarme, auspicando come soluzione la drastica riduzione del numero di dottorati, e di una riconfigurazione strategica della struttura dei gruppi di ricerca, con un aumento delle posizioni fisse di ricercatore, a vantaggio della produttività "pro-capite". Per quanto una posizione simile possa apparire pragmatica in un sistema con scarsi finanziamenti, essa risulta insufficiente in quanto non tiene conto del contesto produttivo in cui l'università è inserita.

La ragione sistemica dell'alto numero di postdoc e dottori costretti a lasciare l'accademia è infatti il bisogno di formare, a spese dello Stato, un esercito di tecnici altamente specializzati a disposizione delle grandi e piccole aziende, ovvero risorse che verranno nella maggior parte dei casi trasferite al contesto privato.

Questi dati vanno dunque letti nel contesto del mondo del lavoro e della produzione moderni riassunti sopra: a fronte di un mercato in espansione enorme nel settore dell'AI vi è un bisogno estremo di manodopera nel campo della Data Science, che viene colmato grazie ad un'incapacità sistemica dell'accademia nell'assorbire i ricercatori.

²⁴ Discorsi a parte vanno fatti per paesi come Cina, India, Polonia, Singapore nei quali, per svariati motivi, solo recentemente si è assistito ad un boom di investimenti pubblici in ricerca).

²⁵ <https://www.nsf.gov/statistics/2018/nsb20181/report>

A riprova di questo fatto, nel già citato rapporto di IBM si stima che entro il 2020 la domanda di tecnici dei dati nel settore privato aumenterà del 28%, con un picco del 39% in settori altamente qualificati dove è necessario un dottorato (e dove gli stipendi superano i 100.000\$). A fronte di questo boom di richieste IBM denuncia una carenza di profili specifici e chiede esplicitamente che “il sistema educativo” si faccia carico della loro produzione²⁶.

6. IL MACHINE LEARNING DALLA FISICA ALLA SOCIETÀ

In questa sezione commenteremo alcuni elementi critici legati all'uso del ML in ambito sociale, paragonandole al caso della fisica. Evidenzieremo in particolare che molti di questi elementi sono parte del processo *già nelle scienze pure*, per ragioni intrinseche.

6.1. Precisione, falsi positivi e fairness

Quando si fa uso di modelli e strumenti matematici, un concetto chiave è la loro capacità di fornire previsioni “precise” in modo quantificabile.

Tuttavia, già all'interno della scienza, senza riferirsi al ML, il concetto o la soglia di “precisione” accettati nelle rispettive comunità variano significativamente. Un caso paradigmatico risale al 2015: la Società Italiana di Fisica (SIF) non sottoscrisse la *Dichiarazione sui cambiamenti climatici* per la Conferenza di Parigi sul clima COP 21, obiettando che il grado di precisione degli studi sul clima non era abbastanza elevato per affermare “incontrovertibilmente” l'impronta antropica sul clima. La definizione di “incontrovertibile”, o la soglia oltre la quale si accetta una misura come evidenza, vengono tuttavia decise dalla comunità e non sono “contenute” nei dati o nella misura; sono il risultato di tecniche statistiche che possono variare significativamente a seconda del campo.

Anche la definizione stessa di “precisione” di un algoritmo può variare in modo sostanziale. Vediamo il caso del ML con un esempio dalla fisica: un rivelatore di particelle. L'algoritmo che regola il *trigger* dell'esperimento può essere programmato per minimizzare, ad esempio, la perdita di “veri segnali” a scapito di avere un maggior numero di

²⁶ Le conclusioni di questa sezione poggiano su dati provenienti dal mercato del lavoro combinati con l'analisi delle pubblicazioni scientifiche delle sezioni precedenti. D'altra parte, sarebbe utile e necessario corredare queste conclusioni con un lavoro di inchiesta nella comunità accademica, e in particolare tra i giovani ricercatori.

