**Claudio Sarra[1]**

# Data Mining and Knowledge Discovery. Preliminaries for a Critical Examination of the Data Driven Society

[1] Università degli Studi di Padova, Padua, Italy, E-mail: claudio.sarra@unipd.it. https://orcid.org/0000-0002-5850-9972.

**Abstract:**
Data Mining (DM) is the analytical activity aimed at revealing new "knowledge" from data useful for further decision-making processes. These techniques have recently acquired enormous importance as they seem to fit perfectly the requests of the so called "Data Driven World". In this paper, first I give an overview of DM, and of the most relevant criticisms raised so far. Then using a well-known case study and the European General Data Protection Regulation as benchmark, I show that there are some specific ambiguities in this use of "knowledge" which are relevant for the ethical and legal assessment of DM.
**Keywords:** Data Mining, knowledge discovery, data protection, law and ICT, GDPR

## 1    Introduction. "From a world full of data to a data-driven world"

According to a common definition, Data Mining (DM) is the crucial step in the process of *knowledge discovery*, brought about through analysis of large quantities of data, aiming at the extraction of non-trivial, potentially useful, implicit and previously ignored information, that can be arranged in patterns to be applied in further decision making processes.[1] It includes steps like problem definition, data preparation, selection and optimization of appropriate analysis techniques, modeling and deployments of results.[2]

Although by no means new, DM techniques are receiving greater attention and relevance nowadays as they lend themselves to application in almost every field of human knowledge and research, as well as in the modernization of business management.[3]

The data revolution is actually an impressive phenomenon. The crucial shift from a cyberspace made on the basis of a neat distinction between data-producers and data-users, to an informational environment where everyone, in every moment, is at the same time user-and-producer of data, is what boosted the so called Big Data era.[4] In 2016 a report by McKinsey Global Group tried to measure this phenomenon estimating, for instance, that every day a Terabyte of data was produced by NASA satellites, that there was 650 millions of websites, with dynamical possibilities of data creations, that there was more or less a billion of *Facebook* users producing three billions of posts every day, that a small to medium company had data warehouses containing data for one hundred million commercial transactions, with everything continuously growing every minute, making those very measurements doomed to become severely obsolete as soon as published.[5]

Moreover, it is on the quite reasonable prediction that this growing trend is far to slow down, that another crucial passing point is announced: that "from a world full of data to a data driven world".[6] What this slogan actually means is, in first place, that the evolution of the economic behavior of private and public agents will be more and more based on the exploitation of this world-data by means of decision-processes backed up by DM, for more profit and economic advantage, something that poses some deep legal and ethical issues.[7] This is the reason why DM is acquiring such a great importance that, for instance, the Italian Public Authority for Data Protection in a recent interview said that it represents one of the major challenges for traditional legal categories in the third millennium.[8]

Now, this situation, obviously, raises a series of technical issues concerning the search for the best tools and algorithms for useful analysis in the ever-changing sea of data: so large an amount that cannot be stored in any physical support nor can it be processed as a whole (*Big Data*).[9] But, on the other hand, it also seems to offer the best chances for previously impossible discoveries: since almost every moment of people life is recorded, datafied and put in relation to many others, we witness the rise of a never seen before power for the human

---

Claudio Sarra is the corresponding author.

intellect to extract unsuspected knowledge that may be needed for the development and the future well-being of mankind.

In this work I try to focus on the main core of DM which is centered around the very idea of "knowledge". Since this is a philosophically "thick" word, it may be the case that its use, especially in non-technical circles, hides some presuppositions or some ambiguities relevant for the ethical and legal assessment of problematic cases. Thus, after a description of the lines of development of the data-driven society and a brief review of the literature about the issues of DM, I will discuss the main topic using a well-known case study and the European General Data Protection Regulation as a benchmark. The discussion will show that different ideas of knowledge implied in DM processes can actually lead to different ethical and legal assessments, raising also a major question about what kind of knowledge we want to build our society on.

Thus this paper is structured as follows: § 2 explains how exactly large scale use of DM techniques can be efficient in determining relevant social changes, and the so called "modeling loop" is introduced; § 3 presents an overview of the criticisms raised so far towards the use of DM by powerful social and economic agents and governments; § 4 focuses on the main core of DM, that is the idea of a "knowledge discovery", and shows how the ambiguity which lies behind it can lead to legal and ethical perplexities; § 5 tries to reconcile the ambiguity and suggests an "epistemological precautionary principle" in dealing with these issues; § 6 draws some conclusions highlighting four points for further deepening the critical evaluation of DM.

## 2 From models to reality and back: the "modeling loop"

After the massive spread of personal computers (the Eighties of the XX cent.), the rise and the diffusion of Internet access, and the perfection of the necessary infrastructures for the broadband connections (Nineties-2000), it was the *smart mobile* phone (2007) that realized a major step forward towards the simultaneous data use and production by every single user. From that time on, a new level of concern for privacy and personal data protection related issues has been reached: as a matter of fact those data-producing tools go along with the user in every moment of his day, recording and transmitting his position, moves, opinions, healthy condition, lifestyle, desires, physical traits, behaviors, sentiments, and much more. Such a real-time life-scanning marks every moment of his ordinary day, offering data for the construction of a virtual personality which may be supposed to be progressively adherent to his real one. Then his virtual projection gives material for further profiling, which aims at the building of a model to be used as reference point in order to evaluate the future behavior of relevant groups of other real subjects.[10]

Those models – built from massive data taken from reality scans – are then used in real-life decision-making processes, in particular by powerful economic agents (or, also by governments),[11] for determining policies and large-scale regulations: for instance, when it comes to take decisions for marketing purposes, to set general rules for mass contracts by bank, insurance corps., or multi-national corporations, or to enhance recruitment practices, as well as to determine policies for crime prevention.[12]

But, on the other hand, this has also consequences on the way people behave, creating an interesting sort of *feedback.* Since people know that they are to accept some kind of regulations if they want to get access to certain services, and they have no power to negotiate them, as long as these services are enjoyed more and more by others, they gradually tend to conform to those standards provided by the author of the regulation. Regulations based on models, presuppose differentiations into groups, in other words they presuppose generalizations. The more these groups become empirically present in society the more this very fact can trigger processes of further social construction. We may use Georg Simmel's account on social construction to further shed some light on this point[13]: with group profiling structuring models, we are in the position to perceive others as "in some degree generalized", which, in Simmel's account, is the first *a priori* of society. But as we are expected to be part of that society, we feel also the possibility to act differently from what expected, that is to say, we experience our being "something more" than that, which is the second *a priori*. Finally, as we take our decision to behave in one way or another, we undergo a process of taking a social stand, contributing to the creation of *that* society made of all the functional differentiations that are "discernible only through the actual doing and experiencing of individual", which is the third *a priori*.[14] As long as those differentiations offer themselves in some datafied form, they also can be structured through similar model construction.

In other words, as long as models for generalized behaviors are introduced, a specific dialectic of social construction is triggered.

