

***On Truth and Data Science:
A Response for Dora¹***

Michael L. Brodie

Computer Science and Artificial Intelligence Lab, Database Group

August 15, 2017

Is There Absolute Truth?

Saying that “There is no absolute truth” sounds intuitively wrong, even shocking. Who would say such a thing? Alternatively, it may sound very philosophical - an intellectually fascinating puzzle. From a rational, scientific perspective you know that it is reasonable but your world view may say that it is wrong or even that you are violating a basic belief.

Our World View comes from Our Family of Origin in our Childhood

When we were young our parents taught us what the color blue was, and red, and green. They also taught us to tell the truth, and may have told us to always tell the truth and to stick to the facts. So, there must be a “truth”! They may have exposed us to religious beliefs that not only emphasize “truth” but demand that we honor truth in our daily lives. By the time we were three or four years old we had learned most of what we know about the world from our family of origin, i.e., parents, siblings, relatives, neighbors, TV, church, etc. By four years old, we knew millions of facts and even implicitly understood physics, chemistry, optics, and much more since you had probably conducted thousands of experiments with water – in the bath, at the beach, with your food – and had a very good intuitive understanding of liquids, viscosity, surface tension, specific gravity, and much more. We also learned what love and affection were. They were what you experienced from those who “loved” you – your parents, grandparents, relatives, etc. Love, as learned from their parents, for some, could be a wonderful, balanced, natural and rewarding feeling of attachment. For others, it could be a feeling of remoteness.

By three or four years old we have a sophisticated “world view” that we understand implicitly, in which facts and truth are so important that there was never any question as to the existence of truth.

Our World View Requires Constant Revision

Our world view becomes the basis of who we are and how we think. We may not even understand that we have specific beliefs until later in life when some beliefs fail us or we reflect on them and find flaws. Some beliefs are more realistic and helpful than others. Some may inhibit our growth “A child should be seen and not heard”, “You are a bad girl”; and some may be supportive “Go Dora, you can do it”; “You are perfect in every way.”

Let’s revisit the color blue. Is the blue that you know “true” blue? Well, simply look up the definition for blue in terms of the part of the visible spectrum that is defined as blue. Wikipedia says: “Blue is the colour between violet and green on the optical spectrum of visible light. Human eyes perceive blue when observing light with a wavelength between 450 and 495 nanometers, which is between 4500 and 4950 ångströms.” While that is a very precise definition, a lady who weaves cotton in Indonesia knows many shades of blue that may fall outside that strict definition due to her Indonesian cultural preference for a variety of shades of blue. Your parents who taught you blue probably did not know either of the above definitions. There are many blues and many “truths”. In any case, all such definitions are human artifacts and “blueness” exists independently of any such definition.

Faulty Reasoning

How can an infant at a few months, at 1 year, at 2, at 3 reason about any of the truths that they are taught or see around them? Fire does burn. Love is very important. Things fall to the ground. While the “laws of nature” seem real because you experience them consistently, many family or social truths are harder to understand. Love may be expressed consistently in the family but very differently in my friend’s family. This may seem confusing

¹ This was a discussion with Dora (Θεοδώρα Μπουκουρά), an inquisitive Greek Biology student interested in Data Science, following my keynote: *Data: The World’s Most Valuable Resource*, 2017 Onassis Lectures in Computer Science on Big Data and Applications, Heraklion, Crete, Greece, July 10, 2017

or uncomfortable. Your beliefs may be so deeply ingrained that you do not even ask “Is love as I learned it true?” but rather reject their form of love. A child may think “Since I seem to get things wrong so much and my mother says that I am wrong so much, I must be very stupid.” This is a typical case of faulty learning. How could a small child reason about complicated things? Or perhaps the child knew exactly what was going on but the parent denied it. This example of being stupid is from my childhood. With two PhDs it must be false, yet at my age those thoughts are still in my world view and require work to overcome. I give you this example to suggest that much of our world view created at a young age; is deeply embedded in our lives; may be based on faulty learning; and requires revision.

Truth: A Convenient Tool

Are the truths we were taught really true? For many practical reasons, it is convenient for parents, the church, governments, and schools to assert truths such as “stay away from the fire or you will get burned”, or “love is forever, so divorce is a sin”, and “life is sacred so abortion is a sin”. What is a sin? It is a violation of a truth asserted by some authority. As we get older we find that we no longer believe in some of the truths that we were taught. But people seldom go beyond individual truths to question *truth* itself. *What is truth?*

What is Truth?