“falsi allarmi” (cioè eventi erroneamente classificati come segnale), oppure per minimizzare la *rate* di falsi allarmi, a costo di rischiare di mancare alcuni eventi reali. È importante sottolineare che è la stessa definizione matematica di questi due casi (“Falso positivo” e “Falso negativo”) a far sì che un compromesso tra i due sia necessario²⁷. Esistono addirittura altre possibili scelte oltre alle due qui citate. Questo tipo di scelte è essenziale nella scienza (e nell’uso degli algoritmi in generale), fa parte del suo stesso statuto e non è superabile con nessun avanzamento tecnologico. Se nel caso di un esperimento di fisica delle particelle una simile decisione potrebbe avere conseguenze tutto sommato rimediabili (magari un tempo più lungo per arrivare ad una scoperta), non lo stesso si può dire di un software che abbia un ruolo decisionale o consultivo nella vita pubblica. Ciò ci porta a riformulare tali concetti in termini di *fairness* (equità) dell’algoritmo. Un esempio famoso è il software COMPAS²⁸, usato negli Stati Uniti per predire il rischio nel rilasciare un detenuto su cauzione prima del processo (Courtland, 2018). Un team di giornalisti ha mostrato come il software non fosse “equo” – al contrario di quanto sostenuto dall’azienda produttrice –, in quanto prediceva un maggior numero di falsi positivi tra la comunità afroamericana. Alla radice della controversia vi è proprio il fatto che i due team usavano definizioni di equità diverse. È importante capire che questa differenza si dà *in primis* a livello matematico: esistono infatti diverse definizioni di *fairness*²⁹ che sono statisticamente impossibili da conciliare, e la cui scelta è in mano al programmatore. Il prof. A. Narayanan, *computer scientist* a Princeton, ha recentemente tenuto un seminario dal titolo *21 fairness definitions and their politic*³⁰ dove mostra appunto che tante possono essere le definizioni di questo concetto applicabili in un algoritmo.

Si potrebbe obiettare che molti classificatori hanno differenze che possono essere minimali (pochi punti percentuali, o meno) tra i diversi tipi di errore. Tuttavia, nel caso di vite umane ogni dettaglio può essere determinante, come nel caso di SKYNET, un algoritmo usato dalla

²⁷ Ad esempio, l’*f-score*, una media ponderata tra le due, o l’*accuracy*.

²⁸ COMPAS (Correctional Offender Management Profiling for Alternative Sanctions), prodotto dall’azienda Northpointe, ora Equivant (<http://www.equivant.com/>)

²⁹ In questo caso la, *predictive parity* – cioè il fatto che la probabilità di essere ri-arrestati tra bianchi e neri non dipenda dal gruppo di appartenenza – e il *false positive rate*, cioè la misidentificazione come soggetto a rischio. A sua volta un minimo *false positive rate* è inconciliabile con un minimo *false negative rate* - in questo caso, la misidentificazione come soggetto non a rischio.

³⁰ Qui il video della conferenza: <https://www.youtube.com/watch?v=jIXluYdnyyk>. Per il report: <https://docs.google.com/document/d/1bnQKzFAzCTcBcNvW5tsPuSDje8WWWY-SF4wQm6TLvQ/edit>

NSA per monitorare la popolazione pakistana assegnando un rischio di collusione col terrorismo. Una fuga di notizie³¹ ha rivelato che il livello di falsi positivi poteva oscillare tra lo 0.008% e il 0.18%. Su una popolazione di 55 milioni di abitanti, questo significa che fino a circa 100mila pakistani potrebbero essere erroneamente messi sotto sorveglianza come sospetti terroristi (Hosni, 2017). Al di là dei numeri, questo esempio mostra anche la necessità di una decisione a monte sull'utilizzabilità dell'algoritmo stesso: è accettabile applicare uno strumento simile al monitoraggio di un gruppo etnico (per quanto piccolo, addirittura nullo, possa essere l'errore)?