Thus, after the process of building *from reality to models,* a consequent adaptive process in the opposite direction *from models to reality* is actually engaged, obviously both on large scales. This is a sort of adaptive modification of real behaviors towards those predicted by the model, that in change produces a progressive verification of those same predictions, thus corroborating the model itself.[15] So, a continuous "modeling loop"

*from reality* (gathering of data) *to a model* (through DM techniques) and *from the model to reality* again (through the consequent adaptive behaviors) appears to be the macro-dynamics of the social construction in the *Big Data* era. And that is the way by which DM is *efficient* in shaping social reality, acquiring its new status as the most sensible point in the business intelligence set of tools.

DM techniques will arguably become more and more important and will be used for every economic decision in every field, making ubiquitous the aforementioned recursive circle *reality-model-reality* eventually leading to a complete *data driven world*.[16]

As the definition itself shows, the crucial aspect of DM is not the simple recovering of previously stored data, instead it consists in the organization of algorithmic procedures for manipulating huge quantities of data in order to highlight correlations as well as regularities which are not visible from disconnected data.

Hence the goal is to determine patterns, qualified groups and regulative schemes to include new occurrences and manage them according to the predictions made by means of the model.[17] The more the consequences of this way of managing will be in line with the predictions derived from the model, the more the model will be corroborated by reality. When things go this way, it is usually said that a new *knowledge* about how things actually are going has been acquired.

From the legal point of view (as well as from the ethical one), an interesting aspect is that procedures like those described so far are often used with the goal of producing *policies* to manage large quantities of future social and juridical relationships. As the development of such work requires very specialized competencies, often corporations hire other companies to do the analytical job, thus those who actually perform the DM happen to be different from those that will *use* the results they achieved for the actual business. In those cases, it is quite unlikely that users are completely aware of all the choices and decisions that it was necessary to take in order to prepare and conduct the analysis. This lack of "readability and legibility",[18] can be seen as an evolution of the "invisibles choices" issue that James Moor saw lucidly more that 30 years ago talking about the Computer Revolution.[19] Moor argued that one of the major features of the so called "Computer Revolution" is the "invisibility factor", that is the fact that whenever a certain task is assigned to a computer in order to have it done more efficiently, there is an exponential increase of hidden operations which are at play: in particular, he talked about invisible choices embedded in a software or, we can say today, in any application, portal or platform, made by the creators, sometimes without any particular critical reflections about it. Using models derived from DM to forge organizations, hides all the choices and decisions that have been taken in the construction of the knowledge presupposed.

Nowadays, decisions and practices based on sophisticated DM procedures are already quite common. Some fairly known examples are: the massive email sending with carefully selected targets for marketing purposes, the introduction – not previously agreed upon – of standards for selecting and/or evaluating workers, the conclusion or refusal of a financial agreement without any real negotiation but on the sole base of previous profiling techniques, and so on.[20] But science itself is more and more engaged in DM,[21] perhaps evolving to a stage where the very idea of a 'scientific method' will be declared obsolete.[22]

Legal scholars are now beginning to look more carefully into practices like these in order to evaluate if the technical mask is hiding a more serious challenge to the rights of people.[23]

## 3   General overview of DM criticisms

Now, in comparison with other fields of study such as Cyberethics or Sociology, a specific legal literature on DM is still relatively poor, especially if – in the name of specialty – we rule out the (huge) general literature dedicated to the legislation on personal data protection. Obviously enough, this literature matters, as those kind of concerns are of the utmost importance in the *data-driven world,* but the point is that this literature usually does not take into account the peculiarities of DM processes, that seems to put new challenges to the familiar tools of protection of personal data. As a matter of fact, those traditional tools are centered on the necessity of various types of legal bases in order for a data processing to be lawful,[24] as well as on the central role of Authorities to give regulation and face specific issues. As we are going to see, if the information acquired through DM are "new", "previously ignored" and "not trivial", as the definition implies, then there are some problematic issues when it comes to evaluate the lawfulness of their processing.

The European General Data Protection Regulation (2016/679, GDPR), establishes the duty of a "privacy by design and by default" approach in every development of applications that have to do with personal data, as well as a set of other demanding duties and safeguards. It does not mention DM explicitly, though it takes into consideration profiling and automatic decision making.[25] Now, contemporary profiling techniques include DM,[26] but since art. 13 and 14 GDPR provide for a duty to inform the data subject only about the existence "of an automatic processing referred to art. 22(1) including profiling", and to give "meaningful information about

the logic involved as well as the significance and the envisaged consequences of the proccessing" as a whole, we may doubt that the data subject could ask for a detailed explanation of the specific DM techniques used.[27]

As a consequence, in this framework, DM is not legally relevant *per se*, but only when included in profiling and when used to support automated decision making.

When the latter is the case, besides the information duties, the Data Controller faces also the prohibition provided for by art. 22, which states that decisions based "solely" on automated data processing, including profiling, which may have legal effects concerning the data subject or "similarly affects" him or her, are, in principle, forbidden. But, again, it is not DM *per se* to be prohibited or otherwise limited, but only the taking of relevant decisions based on complete automation which, of course, may – or may not – include DM.[28] Thus, the specific DM techniques used with all the decisions taken by the analysts to get the job done are *per se* legally irrelevant when used without profiling aims and outside a complete algorithmic decision-making procedure.

Apart from the legal literature, other scholars in different fields have already developed quite a structured criticism on DM massive use.

In order to give a general overview of the criticisms that has been raised so far, we can usefully regroup them in four types. A brief explanation for each will be given afterwards, as well as a suggestion for a philosophical unitary consideration, although, for the sake of present discussion, we are not going to deepen these speculations much further.

So, the most problematic issues raised so far to the contemporary use of DM techniques are related to:

1. privacy/surveillance concerns[29];

2. the introduction of new forms of discrimination as well as the strengthening of social inequalities[30];

3. forms of social exclusion, lack of transparency and participation[31];

4. methodological issues, in consideration of the "performative" efficacy of any social inquiry, in particular when brought about by powerful ICT companies to such a large extent such as those of the more advanced forms of ODM (Online Data Mining).[32]

The first issue is surely the most familiar to jurists, at least when it comes to deal with privacy legal protection, legitimate use of personal data and rights of the legitimate owner. Instead, while the protection towards direct forms of surveillance is related to constitutional rights, a bit less discussed is the transformation of surveillance practices because of the social changes induced by technology. As forms of continuous and reciprocal social monitoring become normal thanks to the widespread use of social media technology, the traditional paradigm of surveillance changes. Once characterized by: a) asymmetry of power between the surveillant and those kept under watch, b) strong hierarchy and c) mono-directional move from the first to the latter, it seems that nowadays, technology has determined the advent of a more complicated model.[33] Since social media allow the ordinary, continuous, highly focused, reciprocal monitoring among the members of a circle of "friends", and since (often in the name of "popularity") those circles are made of people who do not know each other very well, we may see here the rise of a different kind of surveillance practice. In this case, asymmetry, hierarchy and direction of control are not fixed but movable, as they change with reference to the particular moment and context. This kind of 'partecipatory surveillance', as it has been labeled,[34] comes with its proper reinforcement behaviors, like 'sharing', 'likes' and so on, which have they specific potential to trigger rewarding processes even at neurological level.[35] On the other hand, since social media are nowadays the most fruitful sources of usable data, economic companies take part strategically to this reciprocal monitoring play, influencing people behaviour (using techniques such as "astroturfing" or "likes" selling and purchase, and so forth), while learning and improving their economic decisions through DM techniques. Lastly, private companies that, for business reasons, holds large amounts of personal data are more and more engaged in collaboration with governments to support legitimate surveillance and law enforcement practices.[36]

