

Back to the Future: A Framework for Modelling Altruistic Intelligence Explosions

Mahendra Prasad

University of California, Berkeley, Department of Political Science, PhD Candidate
mrprasad@berkeley.edu

Abstract

A necessary condition of an intelligence explosion is that for some large number of possible beliefs, the updated probability of each of those beliefs being true has greatly increased (perhaps close to 1) over a relatively short time. In the 18th century, the French mathematician, Nicolas de Condorcet, proposed a model for how collective intelligence could be used to determine facts with near certainty. That model is today known as Condorcet's jury theorem. Going back to Condorcet's 18th century model, we will update it to provide a proof of concept of how to model intelligence explosions for the social good. The model will pinpoint the kind of social and political institutions that must be in place for the kind of intelligence explosion, described by the model, to occur. The hope is that by providing this kind of proof of concept model, future research will be able to tweak, add, or remove different model assumptions to better fit the circumstances with which researchers are concerned, and figure out what kinds of institutions would need to be necessary to facilitate an altruistic intelligence explosion.

Introduction

In recent years, several philosophers, scientists, and engineers, including Nick Bostrom, David Chalmers, Ray Kurzweil, and Stuart Russell, have put forth models of the intelligence explosion hypothesis. While there are many versions of it, the basic idea is that rapidly growing knowledge and technology will radically change human economic, political, social, and biological structures. Occasionally, Nicolas de Condorcet is mentioned as a proto-discoverer of the hypothesis based on his rough descriptions of it in his philosophical work. What goes unmentioned is that he actually created a mathematical model for an explosion in his mathematical work. Much of Condorcet's work in economic/political/social philosophy can be understood as an attempt to figure out how to live in a world with an intelligence explosion. But sadly, much of this work is unavailable to the English speaking world be-

cause many of his important mathematical and philosophical works, including *Essay on the Application of Analysis to the Probability of Majority Decisions* (1785) and *Essay on the Constitution and Functions of Provincial Assemblies* (1788) have not had complete English translations. In this paper, I hope to elucidate Condorcet's model of an intelligence explosion to encourage translation and study of Condorcet's important contributions by scholars who are tackling important issues related to the intelligence explosion hypothesis.

Background

In his *Sketch of a Historical Picture of the Progress of the Human Mind* (1795), the 18th century French mathematician/philosopher/revolutionary Nicolas de Condorcet argued that human intellectual/technological/ethical knowledge growth was faster than linear. He tried to demonstrate this by historically showing how past and then-contemporary human knowledge was increasing at an increasing rate. *Sketch* would go on to assert that these trends could potentially occur indefinitely leading to radical changes in human social, economic, political, and biological structures (Condorcet 2004).

In response Thomas Malthus wrote *Essay on the Principle of Population*, as it affects the future improvement of society with remarks on the speculations of Mr. Godwin, M. **Condorcet**, and other writers (1798) (bold added for emphasis). In *Essay* (1798), Malthus accused Condorcet's *Sketch* of assuming that past performance implies future results (i.e. past and current geometric growth in knowledge do not imply future geometric growth) (Malthus 1976, 65-75).

But critics, such as Malthus, rarely if ever read Condorcet's work on social choice, like his *Essay* (1785), where through a corollary to his jury theorem, he demonstrated conditions under which human knowledge could asymptotically but quickly approach perfection (i.e. a probability of

1 of being correct). Condorcet's jury theorem showed that if jurors had a probability greater than $\frac{1}{2}$ of correctly judging whether or not some statement is true, and if each juror's judgment of truth is statistically independent of other jurors' judgments, and each juror sincerely expressed their judgment, then the majority of the jury is more probably correct in its judgment about the statement than the minority. The corollary, which I will call Condorcet's asymptote, shows that under these conditions, as the number of jurors increases, the probability of the majority being correct quickly approaches 1 (Baker 1976, 46-57).