6.2. *Bias e ipotesi autoavveranti*

Nella prima sezione abbiamo parlato del concetto di “dati carichi di teoria” introdotto da Kuhn. Questo concetto è legato a quello di *bias*. Nell'applicazione degli algoritmi di DL al contesto sociale abbiamo già molti esempi che dovrebbero fare scuola. Vediamone uno eclatante: lo scorso anno in uno studio dell'università di Stanford (Wang, 2017) una DNN è stata usata per distinguere, tramite riconoscimento facciale, tra eterosessuali e omosessuali. A detta degli autori la precisione oscillava tra l'80% e il 70% – molto maggiore dei risultati di un umano. Molte testate giornalistiche rilanciarono lo studio, che fece scalpore perché giustificava possibili ipotesi neo-lombrosiane e cause genetiche dell'omosessualità. Qualche mese dopo i ricercatori di Google³² mostrarono con analisi più sofisticate che l'algoritmo era affetto da grossolani stereotipi di genere nell'allenamento e più che distinguere la fisionomia, discriminava in base al trucco, gli occhiali e altre caratteristiche riguardanti il look e contesto sociale più che la genetica. Resta un altro grande problema: la domanda posta all'algoritmo e la divisione binaria tra etero e omosessuali sono di per sé controverse e frutto di stereotipi. In questo caso, la non-neutralità della scienza e dell'algoritmo sono palesi. Anche qui possiamo trovare analogie con la fisica. Nel 1572, Tycho Brahe osservò un nuovo oggetto celeste molto luminoso, che interpretò come la formazione di una nuova stella. Oggi sappiamo che l'evento è una Supernova Ia, un evento astronomico che corrisponde alla fine del ciclo di vita di una stella, non conosciuto all'epoca di Tycho Brahe. Lo stesso errore di classificazione potrebbe

³¹ SKYNET: Applying Advanced Cloud-based Behavior Analytics. The Intercept, 8 May 2005.

³²<https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>

avvenire nel caso di un moderno algoritmo di ML allenato, poniamo il caso, a classificare l'immagine di un oggetto celeste come "stella" o "galassia". Il problema è di nuovo che il numero di classi è determinato dall'esterno e guidato dalla teoria.

Un ulteriore spunto di riflessione viene ancora dal software COMPAS. Il fatto che la predizione di rischio tra gli afroamericani sia più elevata e che questo risulti da un bias pregresso, innesca un processo per cui tali pregiudizi vengono confermati e probabilmente enfatizzati.

Un altro esempio sono i casi dei cosiddetti software di *predictive policing* usati negli Stati Uniti per predire in quali aree è più probabile che avvenga un crimine (senza sapere, sottolineiamo, per quale ragione), come *PredPol*³³. L'uso indiscriminato dei dati a disposizione può creare un feedback per cui, se gli agenti sono spinti ad essere presenti in quartieri precedentemente segnati da più reati, compiranno più interventi per il semplice fatto di essere già sul posto, e non perché il rischio sia realmente più elevato (Ensign, 2017). Questo indurrebbe un circolo vizioso paradossale per cui la pretesa "pericolosità" verrebbe confermata! Alla base vi è una doppia fallacia: l'assumere che i dati possano rivelare tendenze valide in futuro in assenza di una spiegazione (quello che in fisica avevamo chiamato "modello"), e il fatto che più interventi polizieschi siano dovuti a maggiore insicurezza (classico caso di correlazione spuria). Non è del resto questa una situazione nuova. In economia la polemica corre da anni su binari simili: l'economia neoclassica non solo fallisce nelle sue previsioni (Sylos Labini, 2016), ma è l'uso stesso dei suoi strumenti, ricavati all'interno di un paradigma teorico senza basi empiriche, che porta gli attori in gioco a comportarsi in maniera tale da far scoppiare le crisi (Bouchaud, 2008). Ancora una volta, è determinante capire come in fisica il fatto che un modello sia presente per fornire previsioni non è abbastanza enfatizzato nella discussione sul ruolo degli algoritmi.

6.3. *La proprietà dei dati.*

In fisica, la replicabilità degli esperimenti e l'accesso ai dati sono considerati un principio inviolabile e condizioni necessarie per trarre qualsiasi conclusione.