About point 2), on one hand, scholars have already repeatedly drawn the attention to the fact that data are not neutral representations of reality. They are, instead, the product of a complex net of social practices, points of view, highly teleological conceptualizations, decisions and consequent elaborations.[37] In other words, data are the result of interactions in the entire "social ecology".[38] Thus, data often embody and reflect social structures made of peculiar distribution of bias and power and their selection and use, under the rhetorical tale of *neutrality*, tends often to confirm and strengthen them. Besides, companies that have the power to dictate standards derived from DM actually act as "disciplinary systems", producing processes of "normalization", reinforcing hidden value choices as well as unjust social structures (*injustice in, injustice out*).[39] And as Jeffrey Johnson has noticed, the full disclosure of data seems to be insufficient to get rid of the injustice embedded, although it is a necessary condition in order to let criticism and awareness begin.[40]

On the other hand, it has been also shown that (at least) some algorithms are necessarily *value-laden:* in other words, their specific design depends totally on value choices that cannot be substituted with factual represen-

tations, thus the operations supported by those algorithms, even if in general considered socially useful, are inevitably *biased*.[41]

About point 3), in contrast with the popular idea about the Internet age as an era marked by a radical boost in social equality and mass-distribution of new potentialities, scholars have noticed that, actually, new forms of inequality and discrimination are spreading, leading to specific forms of *digital divide*.[42] Since three different social blocks are emerging, with different power distribution which can lead to a threat for democracy, a three-fold distinction has been proposed: a) those who produce (whether wittingly or not) data, b) those who actually gather them, and c) those who have knowledge and technical tools to use them for large analysis and further strategic and effective decisions. Among the last group, the big IT companies unmistakably emerge as subjects who hold the key for unevenly ruling the social play.

Point 4) shows that DM share with other methods of social inquiry a peculiar feature. Since the study of social practices cannot but be realized through social practices, it is unavoidable that it ends interfering with the object itself it is supposed to study neutrally. In particular, once brought about on a large scale, it produces results that actually have complying effects on those practices under scrutiny. So, the study of a social segment of society through data analysis in order to find new useful knowledge may be seen in certain measure as a manipulation of society towards the fulfillment of specific interests, as well.

Now, all these four points have received specific attention, but here we can notice that they can be seen as consequences of a more general issue called "datafication". *Datafication* is a term used to talk about the tendency to translate every phenomenon into a quantifiable form, so that it can be measured, stored, tabulated and analyzed.[43] As already mentioned, the raw source of DM, data, is never completely "raw", but it is instead a product.[44] More precisely, it is the by-product of a series of decisions and in particular of the mother of all decisions: that is to translate every aspect of life into an objectivated and accessible *datafied* form, or to consider as relevant only what can be transform into data, that leads to the issue raised in point 1).

Then there is the decision to use DM techniques to have a differentiated consideration of those users under scrutiny, after having inserted them into differentiated lists of *targets* and *wastes,* on the basis of a knowledge they don't have access to, that leads to point 2).[45]

Then again there are all those technical decisions taken by analysts which cannot be fully shared nor discussed with final addressees, who are excluded from the possibility to take part in the determination of the way they will be treated, hence point 3).

Lastly, since the actions of few economic agents (Google, Facebook, Apple, etc.) are extremely more efficient in influencing social practices worldwide, their DM practices even when aimed at pure social inquiry are quite sensible and should be made more transparent as well as be put under critical scrutiny, hence point 4).

Resolving the world into a sea of data to sail is the philosophical point here: the cult for *datafication* which comes along with a highly structured epistemology.[46]

However, there is a more general question in the usual talk about DM which constitutes the most problematic point when it comes to build a *legal* protection against DM abuse.

As I have already hinted at, the problem lies in a specific ambiguity at the heart of the very definition of DM: the meaning of "*knowledge*".

## 4 Knowledge discovery: what exactly does it means?

In 2012 the *New York Times* reported an interesting and quite articulate survey to show how large-scale business is changing thanks to the use of DM techniques along with the clever exploitation of information derived from experimental work taken from behavioral psychology and neuroscience.

The perfect title of the article was *How Companies Learn Your Secrets* and it reported a case study that involved *Target Corp.,* one of the biggest discount store retailers in the USA. The case has become a well-known example in the literature for reflecting critically on DM.[47]

So, through DM techniques, the analysts engaged by Target came to the elaboration of a model, based on the purchase of a combination of 25 products, that assigned to the female customers of *Target* a percentage of probability to be in their first period of pregnancy. Once a high-rate customer of that kind was identified, a personalized marketing treatment was enabled, with the sending of numerous emails advertising maternity-related products. The story goes that the father of a teen-ager found those emails and brought his protests to a *Target* retailer point, receiving the excuses of a surprised executive director, who wasn't aware of the new marketing strategy. After a few days, that caring director called upon his client to apologize once more, just to hear an embarrassed father saying he discovered that his daughter was actually pregnant. The algorithm worked perfectly.

Now, the case so briefly summarized is a good one to make us focus our attention on the critical stance of the story and, in general, on the main ambiguity of the DM issue, namely the very concept of "knowledge" which is at stake.

As a matter of fact, right before the explosion of the so-called Big Data era, some Authors within the field of Computer Ethics had already highlighted the risks associated with the use of DM techniques and Knowledge Discovery in Database (KDD). When starting from a set of personal data legitimately acquired with no particular privacy issue, you can get *new* information about people "derived from implicit patterns in the data, which can suggest "new" facts, relationships, or associations about that person such as that person's membership in a newly discovered category or group".[48] Those "new" personal information were not known in advance, nor they were predictable before the use of such techniques (and so they are said to be "new", "not trivial" and "useful") and, as a consequence of that, no previous *specific* consent could have been obtained.[49] Besides, as the personal facts discovered are "new", even if a generic consent was given trying to "cover" every outcome of data manipulation, it would have been so indeterminate to be considered legally invalid according to many legal systems. For instance, the very definition of "consent" in the new GDPR requires it to be freely given, informed and *specific* (GDPR, art. 4, 11).

In the Target Corp. case we may presume that the girl (or her parents) actually had given a generic consent to the use of her data when registered in the clients' database, but we can easily and legitimately doubt that it could include even the new and quite "sensitive" information of her pregnancy. However, it is a fact that, even if she did not communicate it, that personal information was discovered and used for commercial purposes.

Now, in the context of EU legal regulation, a personal datum is "any information relating to an identified or identifiable natural person" [GDPR, art. 4, 1; Dir. 95/46/CE, art. 2, a], so, in this case, if we refer the idea of a "knowledge discovery" associated to DM techniques to the *fact* of the pregnancy, then we are forced to conclude that a specific legal basis was necessary *before* any legitimate use of that information, at least according to the EU legal systems. In other words, if "knowledge" is about an information that reflects a *fact* related to a specific person, then what *Target* discovered was actually a personal datum, and its use should have been subject to a legitimate previous legal basis allowed [GDPR, art. 6].