Let’s look to science to see a definition of truth. The scientific method is an agreement amongst scientists about how to establish facts and theories about the natural world (the scientific method applies more broadly but let’s keep it simple). The scientific method says that you develop a hypothesis and design and conduct an experiment according to the scientific method and produce a “statistically significant” result that confirms or denies the hypothesis. If the reasoning holds up under scrutiny by the relevant scientific community, e.g., review by a scientific journal, the result is accepted by the community as valid within the defined bounds of the hypotheses, error bounds, precision, etc. – such results are highly qualified or conditional. Formally, no scientist would say that the result was “true” but informally scientists and the broader community treat such results as true, e.g., quantum theory, the standard model of high energy particle physics, and “central dogma of molecular biology” are informally “true”. However, in most natural sciences, the fundamental theories (beliefs) change radically every 50 to 100 years. This generally means that the previous “truths” are now false, replaced by new “truths”. New theories are not always totally false. Although the Neils Bohr’s model of the atom has been supplanted by other models, its underlying principles remain valid.

The scientific method is an example of a community established agreement as to how to define and validate observations (facts) and theories about the natural world, informally an agreement on establishing “truth”. The scientific method contains within it the ability to question and revise established truth (Bohr’s model versus the current standard model of particle physics). This highly respected method established scientific “truth” relative to what we agree that we know. It does not assert absolute truth.

Has the world view of your childhood been revised by you in your consciousness or in your family of origin? This would be ideal, suggesting that there is no absolute truth, yet would be a rare occurrence.

My view is that truth is as a community agreement under specific conditions. Such a truth requires a convention or an agreement in the community to establish and revise as needed, based on our current knowledge and observable facts, what is accepted as “true” by the community according to the conditions of the agreement and of the community. I incorporate this into my definition of Data Science (see below).

To exist, a society or community requires community-wide agreements. A simple example is traffic lights that enable traffic to flow and to avoid accident at intersections. More complex agreements are the laws of the land. *De facto*, these are the truths that enable the society to operate. Totalitarian states and some religions impose the agreements. In principle, societal agreements emerge from the citizens in a “true” a democracy and the agreements must be reviewed and extended or improved as we learn more about the phenomena or the society.

This view of truth is, like science, extremely pragmatic. A richer intellectual pursuit of truth goes back to the beginning of philosophy and will likely continue forever (see [Wikipedia](#)).

Why is an understanding of truth important in Data Science?

The informal, social notion of absolute truth that often lies deep in our consciousness has no place in data science. Such a deeply embedded belief may inhibit our ability to discover systematically observable properties

of the phenomena, i.e., plausible models, that we are exploring since belief in a single truth violates the scientific method by anticipating a specific outcome, i.e., creating a bias towards a specific outcome.

In data science, as in science, we are trying to discover plausible hypotheses that might be proven under specific conditions to systematically recur, e.g., a cancer that occurs in a specific context due to specific factors and might be minimized or eliminated by means of a specific treatment.

Data science involves hypothesizing (theory-driven or top-down) or discovering (data-driven or bottom up) systematically observable properties of a phenomenon, i.e., a model, under specific conditions.[2] Belief in an absolute truth may suggest that there is only one model, one set of properties. As Plato's allegory of the cave taught us 2,400 years ago, we cannot observe the "real" thing, we can observe only an image of the thing (a model) from our perspective. In science, as in life, understanding of a phenomenon may be enriched by observing the phenomenon from multiple perspectives (models). A recent scientific trend, e.g., pursued in biology by [Pardis Sabeti](#) at Harvard involves a shift from understanding a phenomenon with one theory (perspective) to using multiple theories or models in what is called ensemble modeling. Ensemble modeling is the process of running two or more related but different analytical models and then synthesizing the results into a single score or spread in order to improve the accuracy of predictive analytics and data mining applications.

Related Observations

My definition of Data Science is based on the definition of the scientific method as a process of acquiring new knowledge, and correcting and integrating previous knowledge, meaning that it is part of a continuous discovery process.

Data Science is a body of principles and techniques for applying data analytics to accelerate the investigation of phenomena by acquiring new data, correcting and combining it with previous data, with measures of correctness, completeness, and efficiency of the derived results (correlations) with respect to some pre-defined (theoretical, deductive, top down) or emergent (inductive, bottom up) specification (scope, question, hypothesis, requirement).[3]

What versus Why: Data Science can be used to discover correlations (What phenomena occurred) but cannot be used to establish causality (Why the phenomena occurred).