Inoltre, come abbiamo visto, il tipo di dati da raccogliere (in altri termini, la progettazione dell'esperimento) sono guidati dalla teoria e parte del processo di ricerca. La stessa cosa non si può dire delle

³³ PredPol, Predict Prevent Crime. Predictive Policing Software: www.predpol.com/

moderne *data warehouse* e dell'uso di modelli data-driven al di fuori dell'accademia. Particolarmente grave è la questione della proprietà dei dati. Se infatti molti strumenti alla base dell'estrazione di valore tramite previsioni sono oggi pubblici³⁴, la vera risorsa alla base dell'accumulazione sono proprio i dati su cui allenare gli algoritmi e questi sono spesso proprietari: "data is the new oil", recita una battuta, e tanto basta per farne capire il valore. È importante menzionare che questo fenomeno ha un impatto sul processo di migrazione al di fuori dell'accademia, in quanto spesso ricercatori in istituzioni pubbliche hanno serie difficoltà di accesso ai dati per condurre le loro ricerche, mentre il problema non si pone per dipendenti di grandi piattaforme private³⁵. La privatizzazione dei dati genera quindi delle gerarchie di potere rigidissime e violente, aggravate dal fatto che la proprietà tende a concentrarsi nelle mani di poche grandi piattaforme.

6.4. *Il ruolo dello scienziato nella società.*

Abbiamo sottolineato più volte come gli algoritmi di ML prevedano, nel loro utilizzo, una serie di passaggi in cui l'utilizzatore è tenuto a prendere delle decisioni che influenzeranno l'output, come la selezione delle feature o la scelta di una soglia di probabilità³⁶ e che questa arbitrarietà è presente a vari livelli. I moderni algoritmi di DL, invece, funzionano più come *black boxes* difficilmente decifrabili, in cui il processo che porta alla decisione finale può addirittura non essere ricostruibile, ma che comporta comunque scelte come il criterio con cui valutarne la precisione. In entrambi i casi, inoltre, si assume la possibilità di applicare a casi non conosciuti uno strumento predittivo allenato solo su altri esempi, rinunciando completamente ad un modello. L'adoptare questo modo di procedere "induttivista" comporta già una profonda presa di posizione epistemologica, spesso ignorata al di fuori della fisica.

L'uso della modellizzazione e della matematica in contesto sociale ha una doppia implicazione: da un lato rende falsamente oggettivo (o, si potrebbe dire, naturale, neutrale) un processo "social-driven" (parafrasando il "data-driven"), in un certo senso "normalizzandolo".

³⁴ Come nel caso delle librerie python TensorFlow e PyTorch, sviluppate rispettivamente da Google e Facebook per implementare reti neurali profonde.

³⁵ L'esempio più eclatante è il caso di Facebook, che non permette l'accesso ai suoi contenuti tramite API - nemmeno in forma limitata, come nel caso di Twitter. Per una discussione di alcuni casi specifici legati a Google e Facebook, si veda Cozzo (2018).

³⁶ Anche se in questo lavoro siamo interessati al caso del ML, il discorso vale in generale per ogni tipo di algoritmo.

Dall'altro fa ricadere sullo scienziato, tramite la sua ricerca, la decisione politica, ergendolo a tecnico-veicolo del sistema culturale ed economico egemone. Se per la prima implicazione abbiamo già fatto obiezioni sui rischi di una applicazione naive in campo sociale di tecniche matematiche, forzando le ipotesi in cui erano state inizialmente sviluppate (qual è il dominio, come retroagisce il sistema, etc.), ben più sottile, ma altrettanto importante e rischiosa è la seconda implicazione. Lo scienziato, in un contesto di divisione tecnica del lavoro, è spesso convinto o è spinto a credere di avere un ruolo puramente tecnico e di non avere un ruolo politico oppure, paradossalmente, crede di averne uno ancora più alto: essere il disvelatore della verità, davanti alla quale la politica non conta. Feticismo della scienza, alienazione e precarietà lavorativa si mischiano per rafforzare in maniera inconsapevole il sistema socioeconomico e ideologico, a maggior ragione quando nemmeno lo scienziato è pienamente consapevole del lavoro svolto dall'algoritmo, data la sua opacità e *bias*.

