

Crowd & Prejudice: An Impossibility Theorem for Crowd Labelling without a Gold Standard

Nicolás Della Penna
 ANU NICTA
 Locked bag 8001, 2601
 Canberra, ACT, Australia
 nicolas.della-penna@anu.edu.au

Mark D. Reid
 ANU NICTA
 Locked bag 8001, 2601
 Canberra, ACT, Australia
 mark.reid@anu.edu.au

ABSTRACT

A common use of crowd sourcing is to obtain labels for a dataset. Several algorithms have been proposed to identify uninformative members of the crowd so that their labels can be disregarded and the cost of paying them avoided. One common motivation of these algorithms is to try and do without any initial set of trusted labeled data. We analyse this class of algorithms as mechanisms in a game-theoretic setting to understand the incentives they create for workers. We find an impossibility result that without any ground truth, and when workers have access to commonly shared 'prejudices' upon which they agree but are not informative of true labels, there is always equilibria where all agents report the prejudice. A small amount amount of gold standard data is found to be sufficient to rule out these equilibria.

INTRODUCTION

For "the crowd" is untruth.—Kierkegaard

Precedent literature has proposed a large number of algorithms that take a set of data points labeled by a group of agents, and try to estimate both the reliability of agents. These algorithms can be divided into two sets: those that leverage a small amount of gold standard (ground truth) data (Snow, O'Connor, Jurafsky & Ng 2008, Wauthier & Jordan 2011), and those that do not (Dekel & Shamir 2009b, Raykar, Yu, Zhao, Jerebko, Florin, Hermosillo Valadez, Bogoni & May 2009, Raykar, Yu, Zhao, Valadez, Florin, Bogoni & Moy 2010, Kumar & Lease 2011, Dekel & Shamir 2009a, Yan, Rosales, Fung, Schmidt, Hermosillo, Bogoni, Mouy & Dy 2010). These algorithms that attempt to do without the need for gold standard, do so by using agreement among different labellers as indicative of correctness of a label. This agreement is either at the level of how to label of a given datapoint, as in most cases; or in how features map to labels, as in (Dekel & Shamir 2009b). To achieve this they place their trust on agents who provide labels that are consistent with the labels provided by other agents, or in the case where the same datapoint is not labeled twice, where the proposed feature to label mapping is consistent with other agents' mapping of features to labels. It is often the case that labellers want to be seen as informed by those who are collecting the labels, as the labelling tasks soften pay and it is natural for those collecting the data to avoid the unnecessary cost of paying for labels from uninformed labellers.

We analyse the class of algorithms that do not use gold

standard data as mechanisms in a game-theoretic setting, in order to understand the incentives they create for the agent providing the labels. We first present an impossibility result: that without gold label data, and when workers have access to commonly shared 'prejudices' upon which they agree but which are not informative of true labels, then there is always equilibria where the mechanism does not obtain the true labels from the informed workers, but rather all workers report the prejudice. We then consider how a small amount of gold data is generally sufficient to render situations where the prejudice is reported by informed players as outside the equilibrium set.

One possible criticism of our work is that there is little interest in pointing out that when the assumptions of a statistical model (in this case, that agreement among labellers indicates correctness) do not hold the conclusions drawn from such a model can be misleading. Our argument, however, is more subtle than this: the incentives created by the natural applications of the model in its intended task undermine the very assumptions of the model, by creating incentives for players to agree on the labelling with others, irrespective of whether they believe these to be the true labels.