But the point is that this is not the only way to interpret the meaning of "knowledge" as used in these situations.

We can argue, very differently indeed, that the word "knowledge" is not associated in this case to a specific *fact* whatsoever, meaning by "fact" the existence (empirically testable) of a certain "state of affair" (Wittgenstein, *Tractatus Logico-Philosophicus*, prop. 2).[50] Instead – so the argument may go – as the relevant correlations are highlighted thanks to specific algorithms, often built on purpose, the only "knowledge" we may claim to have is about the mathematical consequences of those manipulations we made on existing data, with no direct reference to *any real* state of affairs outside.

In other words, as analytical techniques are used to create a *model* centered on an *ideal* group of elements, by gathering the correlations found in a set, what you discover is not – in our example – the information about the *actual* pregnancy of one (or more) specific girl(s). Instead, you are elaborating an abstract *percentage of probability* to be assigned to any subject-client of certain kind whose data fit – more or less – the model. Thus, in this case, "knowledge" is not referred to a *fact* but to a *property* (a *feature*) of the model, and as it was created in the exclusive interest of the company, that company can do what it wants with that. The fact that the prediction associated to the model that determined the concrete action (sending specific emails to a specific subject) *happens* to be true, does not change the nature of the "knowledge" acquired through DM, which is referred to the model and its properties and not to a specific reality whatsoever outside it. So – the argument may end – no previous consent and no any other legal basis is necessary besides that used to process the data in the first place, since no *actual fact* is the object of the "knowledge" acquired through DM. Moreover, as the prediction can apply to any number of actual people – and can even result to apply to none – it would be impossible to state in advance who they will likely be and ask them for a specific consent.

Thus, according to the first interpretation of "knowledge" *Target* should have counted on two different legal bases: one for the acquisition and processing of the starting data and another in order to use the "new" information it discovered through DM.

But according to the second interpretation, *Target* would not have needed any further legal basis, since no "new" personal data was discovered.

Again, within the specific framework of the new GDPR, if we are for the first interpretation (i. e. "what you discover through DM are new personal data"), then the Data Controller can use the newly acquired data only on a legal basis different from the explicit consent (art. 6), since if the data is really "new" no specific consent could have been previously acquired. Then, he is also obliged to give accurate information according to art. 14 ("*Information to be provided where personal data have not been obtained from the data subject*"), which is slightly more demanding than art. 13 ("*Information to be provided where personal data have been obtained from the data subject*"), highlighting in particular the legal basis for the data processing (art. 14, (1), lett. c).

Instead, if we are for the second interpretation (i. e. "what you discover is not personal data but features of your own model") then the Controller does not need an additional legal base for the new "knowledge" acquired, nor is he under more obligations than he was with reference to the previously acquired data.

It is worth noting that leading to the extreme this second interpretation of "knowledge", as a purely internal affair producing no such thing as "new" personal data, may also pose a threat to the efficacy of the entire GDPR since it applies only to personal data. In other words, all "new" information may be processed and used on their own, with virtually no limits, regardless of the destiny of the personal data used to achieve them. Even a possible "right to explanation" would be of limited protection since it would serve the data subject just to have an analytical representation of the processes applied to the original data without including other processes applied to the new information extracted.[51]

Since, these two different interpretations lead to different consequences as for the legitimacy of the particular use of the information acquired through DM, the ambiguity is legally relevant. And we may legitimately suppose that this is one of the major difficulties for a satisfactory data protection legislation that can face the issues derived from the use of specific analytical techniques in the discovery of "new", "not trivial" information.

Finally, from the point of view of the four points of general criticism highlighted on the previous section, again favoring one interpretation or another can be quite significant, since the less demanding information duties upon the data Controller may contribute to exacerbate issues of social exclusion and lack of transparency (*see above* point 3), as well as to hide even more the social inequalities embedded in biased data (point 2).

Thus, the ambiguity in the meaning of knowledge is legally, ethically and socially relevant at least with reference to two out of four general critical issues raised so far.

## 5    Reconciling the opposites: *knowledge "by correlations"* and *knowledge "by explanations"*

While such a relevant ambiguity is unbearable when it comes to focus on the need for rights protection, it seems that we actually are in the position to find a reconciliation useful to foster the analysis even more.

As a matter of fact, the term "knowledge" has a rich variety of uses in ordinary language, as well as a highly complex philosophical history and some very peculiar definitions in specialized scientific fields. Thus, in cases like this one, it is not a good practice to let ambiguities go on the loose, and, perhaps, we could also suggest not to use such a controversial term and strive to elucidate more clearly what kind of phenomena are at play in situations like the one highlighted above.

Now, it is very important to notice that the ambiguity may be caused by an insufficient clarification of the *ontology* of data: if one thinks that a datum is a neutral representation of reality, some sort of pure, objective "mirror" of the real thing, then he may easily take for granted that any correlation between data shows by itself a (substantial) causal relation among facts. But this is highly questionable and, generally, rejected.[52] Then it is a good advice to distinguish the different aspects involved and use some sort of "epistemological precautionary principle" separating two concepts here: "knowledge of *correlations*" from "knowledge of *causal explanations*".

So, in order to achieve this, we can conveniently notice that those two meanings of "knowledge" are actually referring to different things, and also, from the point of view of our case study, to different *circumstances*, that are all relevant in these complex cases. In other words, what we can criticize is the claim to deal with a one-facet epistemic experience, while, on the opposite, we have to do with different kinds of relevant knowledge.

In first place, what *Target* knows is simply that, given a set of technical procedures of manipulation of data collected in certain ways, some *correlations* can be highlighted: for instance, that some sets of products – of very different kind – are bought together under some circumstances. This can be enough to state a practical rule, by which that same correlation is expressed in a conditional way: "*if p then q*", that is "if *some sort of products in a certain arrangement* are bought, then *some more of different kind or arrangement* will (probably) be either". Obviously, this rule suggests a lot of things to do for a company interested in optimizing its economic performances. But, since the *correlations* are highlighted for what they are, that is *without any real causal investigation*, it is correct to say that "knowledge", in this case, simply express a property of a model and, at this stage, it's an internal affair.[53]

But the point is that that model has a sort of "second life", since it is used to take decisions that affects the life of people and, at this *second stage*, new information, but this time *about the model*, enrich the old ones, leading to a different kind of knowledge, and a new evaluation *of the model itself* is possible. Then, at this stage there is again some "new" knowledge to be acquired which is of a very different kind from the one reached in the first stage, but things are more sensitive, now, exactly because the life of people is crossed. Once put into practice, the hypotheses will reveal their "truth-capability": in other words, they will likely result in being true for some

instances and false for others. Either way the knowledge acquired is of different kind from the one we claimed to have in the first stage.