7. CONCLUSIONI

Questo lavoro è dedicato allo studio dell'utilizzo del ML nella fisica contemporanea. La scelta di questo caso parte dal riconoscimento della natura storicamente e politicamente determinata della scienza e il suo legame con la materialità dell'epoca in cui essa si sviluppa. Il capitalismo contemporaneo si fonda sull'estrazione di valore da grandi moli di dati tramite algoritmi avanzati di ML e grande potenza di calcolo. Questi tre elementi hanno segnato un momento di svolta nel ruolo delle grandi piattaforme, mutando quindi la produzione e il mercato del lavoro. È naturale dunque chiedersi quale sia l'impatto di tale processo sullo sviluppo della scienza, e allo stesso tempo quale ruolo la scienza abbia nell'alimentare e legittimare un certo paradigma nei suoi aspetti politici e sociali, in questo caso il ruolo determinante che gli algoritmi giocano nel mondo contemporaneo.

L'analisi si è sviluppata attorno a tre elementi:

1. *la scienza non è neutrale*: la ricerca scientifica, anche nelle scienze "pure", evolve in modo non lineare seguendo inestricabilmente le trasformazioni del modello produttivo egemone nel resto della società.

2. *gli algoritmi non sono neutrali*, sia in quanto prodotti in un ambito scientifico che non è neutrale, sia perché contengono elementi soggetti a decisioni esterne non eliminabili con alcun progresso tecnologico.

3. Nella prospettiva di un paradigma “data-driven”, diventa tanto più importante affermare l’impossibilità di avere strumenti oggettivi e neutrali addirittura all’interno delle scienze pure.

Abbiamo affrontato il problema secondo linee complementari.

Una analisi storica e bibliometrica dell’utilizzo del ML nella fisica (sezione IIIa) ci ha permesso innanzitutto di mostrare come la sua esplosione coincida con avanzamenti tecnologici avvenuti fuori dall’accademia per ragioni legate a necessità del mercato (estrarre valore da grandi moli di dati) che hanno a loro volta determinato una profonda trasformazione della produzione e del mercato del lavoro. Abbiamo inoltre mostrato come il boom del ML in fisica corrisponda anche a necessità materiali contingenti dei ricercatori, i quali, sottoposti a condizioni di lavoro precarie e a poche prospettive di carriera dentro l’accademia, vengono probabilmente spinti a indirizzare studi e ricerche verso un ambito spendibile nel mercato del lavoro (sezione IV).

In secondo luogo, abbiamo passato in rassegna l’uso del ML nei campi della fisica fondamentale in cui questo è più rilevante: la fisica delle particelle, l’astrofisica, la meccanica quantistica e la fisica della materia condensata (sezioni IIIb-IIIId). Significativamente, l’impatto dei fattori esterni all’accademia sopra menzionato si riscontra direttamente nella letteratura scientifica. Allo stesso tempo, abbiamo mostrato come ipotesi à la *Anderson* che paventino una prossima “fine della teoria” e una scienza che faccia a meno di modelli non sembrano reggere un confronto approfondito con la realtà, al di là di ogni entusiasmo tecnofuturistico³⁷.

Infine, il fatto che problemi di interpretazione dei risultati emergano chiaramente anche in un contesto teorico e sperimentale fortemente controllabile come quello della fisica fondamentale, in cui le variabili in gioco si presumono ben definite, ci ha permesso di trasferire le nostre critiche all’utilizzo degli algoritmi in ambito sociale, dove *bias*, pregiudizi, incontrollabilità del sistema in oggetto e *feedback* assumono un ruolo dirimente e in cui in gioco non ci sono variabili di un esperimento, ma persone (sezione V).

Tuttavia, tornando ad *Anderson* e valutando l’impatto che gli algoritmi hanno oggi, il loro impiego ed una sempre più diffusa retorica basata su presunte “verità scientifiche” e in particolare sui dati, è forte il sospetto che la sua previsione, anche se basata su presupposti falsi, si possa avverare. La possibilità di predizioni sempre più accurate, la

³⁷ Per altri commenti critici interni alla comunità scientifica si veda anche (Zdeborová, 2017).

potenza di calcolo, e i progressi stupefacenti nel DL sembrano poter riaprire acriticamente la strada al ritorno di una sorta di nuovo positivismo, scientifico e sociale, fondato sui dati, e trainato non tanto dagli scienziati, che restano consapevoli della posta in gioco nel metodo scientifico, ma da fattori esterni alla scienza.