To make the situation we have in mind more concrete and to clarify how it defers from standard information cascades studied in economics, consider the hypothetical example of a professor who assigns their teaching assistants to grade exams, without grading any themselves. The TAs may or may not know the topic at hand (be informed or uninformed) and they must provide a grade (label) to each exam they are assigned. If the TAs each grade a question on each exam and do so sequentially, so that for each previous answer in a given exam they can observe the grades other TAs have assigned to them, what in economics is referred to as an information cascade can occur. TAs grading later questions can look at the grade a student received for initial exam questions, and guess that the question they were assigned will receive a similar grade, instead of having to understand the answer to the question the student gave and how it relates to the correct answer. In contrast, we study a related but different situation, analogous to one, in which each TA grades (possibly overlapping) full exams, and they do so simultaneously without access to what the others are assigning. Note that if the TAs expect to be rewarded for agreement with others and if they believe others may use some prejudice to grade the exam, such

as assigning higher grades to, say, students with neat hand writing or who use longer words, they might be motivated to also used said prejudice.

An equivalent example can be considered in a crowd sourcing context. Suppose we ask for translations of a given word in language A to speakers of language B; a common prejudice for speakers of language B would be that if a word in language A sounds like the word in language B it must translate to that word. Even if bilingual speakers of A and B are present in the worker pool, if they believe this, the consensus label will be the similar sounding (but possibly incorrect) translation they may choose to report this prejudice as to continue to be employed in the translation task.

RELATED LITERATURE

The study of learning algorithms from a mechanisms design perspective was initiated by (Dekel, Fischer et al. 2008) on a task they term “incentive compatible regression learning”, where agents care about the function that is learned and can report their observations to strategically manipulate it. In contrast to that model, our agents are not motivated to manipulate the learned function mapping examples to labels but are instead are motivated to be seen by the mechanism as capable and thus to continue to be employed as a source of labels for the task. In (Meir, Procaccia & Rosenschein 2010) a strategy-proof mechanisms where agents report labels and their objective is to maximize the accuracy of the learned classifier only on their subset of the data.

Mechanisms designed to elicit subjective probabilities truthfully exploit richer action sets, where the action is not just reporting a label, but also reporting the distribution of labels the population will report. Examples of this type of mechanism is the Bayesian Truth Serum introduced in (Prelec 2004), and the extension of the Peer Prediction Method (Miller & Zeckhauser 2008) proposed by (Witkowski & Parkes 2011).

Our ‘prejudice’ can be thought of as ‘extrinsic random variables’ which allow agents to coordinate their decisions, models for the equilibrium of these have been extensively studied in the economics literature. For a recent review of the literature see (Shell 2008), for experimental evidence for the laboratory see (Duffy & Fisher 2005)

A rich literature on herding behaviour exists in economics, and is closely related to the model we examined but in a sequential instead of simultaneous setting, thus the externality that encourages the herding is of an informational nature instead of in our case where it directly affects payoffs. When agents arrive exogenously ordered sequence and can observe previous agents choices and can follow these, they can either avoiding paying the cost of acquiring information about the payoff of actions themselves or disregarding private information which they may possess, the classic papers in this stream are (Banerjee 1992, Bikhchandani, Hirshleifer &

Welch 1992). Experimental laboratory studies have been carried out looking at herding and information cascades, both in the laboratory (Cipriani & Guarino 2005) and in the internet (Drehmann, Oechssler & Roeder 2005) .

For a recent multidisciplinary review of herding in humans from a cognitive neuroscience perspective see (Raafat, Chater & Frith 2009)

SETTING

Let $\mathcal{Y} = \{1, \dots, K\}$ a set of labels. We consider a game between the *world*, a *mechanism* M , and a set of *agents* \mathcal{A} . Each agent $a \in \mathcal{A}$ falls is of one of two types: *informed* or *uninformed*. We denote the set of informed agents by $\mathcal{A}_I \subseteq \mathcal{A}$. The goal of the mechanism is to identify which agents that are informed. The goal of the agents is to be identified as informed by the mechanism, even if they are not.