In a longer piece, I will show that while Condorcet acknowledged the difficulties involved in attainment of such conditions, Condorcet's works on philosophy, like *Sketch*, consistently attempted to advocate principles that would bring humanity closer to fulfilling those necessary conditions for human knowledge to grow towards perfection. These principles include improving education and making it universally accessible in order to improve the probability of individuals being able to discern the truth (McLean and Hewitt 1994, 22-33). Insistence that individuals be allowed to come to their own conclusions through their own reason to help ensure statistical independence of judgments (McLean and Hewitt 1994, 37-63). The promotion of ethical principles such as altruism and honesty to help secure sincerity in voting. Universal suffrage (regardless of gender or race) (McLean and Hewitt 1994, 170) and population growth (Condorcet 2004) to help increase the size of voting populations, and thus make them more likely to make correct judgments.

While Condorcet is sometimes discussed as an early discoverer of the intelligence explosion hypothesis; this assertion is made solely on his argument in *Sketch*, which *metaphorically speaking* might be understood as a non-linear regression of a scatterplot showing that knowledge grows at an increasing rate over time. But *Sketch* does not provide a quantitative model for why this faster than linear growth is occurring. But when we look at Condorcet's asymptote, one clearly sees a quantitative model for a technological singularity.¹ By connecting Condorcet's mathematical work (i.e. his 1785 *Essay*) with his philosophical work (i.e. his *Sketch*), we can demonstrate that Condorcet is the earliest known thinker to model in detail an intelligence explosion hypothesis, centuries before others did so.

A few potential criticisms of this claim may be as follows. First, if Condorcet's social choice work is the model for the accelerating increases in human knowledge and

¹ We should not conflate the notion of a technological singularity with a mathematical singularity. While some models of a technological singularity can be expressed with a mathematical singularity, many models do not use a mathematical singularity. Among existing taxonomies of kinds of technological singularities (e.g. (Yudkowsky 2007) and (Sandberg 2010)), Condorcet's asymptote is probably best described as an "intelligence explosion".

lifespans that Condorcet laid out in *Sketch*, then why did he not mention them in *Sketch*? Second, if Condorcet's intelligence explosion is dependent upon population growth to approach perfect knowledge, is not his argument still susceptible to Malthus' counterarguments that there are limits to population growth?

With respect to the first potential criticism, there are a few things to note. First, Condorcet wrote *Sketch* as a non-technical and accessible summary of his ideas. Including math may have made *Sketch* too technical and inaccessible for the readership Condorcet wanted. Second, Condorcet wrote *Sketch* under extreme duress, while he was in hiding from the French Reign of Terror. Under those conditions, he wrote it as a sketch of his ideas, as its title suggests, perhaps with the hope that he could fill in details later. Eventually the Terror caught up with him, and he was captured and sent to prison where he died under mysterious circumstances. With his death, he was never able to fill in the details. Finally, Condorcet does not seem to have resolved a potential problem with his model, a problem which we today call Condorcet's paradox, which is the intransitivity of majority preference.² However, since Condorcet's time, several scholars have shown how this problem can be overcome (Young 1988; List and Goodin 2001; Ben-Yashar and Kraus 2002; Prasad 2012; Brams and Kilgour 2014).

With respect to the second potential criticism, it is possible with mathematical knowledge available today, to show how Condorcet's asymptote can be modified to allow for a finite population of voters, that still asymptotically approaches perfect knowledge. The following model is constructed to be as simple as possible while remaining close to Condorcet's original asymptote. By keeping the model close to Condorcet's, I hope to show how close Condorcet's asymptote was to resolving Condorcet's paradox and the Malthusian criticism of limits to population growth.³

² A simple example of Condorcet's paradox is the case with $3n$ number of voters and three alternatives (e.g. x , y , and z). Suppose the first n voters prefer x over y over z , the second n voters prefer y over z over x , and the third n voters prefer z over x over y . Note that majority prefers x over y , and another majority prefers y over z . If majority preference were transitive, then this would imply that the majority prefers x over z , but when we look at the preferences of the voters, a majority in fact prefers z over x . This was a problem for Condorcet's asymptote that Condorcet recognized because, for example, if for any given voter and any alternatives a_1 and a_2 , a given voter prefers a_1 over a_2 iff she believes a_1 is more probably true than a_2 , then as $n \rightarrow \infty$, Condorcet's asymptote implies that in the case of Condorcet's paradox with $3n$ voters that x is more probably true than y which is more probably true than z which is more probably true than x .