We must be aware of the fact that it is at this point, when the building of "knowledge" is used to take decisions, that practices based on DM may encounter other specific limitations, in particular when decisions are supposed to be based solely on automated processing. As already mentioned, the GDPR, art. 22, prohibits decisions which may have legal consequences on the data subject, or may "similarly" affect him or her, based on complete automation. This restrictive approach, which is not a novelty in itself,[54] is based on the important rationale that any decision which bears significant consequences onto others needs to be rooted not only in the analytical power of machines but also in a human intercourse.[55] However, the prohibition suffers of three relevant exceptions namely when the complete automation is necessary to perform or to enter a contract, when it is authorized by the EU law or by the law of the member state to which the controller is subject, and when it is based on the data subject's explicit consent. As for its predecessor in the Directive, we may suspect that the amplitude of these exceptional cases can severely restrict the concrete application of the general prohibition. In any case, as we are going to see, whether a model is suggesting true or false new information about a subject is quite relevant for the discipline on automated decision-making.

Anyway, the point is that models serve to have structured hypotheses to approach future instances in an organized way, by being prepared in advance to respond the needs that will arise. From a logical point of view, when facing real occurrences, models and profiles may be totally true, totally false, partially true, partially false, that is the properties assigned to the items of the classes built from analysis may match or not the single real new case totally or partially.

Since this is an *inductive* way to proceed, there is always the possibility to fail, totally or partially in single cases.

The sensitive issue here is that, at this second stage, succeeding equals to having knowledge of facts, that is, of personal data of people, while failing in single cases actually means acting towards them in a *bias,* treating like they were something different, *assigning different identities to them,* something that hurts the rights of people. Obviously, a model that happens to be partially true and false shows both kinds of problems.

It should not be underestimated that giving a valid consent to the processing of personal data, even for the elaboration of group profiles or general models through DM, does not logically imply *per se* also the legitimacy *to use* those models back towards the data subject. In other words, since *giving* data to be processed is a different thing from accepting to *be treated* like the model would predict, a consent given (or another legal base) for the first should not be supposed to include *ipso facto* also a consent for the second. This is a relevant point when it comes to discuss the case in which the models elaborated through DM happen to be actually true.

In this case, as already hinted at, there is evidence enough that some more information about the data subject have been discovered and the legitimacy of their use must undergo a specific evaluation. In particular, it must be verified if the legal basis used for processing the data in first place actually covers also the intended use for the new ones. This should be denied in principle when the legal base was the free and informed consent since it could not have been "specific" (art. 4, n. 11) with reference to data actually "new". The data subject can count on all the rights and protection ordinarily given by the GDPR and on safeguards particularly useful for this situation, for instance: the verification of the compatibility of the use the data controller is going to make of the information acquired with the original purpose of the data collection (art. 6, § 4), the right to oppose the processing when the new data are used for marketing purposes, or when the data controller invokes lett. e) or f) as legal base for the new processing (art. 21, GDPR).

Moreover, with reference to the discipline provided for by art. 22, it is clear that the consent given to the data processing (art. 6, n.1, a, art. 7) is different from (and cannot include *per se*) the consent to automatic decision making (art. 22, § 2, c). Nor the latter includes *per se* the approval of the outcome. In this case, as well as in the case of legitimate automatic decisions for performing or entering a contract, some extra rights should be guaranteed[56], in particular the highly important "right to contest" the decision (art. 22, § 3).[57]

On the other hand, when models happen to be false, the core of possible discriminating outcomes, or of unfair data treatment lies in the fact that the failure of the model result in "treating like you were somebody else". Here the "knowledge" acquired has nothing to do with real situations: it is just knowledge of features of a statistical modeling with no truth claims (*knowledge of pure correlations – among data*). Any use of it to act upon a subject is based on an identity mistake. As a consequence of group mismatching, unnecessary data can be requested, or used, or disproportionate treatments of data can be performed, with massive cascading effects onto the rights of people involved and onto the general principles of legitimate data processing, such as "fairness and transparency", "minimalization", "purpose limitation" etc. (art. 5 GDPR). In these cases, further processing of "new" (false) data is unlawful and the data subject is entitled to request either the erasure (art. 17) or the restriction of processing (art. 18, lett. b).[58]

In this situation, we may also legitimately doubt of the very legal possibility of any kind of fully automated decision making. Since art. 22, prohibits in principle decisions based solely on automated data processing, we

need to consider if it is possible at all for any of the exceptions to apply. And since it is certainly not "necessary" to conclude or perform a contract on the base of false information, nor it seems credible that a Ue or member state law can authorize automatic decision-making based on false information, the only dubious situation is the explicit consent of the data subject. Is it legal for him or her to give a valid consent to be subject to a fully automated decision on the base of knowingly wrong information about him or her? To this question my answer is no. According to art. 5, lett. d) "accuracy" of data is a general principle which must be pursued by itself, thus it is not entirely put in the hands of the data subject. As a matter of fact, art. 5 says that "every reasonable step must be taken to ensure that personal data that are inaccurate [...] are erased or rectified without delay", making "accuracy" one of the core aims of the more general accountability principle (§ 2). Of course, the data subject has a right to rectification (art. 16), but this provision is meant to reinforce the general principle and not to legitimate the processing of data which are knowingly false or inaccurate in case he or she does not exercise it. The Data controller would face responsibilities anyway for failing his or her accountability duties.

Thus, no consent can overcome this principle, no exception to the prohibition provided for art. 22 is logically possible, and therefore no fully automated decision making is legally possible in this case.

To conclude, "knowledge" is used in these cases with two significant meanings, the first being referred to "correlations" (*when some kind of products are bought, some others are as well*), and the second one to "explanations" (*somebody buys those kinds of product* because of *a current pregnancy*). The first is linked to the useful building of models and grouping, the second with peculiar causal relations between state of affairs pertaining to people, thus leading to a problem of discrimination and rights protection. Since there is no mutual exclusion between the two, they are simply looking at two different stages in the DM lifecycle endorsing two different points of view, then the problem of defending people rights from potential abuse is real and urgent.

# 6 Conclusions

It should be clear that the criticisms raised towards the contemporary spread and use of DM techniques high-light a general justified worry towards the transformations that our societies are undergoing because of the Information Revolution and, more precisely, of the *datafication* of private and social life. Those transformations go deep into the way we think about ourselves, we come into contact and relate to each other, thus touching the basic categories of construction of social reality. Nowadays the perception of the depth of those changes is stimulating the rise of a more global point of view, called "dataism", based on the "*belief* in the objective quantification and potential tracking of all kind of human behavior and sociality through online media technology" involving peculiar level of *trust in* the (institutional) agents that collect, interpret and share (meta)data culled from social media, internet platforms, and other communication technologies.[59] This is of course a point that deserves some more philosophical attention, something that cannot be offered here.

Anyway, in this situation, DM is expected to acquire more and more importance in our contemporary times and, on the basis of the analysis given so far, we can try to resume the main reasons of that in (at least) five points:

1. *The amount of Data produced is constantly growing and changing.*

   The more interesting and the most promising challenges for the DM are to be found in the exploitation of the World Wide Web as the biggest data base of all time. But the problem is that this is far from being a static data base, instead it is constantly growing and changing. As a consequence, any correlation found is used under the never-certain presupposition of its stability,[60] but it is in constant need of fine-tuning and improvement. This fact alone leads to an intensification of DM to sharpen and adjust the knowledge acquired.