Data Science involves discovering **What** – significant facts or patterns concerning phenomena. These are called correlations amongst variables. Ideally, data science methods will help us identify highly probably (plausible) hypotheses (correlations) that will be proven causal by other means. Data Science involves accelerated methods of discovering THAT correlations occur under certain conditions and with certain probabilities; it cannot discover **Why** – whether the correlations between variables are causal, i.e., explain why the observed correlations occurred. Once data science has been used to establish one or more highly probable hypotheses (correlations), we put aside data science and turn to the conventional methods of the domain in question to establish causality or Why the observed phenomena occurred.

Single Version of Truth: Banks must maintain a single version of truth for your bank account, not multiple versions, since you want the bank to make sure that every euro you put in is credited to you and every euro taken out is credited to the person you are paying. "One version of truth" applies to most businesses that want a persistent, reliable record of all business transactions. Databases were first developed for banking and business; hence they claim to support a single version of truth. While this is critical for some problems, e.g., business transactions, it is not true for most of the rest of the world. Hence, database products do not support multiple models, i.e., the reality of science and life in general. For over 40 years, researchers have tried but failed to develop databases that support multiple perspectives or multiple semantic models.

Most assertions are unprovable: 98% of what people say are opinions that are impossible to prove as "true". The previous sentence is an opinion, hence unprovable. However, it suggests that almost all assertions are mere opinions and should be considered as opinions.

What is a bias? Understanding a phenomenon means that we have knowledge of the phenomenon. Following the above discussion of truth, our knowledge – ideally verifiable, systematic observations under specific conditions – is relative to the data we have and the models (perspectives) that we have used to establish the knowledge (informally truths of the phenomenon). Recently, it has been observed that algorithms used in many areas (mortgage and loan approvals, hiring and promotion, parole and sentencing) are biased. To be biased

means to be prejudiced in favor of or against one thing, person, or group compared with another, usually in a way considered to be unfair. For example, automated parole systems have been shown by [ProPublica](#) to be biased. Specifically, ProPublica showed that automated parole systems systematically made parole decisions that disadvantaged minorities, i.e., blacks and females, with all other factors being equal. In the terms used above, the automated parole system is based on a model that is inconsistent with a community model for fairness towards minorities. This could be that the firms that designed the systems believe based on some evidence (i.e., knowledge) that minorities recommit more than whites and males. Even though the political disposition of the community is that minorities need to be treated fairly just like non-minorities, i.e., receive the same sentences. When is knowledge biased? When the knowledge used to produce a model is in conflict with another model, then the two models are biased with respect to each other. Assuming the knowledge on which the parole system model is based is verifiable under the conditions in which it is applied, it is biased with respect to a model based on fairness to minorities that may be a political aspiration rather than a reality. How do you prove veracity of a model? In this case what is the recidivism rate for the automated parole system versus a parole system based on a model that supports political fairness to minorities? If the original parole system model has a better recidivism rate than that of the fairness model, does society select better recidivism over fairness? This is a modelling question that is outside the realm of data science. A deeper question is how do you detect bias in algorithms? You need to evaluate and compare the models underlying the algorithm versus some other model. Only models can be biased with respect to each other.

Quotes

Marcus Aurelius, 121-180 AD.

- "Everything we hear is an opinion, not a fact. Everything we see is a perspective, not the truth."

Shakespeare 1564-1616: The Tragedy of Hamlet, Prince of Denmark, Act 2, scene 2

- Hamlet: Why, then, 'tis none to you, for there is nothing either good or bad, but thinking makes it so. To me it is a prison. [A reference to Hamlet's earlier "Denmark's a prison."]
- Modernized: Well, then it isn't one to you, since nothing is really good or bad in itself—it's all what a person thinks about it. And to me, Denmark is a prison.

References

- [1] M. Braschler, T. Stadelmann, K. Stockinger (Eds.), "Applied Data Science - Lessons Learned for the Data-Driven Business", Berlin, Heidelberg: Springer, expected 2018
- [2] M.L. Brodie, Necessity is the Mother of Invention: On Developing Data Science, to appear in [1]
- [3] M.L. Brodie, What is Data Science? to appear in [1]