³ It is important to note that Malthus' criticism that there are limits to population growth was not directed at Condorcet's asymptote, which Malthus was not aware of. Because in *Sketch*, Condorcet asserted that lives could go on indefinitely, Malthus asserted that this could not be true due to resource limitations which would cause death and limit population growth. In *Sketch*, Condorcet expresses awareness that humans will have to limit population growth due to resource limitations (Condorcet 2004, 74).

While the model could use many different voting systems to overcome the paradox and the Malthusian criticism, I use approval voting because it seems like the one most similar to Condorcet's work.⁴ In the model, imagine a finite set of voters being presented with a series of statements. After they vote on each statement, they learn and move on to the next statement.

A Model for an Intelligence Explosion

Definitions

Here are the definitions:

Let the set of voters be V , where the n voters are v_1, v_2, \dots, v_n and $n > 1$.

Let the set of m statements be $S: s_1, s_2, \dots, s_m$.

Any given statement s_i is in exactly one of two states: true or false.

For any given statement s_i , each voter has a probability $0 < p_i < 1$ of correctly determining the state of s_i , which is conditional on what voters learned from voting on previous statements.

The probability, that the majority of voters correctly determine the state of s_i , is a_i .

The probability, that exactly half of voters correctly determine the state of s_i , is b_i .

The probability, that the majority of voters incorrectly determine the state of s_i , is c_i .

Sincerity Axiom

In an election on s_i , each voter votes by stating which of the two states she believes s_i to be in.

Independence Axiom

Define independence as follows:

Define $V \setminus v_j$ as the set of all voters in V except for v_j .

For any v_j , how v_j votes on any given s_i is independent of how any subset of $V \setminus v_j$ votes on that s_i .

⁴ Condorcet's last work on voting systems prior to his 1794 death was a brief piece called *On Elections*. Of this work, McLean and Hewitt say "The [*On Elections*] manuscript suggests that Condorcet was moving away from rank-ordering procedures to approval-votes ones..." (McLean and Hewitt 1994, 48).

Learning Axiom

Define the learning axiom as follows:

Let $1 < q_i < 1/p_i$

If the majority of voters correctly determine the state of s_i , then $p_{i+1} = p_i q_i$

If exactly half of all voters correctly determine the state of s_i , then $p_{i+1} = p_i$

If the majority of voters choose the wrong state of s_i , then $p_{i+1} = p_i / q_i$

Let $Q(V, S, p_1)$ specify the values of all possible q_i given V, S , and p_1 . For brevity, we will use Q to refer to $Q(V, S, p_1)$.⁵

Discussion

From Condorcet's jury theorem we know that $b_i = 1 - a_i - c_i$, where a_i and c_i are:

All $k > n/2$

$$a_i = \sum [(n! / [k!(n-k)!]) ([p_i]^k [1-p_i]^{n-k})],$$

All $k < n/2$

$$c_i = \sum [(n! / [k!(n-k)!]) ([p_i]^k [1-p_i]^{n-k})]$$

Now define $\bar{p}_{i+1}(p_i, n, q_i)$ as the expected value of p_{i+1} given p_i, n , and q_i . Note that $\bar{p}_{i+1}(p_i, n, q_i)$ can be expressed as: $\bar{p}_{i+1}(p_i, n, q_i) = (a_i)(p_i q_i) + (b_i)(p_i) + (c_i)(p_i / q_i)$. Furthermore, define $\bar{p}_i(p_1, V, S, Q)$ as the expected value of p_i given p_1, V, S , and Q . Using algebra, the following theorem can be proven.

Theorem: If $(p_1 > 1/2)$ and (for all $1/2 < p_i \leq 2^{-0.5}$, $1 < q_i < 2p_i$) and (for all $2^{-0.5} \leq p_i < 1$, $1 < q_i < 1/p_i$), then as $m \rightarrow \infty$, $\bar{p}_m(p_1, V, S, Q) \rightarrow 1$.