Un vero cambio di paradigma, questo sì, nel senso di Kuhn. Non si tratta più di un positivismo basato sull'oggettività del metodo, ma sulla presunta neutralità dell'algoritmo, sulla possibilità di spiegare tutto attraverso i dati e sulla non-necessità di costruire modelli predittivi (nella scienza) e in ultima istanza immaginarne di migliori nella società.

Se restiamo fedeli all'esercizio di prendere Anderson sul serio, lo scenario futuro è uno sviluppo dominato dai dati, una distopia destinata a riprodurre in maniera circolare lo stato di cose esistenti.

Tocca allora anche al mondo scientifico riportare il problema alla giusta dimensione, ricontestualizzando la scienza stessa e contribuendo a decostruirne l'aura di assoluta oggettività, neutralità ed astoricità che oggi ne caratterizza dati ed algoritmi. Con questo lavoro abbiamo cercato di fare un passo in questa direzione.

RIFERIMENTI BIBLIOGRAFICI

- ABBOTT, B. P. *et al.* (2016). LIGO Scientific Collaboration and Virgo Collaboration. *Physical Review Letters*. 116, 061102.
- ADAM-BOURDARIOS, C. *et al.* (2014). The Higgs boson machine learning challenge., *Journal of Machine Learning Research*, 42, 19-55.
- ANDERSON, C. (2008). The End of Theory, Will the Data Deluge Makes the Scientific Method Obsolete?. *Wired Magazine*, June 23.
- AYRES, D.S. *et al.* (2007), *The NOvA technical design report*. NOvA Collaboration.
- BALDI, P., SADOWSKI, P., WHITESON, D. (2014). Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5, 4308.
- BIAMONTE, J. *et al.* (2017), Quantum Machine Learning, *Nature* 549, 195-202.
- BOUCHAUD, J.-P. (2008). Economics needs a scientific revolution. *Nature*, 455, 1181-1181.
- Carrasquilla, J., Melko, R. G. (2017). Machine learning phases of matter. *Nature Physics*, 13, 431-434.
- CERN (2017). ATLAS Collaboration. *Identification of Jets Containing*
-

- b-Hadrons with Recurrent Neural Networks at the ATLAS Experiment*. Report No. ATL-PHYS- PUB-2017-003.
- CICCOTTI, G., CINI, M., DE MARIA, M., JONA-LASINIO, G. (1976). *L'Ape e l'architetto: paradigmi scientifici e materialismo storico*. Roma: Feltrinelli.
- COURTLAND, R. (2018). Bias detectives: the researchers striving to make algorithms fair. *Nature* 558, 357-360.
- COZZO, E. (2018). *Pratiche scientifiche ai tempi del capitalismo di piattaforma*. In D. Gambetta (a cura di), *Datacrazia. Politica, cultura algoritmica e conflitti al tempo dei big data*. Roma: Escathon.
- CRAWFORD, K., BOYD, D. (2011). *Six provocation for big data*. Oxford Internet Institute's "A Decade in Internet Time: Symposium on the Dynamics of the Internet and Society".
- CYRANOSKI, D., GILBERT, N., LEDFORD, H., NAYAR, A. YAHIA, M (2011). The Ph.D. factory. The world is producing more PhDs than ever before. Is it time to stop?. *Nature*, 472, 276-279.
- DUNJKO, V., BRIEGEL, H. J. (2018). Machine learning & artificial intelligence in the quantum domain: a review of recent progres. *Reports on Progress in Physics*, 81, 074001.
- ENSIGN, D. FRIEDLER, S. A., NEVILLE, S., SCHEIDEGGER, C., VENKATASUBRAMANIAN, S. (2017). Runaway feedback loops in predictive policing. *arXiv*, arXiv:1706.09847
- FEYERABEND, P. K. (2002). *Contro il metodo. Abbozzo di una teoria anarchica della conoscenza*. Milano: Feltrinelli.
- GABBARD, H., WILLIAMS, M., HAYES, F., MESSENGER, C., (2018). *Physical Review Letters*, 120, 141103
- GILLET, N. et al., (2018). Deep learning from 21-cm images of the Cosmic Dawn. *arXiv*, arXiv:1805.02699
- GOLDTHORPE, J. H., (2016). *Sociology as a population science*. Cambridge University Press.
- GOMEZ-URIBE, C. A., HUNT, N. (2016). The Netflix recommender system: Algorithms, Business value, and Innovation. *ACM Transactions on Management Information Systems (TMIS)*, 6 (4), 13.
- GOODFELLOW, I., BENGIO, Y., COURVILLE, A. (2016). *Deep learning*. Cambridge: MIT press.
- Gould, J. (2015). How to build a better PhD, *Nature* 528, 22-25
- HEZAVEH, Y. D., LEVASSEUR, L. P., MARSHALL, P. J. (2017). Fast automated analysis of strong gravitational lenses with convolutional neural networks, *Nature*, 548, 555-557
-