2. *Data are personal footprints giving hints about subjects we may want to have some more knowledge about.*

   Data are the product of man. They say something about him, even if it can be quite tricky to specify exactly what.[61] They are traces of his actions left in the digital environment which has become the main *habitat* and global ecosystem for contemporary civilized communities. Thus the more he deploys his own life in it, the more he leaves a datafied projections of himself. So, making unitary sense of these projections needs appropriate tools. That's why DM is for.

3. *Models based on data are "somewhat" justified tools for decisions making.*

   DM offers a kind of "knowledge" (in both senses of "correlations" and "causal relations") highly dependent on the quality of data as well as the techniques used to get it, so it actually offers the possibility to back a certain decision with reference to a complex series of informations: such as the data sets used as well as the technical procedures applied. Two points are important here:

    **a.** since those procedures are in the hands of qualified experts who master them but are not necessarily the subject in charge for relevant decisions about future actions to be taken, the ability to give a complete "back up argument" is split into a wide array of people, making it difficult to be re-collected promptly. As a result, very often the decisions taken is "opaque" to people whose life is touched by them;

    **b.** even if at the disposal of a "collective mind", because of the constantly changing amount of relevant data, the reasons behind the decisions taken are highly defeasible and can be falsified from any variation on the dataset.

**4.** *Models based on group profiling, create the basis for general policies so…*

**5.** *… they can generate a "modeling loop".*

As we saw earlier, DM is often used by powerful economic agents, in particular IT companies with the aim to build or to rely – for their business – on general models of behavior of masses of people interacting with each other. This fact creates the "modeling loop" we talked about, raising concerns for surveillance issues as well as for human auto-determination and freedom in general.

Finally, we are now in the position to list the points in need of future organic research in order to deepen the ethical, legal and social analysis on DM:

    **a.** deepening the dynamics of the "modeling loop" with particular reference to "virtual" groups;

    **b.** introducing a "epistemological precautionary principle" (EPP) in defining DM and data-related knowledge;

    **c.** stepping into specific DM techniques and algorithms to classify concrete ethical and legal risks;

    **d.** creating within any medium-to-large business administration interdisciplinary committees to monitor the elaboration of models and cooperate with academic researchers on the subject.

We may start from here to open new interdisciplinary fields of research on the evolution of our "quantified selves".[62]

## Notes

1  M. Bramer, *Principles of Data Mining* (3rd ed. Springer-Verlag London Ltd, 2016), 2; A. K. Maheshwari, *Business Intelligence and Data Mining* (Business Expert Press, 2015), 4–5; F. Coenen, "Data Mining: Past, Present and Future," *The Knowledge Engineering Review* 00, no. 0 (2004): 1–24.

2  G. Myatt, *Making Sense of Data* (Wiley and Sons pub, 2007).

3  S. H. Liao, P. H. Chu, and P. Y. Hsiao, "Data Mining Techniques and Applications – A decade Review from 2000 to 2011," *Expert Systems and Applications* 39 (2012): 11,303–11; C. Vercellis, *Business Intelligence: Data Mining and Optimization for Decision Making* (John Wiley & Sons, 2009).

4  V. Mayer-Schoenberger and K. Cukier, *Big Data: A Revolution That Will Transform How We live, Work and Think* (J. Murray, 2013).

5  McKinsey Global Institute, "The Age of Analytics: Competing in a Data-Driven World," 2016, https://www.mckinsey.com/business-functions/mckinsey-analytics/our-insights/the-age-of-analytics-competing-in-a-data-driven-world.

6  *Ib.*

7  L. L. Visinescu et al., "Improving Decision Quality: The Role of Business Intelligence," *Computer Information Systems* 57, no. 1 (2017): 58; J. Pl. Shim et al., "Past Present and Future of Decision Support Technology," *Decision Support Systems* 33, no. 2 (2002): 111; U. Bose, "An Ethical framework in Information Systems Decision Making Using Normative Theories of Business Ethics," *Ethics and Information Technology* 14 (2012): 17; L. van Wel and L. Royakkers, "Ethical Issues in Web Data Mining," *Ethics and Information Technology* 6 (2004): 129.

8  Cfr. Interview in, *La Stampa,* August, the 1st, 2016.

9  M. Bakaev, T. Avdeenko, "Prospects and Challenges in Online Data Mining. Experiences of Three-Year Labour Market Monitoring Project," in *Data Mining and Big Data. Proceedings of the First International Conference on Data Mining and Big Data. Bali, Indonesia, June 25–30, 2016*, ed. Y. Tan and Y. Shi (Springer International Publishing, 2016).

10  For an introduction to profiling and for the differences between individual and group profiling see F. Kaltheuner, E. Bietti, "Data is Power: Towards Additional Guidance on Profiling and Automated Decision-Making in the GDPR," *Journal of Information Rights, Policy and Practice* 2, no. 2 (2017).

11  H. Boghosian, *Spying on Democracy* (San Francisco: City Light Books, 2013).

12  B. Mittlestadt, "From Individual to Group Privacy in Big Data Analytics," *Philosophy and Technology* 30, no. 4 (2017): 475.

13  G. Simmel, "How is Society Possible," *American Journal of Sociology* 16, no. 3 (1910): 372–91.

14  *Ibidem,* p. 388.

15  See M. Hildebrandt, "Defining Profiling: A New Type of Knowledge?" in *Profiling the European Citizen. Cross-disciplinary Perspectives,* eds. M. Hildebrandt and S. Gutwirth (Springer, 2008): 20. The conditioning effect of group pressure upon individual judgment is a traditional sociological and psychological topic, see S. Asch, "Effects of Group Pressure Upon the Modification and Distortion of Judgment," in *Groups, Leadership, and Men,* ed. H. Guetzkow (Carnegie Press, 1951). When groups and social roles are also presented with differentiations in authority some more specific and predictable dynamics of moral behaviours are possible. Here the studies on obedience by Stanley Miligram, and the famous *Stanford Prison Experiment* by Philip Zimbardo come to mind.

16  McKinsey Global Institute, "A Future that Works. Automation, Employment, Productivity," 2017, https://www.mckinsey.com/mgi/overview/2017-in-review/automation-and-the-future-of-work/a-future-that-works-automation-employment-and-productivity.

17  A. Lai, "What Data Scientist Can Learn From History," in *Real World Data Mining Applications,* eds. M. Abou-Nasr et al. (Springer, 2015).

18  G. Malgieri and G. Comandé, "Why a Right to Legibility of Automated Decision-Making Exists in the General Data Protection Regulation," *International Data Privacy Law* 7, no. 4 (2017): 243, 244.

