

Modeling idea transfer between people

Introduction

Since the 2016 U.S. Presidential election, “fake news” has become a mainstream idea within politics, circles, families, etc., although how exactly our society has achieved this state has been a topic of discussion for quite some time [1]. One of the biggest controversies of that year was the allegation that the Russian government deliberately manufactured fake news, targeting Facebook groups, Twitter, Reddit, etc. to spread fake news. This is a threat with relevance today, mostly because we don’t have any idea of how to prevent it. More generally, the spread and prevalence of fake news, is incredibly hard to prevent.

One thing that may be worth studying is formulating an optimal policy for spreading any arbitrary idea amongst a community of people. But for this we need a **model** for such a community of people, for ideas, and for their transfer. The rest of this proposal will try to formulate pieces of this model, and for formulating some interesting problems for an agent to solve when placed inside this model.

The model

Suppose A and B are two people who have pretty similar ideas. If they communicate, they might reinforce each other’s ideas. Let’s represent ideas with sparse feature vectors ϕ_A and ϕ_B , where each component relates to some idea, and their value ranges between $+1$ and -1 . The **agreement** between A and B , we can model with a dot product.

$$\text{Agreement}(A, B) = \frac{\phi_A \cdot \phi_B}{\|\phi_A\| \|\phi_B\|}. \quad (1)$$

Idea transfer works as follows – if A and B share similar ideas, i.e. $\text{Agreement}(A, B)$ is within a **threshold** ϵ of 1, then we might expect for their ideas to blend, by the following update:

$$\begin{aligned} A &\leftarrow A + \gamma B \\ B &\leftarrow B + \gamma A \end{aligned}$$

where γ is a coefficient of idea spread. However, if A and B don’t share similar ideas, then we might expect for their ideas to drift further apart:

$$\begin{aligned} A &\leftarrow A - \gamma B \\ B &\leftarrow B - \gamma A. \end{aligned}$$

After update, we re-normalize ϕ_A and ϕ_B to their original magnitudes.

In addition to accepting or rejecting ideas, we may suppose that they do neither, and their ideas remain the same. For every person A in the community \mathcal{X} , we randomly assign probabilities of accepting, rejecting, or passing on a new idea, $P_a, a \in \text{Actions} = \{\text{accept, reject, pass}\}$. To be totally unbiased, we randomly assign agreement threshold ϵ_A for each A .

Let’s also assign positions x_a , for every $A \in \mathcal{X}$. Ideas nowadays travel much distance thanks to the internet, but we can try to model the clustering of ideas in small regions of space, perhaps with the following:

$$\text{Prob}_{\text{Interaction}}(A, B) = f(\text{Prob}_{\text{Interaction}}(A), \text{Prob}_{\text{Interaction}}(B), |x_A - x_B|),$$

by increasing the probability of idea transfer as the distance between two people decreases. Such clustering may be able to capture the relative closeness of two people in the same Facebook group, for example.

Some more details, such as whether A and B could partially agree, or whether A or B could move around in physical space, may be interesting to model. For now, let’s consider only the case where A and B do not move in space, but instead are randomly assigned coordinates initially, and hope to capture partial agreement with good selection of some threshold agreement value.

The Problem

Let A be a rational agent in a community of people, \mathcal{X} . We distinguish A from the rest of $X \in \mathcal{X}$, by giving A capability to change its policy $\text{Prob}_{\text{interaction}}$, to achieve some goal. More specifically, we want to perform reinforcement learning with A to find an optimal policy for:

1. Maximizing agreement of the community to A .
2. Keeping A 's ideas alive for as long as possible.

Both of these rewards are “global”, in the sense that they don't represent some state of A , but instead are states of the entire community. If the community is on average, in agreement with A , then A is rewarded. If a member of the community is in agreement with A , then A is also rewarded.

What kind of information will A need to make this informed policy? I think for A to create an optimal policy, it will need lots of information, such as the relative distances $d(X, Y) \forall X, Y \in \mathcal{X}$, and the agreement thresholds ϵ_x for all $X \in \mathcal{X}$.

My hypothesis for the optimal policy is something similar to what we've observed since 2016 – the formation of niche online communities that are quite gullible, ie low threshold agreement, who are fed some idea, and keep it alive by ping-ponging the idea between themselves. The idea then spreads to the outside world through members of the community who are more gregarious.

Challenges / Problems

1. I don't know if the problem in its current form is well-formulated.
2. I don't know whether the current representation of ideas as vector, agreement as dot products, idea spreading as simple updates, is reasonable / good.
3. I am unsure if every member in this community is modeled as an individual agent.
4. I do not know how large the community has to be to see interesting policy development.

References

1. Debord, Guy. 1967 *The Society of the Spectacle*