- HINTON, G. E., SRIVASTAVA, N., KRIZHEVSKY, A., SUTSKEVER, I., SALAKHUTDINOV, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. *arXiv*, arXiv:1207.0580.
- HOSNI, H., VULPIANI, A. (2017). Forecasting in the light of Big Data. *arXiv*, arXiv:1705.11186
- HUEMBELI, P., DAUPHIN, A., WITTEK, P. (2018). Identifying quantum phase transitions with adversarial neural networks. *Physical Review B*, 97(13), 134109.
- KOLANOVIC, M., KRISHNAMACHARI, R. T. (2017). *Big Data and AI Strategies Machine Learning and Alternative Data Approach to Investing*. JP Morgan Global Quantitative & Derivatives Strategy Report.
- KRIZHEVSKY, A., SUTSKEVER, I., HINTON, G. (2012). ImageNet classification with deep convolutional neural networks. *NIPS'12 Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, 1097-1105.
- KUHN, T. S. (1972). *La struttura delle rivoluzioni scientifiche*. Torino: Einaudi.
- LAKATOS, I., FEYERABEND, P., MOTTERLINI, M. (1995). *Sull'orlo della scienza. Pro e contro il metodo*. Milano: Raffaello Cortina.
- LAM, C., KIPPING, D., (2018). Transit clairvoyance: enhancing *TESS* follow-up using artificial neural networks, *Monthly Notices of the Royal Astronomical Society*, 476(4), 5692-5697.
- LATOUR, B. (2010). *Tarde's idea of quantification*. In M. Candea (Eds.) *The Social After Gabriel Tarde: Debates and Assessments* (pp. 145-162). New York: Routledge.
- LLOYD, S., MOHSENI, M, REBENTROST, P. (2013). Quantum algorithms for supervised and unsupervised machine learning. *arXiv*, arXiv:1307.0411
- MAAS, A. L., HANNUN, A. Y., NG, A. Y. (2013). Rectifier nonlinearities improve neural network acoustic models. Proceedings of the 30 th International Conference on Machine Learning, 1-6.
- MOSTAFA, A., Y. S., MALIK M.-I., LIN, H.-T. (2012), *Learning from data*. Vol. 4. New York: AMLBook.
- POINCARÉ, H., (1905). *Le valeur de la science*. Paris: Flammarion
- POPPER, K. R. (1991). *Scienza e filosofia*. Torino: Einaudi.
- POWELL, K. (2015). The future of the postdoc. *Nature*, 520, 144-147.
- RADOVIC, A. *et al.* (2018). Machine learning at the energy and intensity frontiers of particle physics, *Nature* 560, 41-48.
- RAINA, R., MADHAVAN, A., NG, A. Y. (2009, June). Large-scale deep
-

- unsupervised learning using graphics processors. *Proceedings of the 26th annual international conference on machine learning*, 873-880).
- ROSENBLATT, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6), 386.
- SAMUEL, A. (1959). Some Studies in Machine Learning Using the Game of Checkers. *IBM Journal of Research and Development*, 3(3), 210-229.
- SHALLUE, C. J., VANDERBURG, A. (2018). Identifying Exoplanets with Deep Learning: A Five-planet Resonant Chain around Kepler-80 and an Eighth Planet around Kepler-90. *The Astronomical Journal*, 155(2), 1-22.
- SYLOS LABINI, F. (2016). *Rischio e previsione. Cosa può dirci la scienza sulla crisi*. Roma-Bari: Laterza.
- SZEGEDY, C. *et al.*, (2014). Going deeper with convolutions. *arXiv*, arXiv:1409.4842.
- VAN NIEUWENBURG, E. P. L., Liu, Y. H., Huber, S. D. (2017). Learning phase transitions by confusion. *Nature Physics*, 13(5), 435.
- WANG, Y., KOSINSKI, M. (2017) Deep neural networks are more accurate than humans at detecting sexual orientation from facial images. *Journal of personality and social psychology*, 114(2), 246.
- ZDEBOROVÁ, L. (2017). Machine learning: new tool in the box. *Nature Physics*, 13(5), 420.
-