⁵ Informally speaking, note that the path of the society of voters through the statements can be visualized with a ternary tree, where each non-leaf node sprouts exactly three children nodes: the majority is correct, exactly half are correct, and the majority is incorrect on the given statement. So if S has m statements, Q specifies $3^0 + 3^1 + \dots + 3^{m-1}$ possible q_i values. This is because, though Q generates 3^{i+1} possible q_i values for any given a_i , which one of those 3^{i+1} possible q_i values is the one that is actualized is dependent on the path the society of voters takes from the root to the depth of a_i . (Each node at depth $i-1$ is at the depth of a_i)

In other words, If $(p_1 > \frac{1}{2})$ and (for all $\frac{1}{2} < p_i \leq 2^{-0.5}$, $1 < q_i < 2p_i$) and (for all $2^{-0.5} \leq p_i < 1$, $1 < q_i < 1/p_i$), then as voters vote on more and more statements, the expected value of their probability of being correct asymptotically approaches absolute correctness.

Conclusions

Admittedly, this model is crude by contemporary standards. It assumes voters all have the same probability of correctly discerning the truth of statements; it assumes that voters produce their beliefs about the veracity of statements independent of one another; it also assumes that all voters express their beliefs about the statements honestly.⁶ These are highly unlikely in the empirical world, but it emphasizes what kinds of social and political institutions need to be in place for an altruistic intelligence explosion to occur given the assumptions of the model.

First, there have to institutions to encourage altruism in voters, so when they vote, it is directed towards altruistic purposes. Otherwise, they could in theory use the explosion towards malevolent purposes.

Second, institutions must be in place to encourage honesty in voting. Otherwise, the aggregation of votes can fail to produce an intelligence explosion.

Third, institutions must encourage statistically independent votes. This can be partially done through encouragement of independent thinking by voters. When voters simply adopt the opinions of other voters, the autocatalytic explosion fails to take off.

Finally, institutions must encourage knowledge and intelligence growth of voters, for example through education. This is necessary to get intelligence levels high enough for an intelligence explosion to occur.

Of course, with a different model with different assumptions, perhaps more realistic ones, how institutions must be directed might be different. The key thing is that for any given model, we can find out what kinds of institutions are necessary to produce the model's altruistic intelligence explosion. Importantly, if the model's assumptions map well to what is possible in reality, then we have a good idea of what institutions must be directed towards in order for an altruistic intelligence explosion to occur.

References

Baker, Keith Michael. 1976. *Condorcet: Selected Writings*. Indianapolis: Bobbs-Merrill.

Ben-Yashar, Ruth, and Sarit Kraus. 2002. "Optimal collective dichotomous choice under quota constraints." *Economic Theory* 839-852.

Brams, Steven J., and D. Marc Kilgour. 2014. "When Does Approval Voting Make the "Right Choices"?" Edited by Michael A. Jones and Karl-Dieter Crisman. *Proceedings of the AMS Special Sessions on the Mathematics of Decisions, Election, and Games*.

Condorcet. 2004. "Sketch for a Historical Picture of the Progress of the Human Mind: Tenth Epoch." Edited by Keith Michael Baker. *Daedalus* 65-82.

List, Christian, and Robert E. Goodin. 2001. "Epistemic Democracy: Generalizing the Condorcet Jury Theorem." *Journal of Political Philosophy* 277-306.

Malthus, Thomas Robert. 1976. *Essay on the Principle of Population*. Edited by Philip Appleman. New York: W.W. Norton and Company.

McLean, Iain, and Fiona Hewitt. 1994. *Condorcet: Foundations of Social Choice and Political Theory*. Aldershot: Edward Elgar Publishing Limited.

Prasad, Mahendra. 2012. "Condorcet, Preference, and Judgment." *11th Meeting of the Society for Social Choice and Welfare*. New Delhi.

Sandberg, Anders. 2010. "An overview of models of technological singularity." *Third Conference on Artificial General Intelligence*. <http://agi-conf.org/2010/wp-content/uploads/2009/06/agi10singmodels2.pdf>.

Young, Peyton. 1988. "Condorcet's Theory of Voting." *American Political Science Review* 1231-1244.

Yudkowsky, Eliezer S. 2007. *Three Major Singularity Schools*. Machine Intelligence Research Institute, September 30. <https://intelligence.org/2007/09/30/three-major-singularity-schools/>.

⁶ There are works in the literature that take into account differing knowledge levels, non-independence, and insincerity when modelling the jury theorem. For simplicity of the proof of concept, I have not taken these details into account.