19  H. Moor James, "*What is Computer Ethics?*," in *Computer Ethics,* ed. T. W. Bynum, (Blackwell, 1985).

20  J. Turow, *The Daily You. How the Advertising Industry is Defining Your Identity and Your Worth* (Yale University Press, 2011); C. F. Chien and L. F. Chen, "Data Mining to Improve Personal Selection and Enhance human Capital: A Case Study in High Technology Industry," *Expert Systems with Applications* 34, no. 1 (2008): 280; D. Zhang and J. Deng, "The Data Mining of the Human Resources Data Warehouse in University Based on Association Rule," *Journal of Computers* 6, no. 1 (2011): 136; H. T. Tavani, "Informational Privacy, Data Mining and the Internet," *Ethics and Information Technology* 1, no. 2 (1999): 137.

21  Just to make some few examples, see Y. Chen et al., "Reality Mining: A Prediction Algorithm for Desease Dynamics Based on Mobile Big Data," *Information Sciences* 379 (2017): 82–93; M. Iniadat et al., "Data Mining Techniques in Social Media: A Survey," *Neurocomputing* 214 (2016): 654–70; K. Lakiotaki et al., "BioDataome: A Collection of Uniformly Preprocessed and Automatically Annotated Datasets for Data-Driven biology," 2018, Database, https://doi.org/10.1093/database/bay011; J. Rung and A. Brazma. "Reuse of Public Genome-Wide Gene Expression Data," *Nature Revue Genetics* 14 (2012): 89–99.

22  C. Anderson, "The End of Theory: The Data Deluge Makes the Scientific Method Obsolete," *Wired* 2008.

23  See, for instance, recently D. Lupton and B. Williamson, "The Datafied Child: The Dataveillance of Children and Implications for Their Rights," *New Media and Society* 19, no. 5 (2017): 780–94.

24  In the EU General Data Protection Regulation, for instance, the legal basis for a lawful processing are to be found in art. 6.

25  See art. 4, 13, 14 and 22. See P. de Hert and V. Papakonstantinou, "The new General Protection Regulation: Still a Sound System for the Protection of Individuals?" *Computer Law and Security Review* 32 (2016): 179.

26  M. Hildebrandt, "Profiling and the Rule of Law," *Identity in the Information Society* 1, no. 1 (2008): 55–70.

27  S. Wachter, B. Mittlestadt, and L. Floridi, "Why a Right to Explanation of Automated Decision-Making does not Exist in the General Data Protection Regulation", *International Data Privacy Law* 7, no. 2 (2017): 76.

28  More on the relevance of art. 22, *infra* § 5.

29  A. Marwick, "The Public Domain: Social Surveillance in Everyday Life," *Surveillance and Society* 4 (2012): 378; C. Fuchs, "New Media, Web 2.0 and Surveillance," *Sociology Compass* 5, no. 2 (2011): 134–47.

30  B. Goodman and S. Flaxman, "European Union Regulations on Algorithmic Decision-Making and a "Right to Explanation"" *AI Magazine* 38, no. 3 (2017): 50; S. Barocas and A. D. Selbst, "Big Data Disparate Impact," *California Law Review* 104 (2016): 671; C. O'Neal, *Weapons of Math Destruction. How Big Data Increases Inequality and Threatens Democracy* (Crown, 2016).

31  D. Boyd and K. Crawford, "Critical Questions for Big Data: Provocations for a cultural, Technological, and Scholarly Phenomenon," *Information, Communication, Society* 15, no. 5 (2012): 662.

32  H. Kennedy, *Post, Mine, Repeat. Social Media Data Mining Become Ordinary* (Palgrave McMillan, 2016); J. Turow, *The Daily You. How the Advertising Industry is Defining Your Identity and Your Worth* (Yale University Press, 2011); J. Law and J. Urry, "Enacting the Social," *Economy and Society* 33, no. 3 (2004): 390.

33  A. Marwick, "The Public Domain: Social Surveillance in Everyday Life," *Surveillance and Society* 4 (2012): 378.

34  A. Albrechtslund, "Online Social Networking' as Partecipatory Surveillance," *First Monday* 13, no. 3 (2008), https://doi.org/10.5210/fm.v13i3.2142.

35  L. E. Sherman et al., "What the Brain 'likes': Neural Correlates of Providing Feedback on Social Media," *Social Cognitive and Affective Neuroscience* 13, no. 7 (2018): 699–707; A. L. Sherman et al., "The Power of the *Like* in Adolescence: Effects of Peer Influence on Neural and Behavioral Responses to Social Media," *Psychological Science* 27, no. 7 (2016): 1027–35.

36  A. Haase and E. Peters, "Ubiquitous Computing and Increasing Engagement of Private Companies in Governmental Surveillance," *International Data Privacy Law* 7, no. 2 (2017): 126.

37  L. Portness and S. Tower, "Data Barns, Ambient Intelligence and Cloud Computing: The Tacit Epistemology and Linguistic Representation of Big Data," *Ethics and Information Technology* 17, no. 1 (2015).

38  D. Boyd and K. Crawford, "Critical Questions for Big Data: Provocations for a Cultural, Technological, and Scholarly Phenomenon," *Information, Communication, Society* 15, no. 5 (2012): 662, 664.

39  Jeffrey A. Johnson, "From Open Data to Information Justice," *Ethics and Information Technology* 16 (2014): 263.

40  *Ib.*

41  F. Kraemer, K. van Overveld, and M. Peterson, "Is There an Ethics of Algorithm?" *Ethics and Information Technology* 13 (2011): 251.

42  H. Kennedy, *Post, Mine, Repeat. Social Media Data Mining Become Ordinary* (Palgrave McMillan, 2016), 52–53; D. Boyd and K. Crawford, "Critical Questions for Big Data: Provocations for a Cultural, Technological, and Scholarly Phenomenon," *Information, Communication, Society* 15, no. 5 (2012): 662, 674–75.

43  V. Mayer-Schoenberger and K. Cukier, *Big Data: a Revolution That Will Transform How We live, Work and Think* (J. Murray, 2013), ch. 5.; S. Baak "Datafication and Empowerment: How the open data Movement Re-articulates Notions of Democracy, Partecipation and Journalism," *Big Data and Society* (2015): 1–11; D. Holtzhausen, "Datafication: Threat or Opportunity for Communication in the Public Sphere?" *Journal of Communication Management* 20, no. 1 (2016): 21–36; R. Kitchin, "The Real time City? Big Data and Smart Urbanism," *Geojournal* 79 (2014): 1–14; M. Lycett, "Datafication': Making Sense of (big) data in a Complex World," *European Journal of Information Systems* 22 (2013): 381–86; S. Mattern, "Methodolatry and the Art of Measure. The New Wave of Urban Data Science," *Places Journal* (2013),https://doi.org/10.22269/131105; J. Van Dijk, "Datafication, Dataism and Dataveiilance: Big Data between Scientific Paradigm Ideology," *Surveiilance & Society* 12, no. 2 (2014): 197–2008.

44  L. Gitelman, ed., "*Raw Data" is an Oxymoron* (The MIT Press, 2013).

45  M. Hildebrandt, "Defining Profiling: A New Type of Knowledge," in *Profiling the European Citizen. Cross-Disciplinary Perspectives,* eds. M. Hildebrandt and S. Gutwirth (Springer, 2008).

46  L. Portness and S. Tower, "Data Barns, Ambient Intelligence and Cloud Computing: The Tacit Epistemology and Linguistic Representation of Big Data," *Ethics and Information Technology* 17, no. 1 (2015).