Numero chiuso il 30 marzo 2019



ULTIMI NUMERI

2018/2 (aprile-giugno):

1. ILARIA IANNUZZI, L'ebraismo nella formazione dello spirito capitalistico. Un excursus tra le opere di Werner Sombart;
2. NICOLÒ PENNUCCI, Gramsci e Bourdieu sul problema dello Stato. Dalla teoria della dominazione alla sociologia storica;
3. ROSSELLA REGA, ROBERTA BRACCIALE, La self-personalization dei leader politici su Twitter. Tra professionalizzazione e intimizzazione;
4. STEFANO SACCHETTI, Il mondo allo specchio. La seconda modernità nel cinema di Gabriele Salvatores;
5. GIULIA PRATELLI, La musica come strumento per osservare il mutamento sociale. Dylan, Mozart, Mahler e Toscanini;
6. LUCA CORCHIA, Sugli inizi dell'interpretazione sociologica del rock. Alla ricerca di un nuovo canone estetico;
7. LETIZIA MATERASSI, Social media e comunicazione della salute, di Alessandro Lovari.

2018/3 (luglio-settembre):

1. RICARDO A. DELLO BUONO, Social Constructionism in Decline. A "Natural History" of a Paradigmatic Crisis;
2. MAURO LENCI, L'Occidente, l'altro e le società multiculturali;
3. ANDREA BORGHINI, Il progetto dei Poli universitari penitenziari tra filantropia e istituzionalizzazione;
4. EMILIANA MANGONE, Cultural Traumas. The Earthquake in Italy: A Case Study;
5. MARIA MATTURRO, MASSIMO SANTORO, Madre di cuore e non di pancia. Uno studio empirico sulle risonanze emotive della donna che si accinge al percorso adottivo;
6. PAULINA SABUGAL, Amore e identità. Il caso dell'immigrazione messicana in Italia;
7. FRANCESCO GIACOMANTONIO, Destino moderno. Jürgen Habermas. Il pensiero e la critica, di Antonio De Simone.
8. VINCENZO MELE, Critica della folla, di Sabina Curti.

2018/4 (ottobre-dicembre):

1. ENRICO CAMPO, ANTONIO MARTELLA, LUCA CICCARESE, Gli algoritmi come costruzione sociale. Neutralità, potere e opacità;
 2. MASSIMO AIROLDI, DANIELE GAMBETTA, Sul mito della neutralità algoritmica;
 3. CHIARA VISENTIN, Il potere razionale degli algoritmi tra burocrazia e nuovi idealtipi;
 4. MATTIA GALEOTTI, Discriminazione e algoritmi;
 5. BIAGIO ARAGONA, CRISTIANO FELACO, La costruzione socio-tecnica degli algoritmi;
 6. ANIELLO LAMPO, MICHELE MANCARELLA, ANGELO PIGA, La (non) neutralità della scienza e degli algoritmi;
 8. LUCA SERAFINI, Oltre le bolle dei filtri e le tribù online;
 9. COSTANTINO CARUGNO, TOMMASO RADICIONI, Echo chambers e polarizzazione;
 10. IRENE PSAROUDAKIS, Mario Tirino, Antonio Tramontana (2018), I riflessi di «Black Mirror»;
 11. JUNIO AGLIOTTI COLOMBINI, Daniele Gambetta (2018), Datacracia;
 12. PAOLA IMPERATORE, Safiya Umoja Noble (2018), Algorithms of Oppression;
 13. DAVIDE BERALDO, Cathy O'Neil (2016), Weapons of Math Destruction;
 14. LETIZIA CHIAPPINI, John Cheney-Lippold (2017), We Are Data.
-