47  A. K. Maheshwari, *Business Intelligence and Data Mining* (Business Expert Press, 2015), 46; H. Kennedy, *Post, Mine, Repeat. Social Media Data Mining Become Ordinary* (Palgrave McMillan, 2016), 1; Luciano Floridi, "Big Data and Their Epistemological Challenge," *Philosophy and Technology* 25 (2012): 435, 436.

48  E. H. Tavani, "KDD, Data Mining, and the Challenge for Normative Privacy," *Ethics and Information Technology* 1 (1999): 265.

49  H. T. Tavani, "Informational Privacy, Data Mining and the Internet," *Ethics and Information Technology* 2 (1999): 137, 128; M. Sax, "Big Data: Finders keepers, loosers weepers?" *Ethics and Information Technology* 18 (2016): 25, 29.

50  Ludwig Wittgenstein, *Tractatus Logico-Philosophicus* (Routledge, 2001), 5.

51  That of the existence of a full "right to explanation" is perhaps the major point of discussion in the literature interested in the protection towards automated processing of personal data. At the moment four positions are on the table: the first denies such an existence: see S. Wachter, B. Mittelstadt, and L. Floridi, "Why a Right to Explanation of Automated Decision-Making Does Not Exist in the General Data

Protection Regulation," *International Data Privacy Law* 7, no. 2 (2017): 76–99. The second one admits it instead: Bryce Goodman and e Seth Flaxman, "European Union Regulations on Algorithmic Decision-Making and a "Right to Explanation," *AI Magazine* 38, no. 3 (2 ottobre 2017): 50–57; Maja Brkan, "Do Algorithms Rule the World? Algorithmic Decision-Making in the Framework of the GDPR and Beyond," SSRN Scholarly Paper. Rochester, NY: Social Science Research Network, 1 agosto 2017. https://papers.ssrn.com/abstract=3124901. The third one introduces some sort of more limited kind of rights: G. Malgieri and G. Comandé, "Why a Right to Legibility of Automated Decision-Making Exists in the General Data Protection Regulation," *International Data Privacy Law* 7, no. 4 (1 novembre 2017): 243–65. Finally there are scholars who think that this is not a major point since the best interest of the data subject can be pursued by other tools introduced by the GDPR: L. Edwards and M. Veale, "Slave to the Algorithm? Why a "right to an Explanation" Is Probably Not the Remedy You Are Looking For," *Duke Law and Technology Review* 16: 18–84.

52  "Facts are ontological, evidence is epistemological, data is rhetorical [...] When a fact is proven false it ceases to be a fact. False data is data nonetheless", D. Rosenberg, "Data before the fact," in "*Raw Data*" *is an Oxymoron*, ed. L. Gitelman (The MIT Press, 2013), 18.

53  B. Goodman and S. Flaxman, "European Union Regulations on Algorithmic Decision-Making and a "Right to Explanation," *AI Magazine* 38, no. 3 (2017): 50; M. Hildebrandt, "Defining Profiling: A New Type of Knowledge," in *Profiling the European Citizen. Cross-Disciplinary Perspectives*, eds. M. Hildebrandt and S. Gutwirth (Springer, 2008), 18.

54  As a matter of fact Directive 95/46/CE included a similar provision (art. 15) which did not receive particular attention nor it had any significant practical application, see Lee A. Bygrave, "Automated Profiling: Minding The Machine: Article 15 of The EC Data Protection Directive And Automated Profiling," *Computer Law & Security Review* 17, no. 1 (2001): 17–24.

55  The current interpretation of this provision highlights the need for a significant human intervention, in other words, the decision is considered based "solely" on automated processing not only when no human takes part in it but also when the human participation is completely formalistic, lacking the authority to change the course of action indicated by the machine, see ART29WP, *Guidelines on Automated individual decision-making and Profiling for the Purposes of Regulation 2016/679*, last Revised and Adopted on February 6, 2018, p. 21; Ben Wagner, "Liable, but Not in Control? Ensuring Meaningful Human Agency in Automated Decision-Making Systems," *Policy & Internet* 11, no. 1 (2019): 104–22.

56  In particular the right to ask for a human intervention and the right to express his or her own opinion (art. 22, § 3).

57  The relationships among these specific rights is yet to be studied thoroughly. For instance, we may ask if they are supposed to be used cumulative or in a mutually exclusive way: in particular, once the data subject makes use of the right to human intervention, even if the decision taken by the machine is confirmed, the entire process is no more completely automated, and then it is no more subject to the discipline of art. 22. As a consequence, it may be argued, the data subject looses his right to contest the decision which is a special safeguard in the case of complete automation. Symmetrically, once the data subject uses his right to contest, there is no apparent limitation in the GDPR to a complete automation in the management of the objections raised as long as it falls within the exceptional cases of art. 22, § 2. The only way to avoid this chain of automations would be to exercise the right to human intervention, with the above mentioned consequence of loosing the right to contest. For a complete discussion on this point and in particular on the "right to contest", cfr. C. Sarra, "Il diritto di contestazione delle decisioni automatizzate nel Regolamento Europeo sulla Protezione dei Dati Personali," *Anuario de la Facultad de Derecho de L'Universidad de Alcalà de Henares (forthcoming)* (2019).

58  It must be taken into consideration that, within the legal framework we are considering, inaccurate information about an identifiable subject are personal data nonetheless: as a matter of fact art. 4, n. 1, does not consider the accuracy or full truthfulness in the definition of "personal data", and that is why there is a "right to rectification" included in the GDPR (art. 16). Of course "accuracy" is a general principle for lawful processing, as we are going to see, but as far as the definition of personal data goes, it is not a necessary requirement.

59  J. Van Dijk, "Datafication, dataism and dataveiilance: Big Data between scientific paradigm ideology," *Surveiilance & Society* 12, no. 2 (2014): 197–208, 198. See also S. Lohr, *Data-ism. The revolution transforming decision making, consumer behavior and everything else* (London: Oneworld Publication, 2015); Y. N. Harari, *Homo Deus. A Brief History of Tomorrow* (HyperCollins, s/l 2017); P. Ted, "The Deluge of Dataism: a New Posthuman Religion?" *Dialog: A Journal of Theology* 53, no. 3 (2017): 211–13; B. Nicholls, "Everyday Modulation: Dataism, Health Apps, and the Production of Self-Knowledge," in *Security, Race, Biopower. Essays on Technology and Corporeality*, eds. H. Randell-Moon and R. Tippet (London: Palgrave Mcmillan, 2016), 101–20.

60  Maheshwari, *Business Intelligence*, cit., 45.

61  For a de construction of both 'data' and 'digital traces' see T. Reigeluth, "Why data is not enough," *Surveillance and Society* 12, no. 2 (2014): 243–54.

62  M. Ruckenstein and M. Pantzar, "Beyond the Quantified Self: Thematic exploration of a dataistic paradigm," *New Media and Society* 19, no. 3 (2015): 401–18; M. Swan, "The Quantified Self: Fundamental Disruption in Big Data Science and Biological Discovery," *Big Data* 1, no. 2 (2013): 85–99.