Each game is determined by a distribution P over \mathcal{Y}^3 with the random variables $(Y, U, I) \sim P$. Letting $I(\cdot; \cdot)$ denote mutual information, we require that P satisfy three conditions:

$$I(Y; U) = 0 \tag{1}$$

$$I(Y; I) > 0 \tag{2}$$

$$P(Y) \neq P(U) \tag{3}$$

We note that conditions 1 and 2 are equivalent to requiring $P(Y, U, I) = P(U)P(Y)P(I|Y)$ and $P(Y, I) \neq P(Y)P(I)$, respectively. Intuitively, Y is to be interpreted as the “true” label the mechanism is trying to learn, U is some uninformative signal about Y that has a different distribution to Y , and I is an informative signal about Y . It is assumed that $P(Y)$ is common knowledge to the mechanism and all the agents.

The game is played by the world first secretly drawing $y \sim P(Y)$. Every agent $a \in \mathcal{A}$ then receives an i.i.d. draw u_a from $P(U|Y = y) = P(U)$. In addition, informed agents receive a draw i_a from $P(I|Y = y)$. The agents then each decide to report some $y_a \in \mathcal{Y}$ to the mechanism and from this the mechanism must try to determine which agents are informed and which are not.

Strategically, uninformed agents have two choices: they can play a *prejudiced* strategy and report $y_a = u_a$, or they can *randomise* and draw a new y_a from $P(Y)$. Informed agents can also play prejudice or randomise but, in addition, can also play a *truthful* strategy and report $y_a = i_a$. The decision for an agent to be truthful, prejudiced, or random depends on the mechanism. An informed agent may strategically decide to not be truthful in order to maximise its chances of being identified as informed.

For our purposes, a mechanism is a function from a set of reports $R = \{y_a : a \in \mathcal{A}\}$ to set of agents $\mathcal{A}_M \subseteq \mathcal{A}$ that the mechanism identifies as informed and truthful. That is, $M : \mathcal{Y}^{\mathcal{A}} \rightarrow 2^{\mathcal{A}}$. The goal of M is to maximise the probability that \mathcal{A}_M coincides with $\mathcal{A}_{I,T}$ the set of informed agents. It suffices for M to ensure that

$$p_{I,I} = P(a \in \mathcal{A}_M | a \in \mathcal{A}_{I,T}) > \frac{1}{2} \quad (4)$$

$$p_{U,U} = P(a \notin \mathcal{A}_M | a \notin \mathcal{A}_{I,T}) > \frac{1}{2} \quad (5)$$

since repeated independent samples will guarantee that $\left(\frac{p_{I,I}}{1-p_{I,I}}\right)^n \rightarrow \infty$ and $\left(\frac{p_{U,U}}{1-p_{U,U}}\right)^n \rightarrow 0$ as the number of labeled examples n goes to infinity.

RESULTS

We consider the class of mechanisms satisfying the above which succeed when a majority of players are informed and their reports truthful and uninformed players randomise. All proposed algorithms in the literature, to our knowledge, satisfy this elementary criterion.

Equilibria

Three Bayes-Nash equilibria of the game induced by these mechanisms are:

All randomise. When all other agents are randomizing, an agent is indifferent among all labels and thus also about randomising over them. This is not a particularly robust equilibrium as deviation of two agents to prejudice is sufficient to cause all others players to deviate it, and a deviation of two informed agents to truthfulness is also sufficient to cause all other informed agents to improve their payoff by deviating to truthfulness.

Informed are truthful and uninformed randomise.

In this equilibrium the mechanism by and large works in the sense that the set of players that acts consistently in the same manner as the set of informed and truthful agents.

Both play prejudice. When all other plays play prejudice an agent's probability of having their labels found to coincide with others is maximized by playing prejudice, rewardless of whether they are informed or uninformed.

An interesting open question is how to select among these equilibria; while the fragility of the equilibrium where all agents randomize irrespective of type makes it an unlikely candidate it is unclear how to select among the other two. We conjecture that given equal sized populations of informed and uninformed agents the equilibrium selected will depend on the relative entropy of the prejudice and informed distribution, with the lower entropy distribution being more likely to be selected.

Impossibility

Any mechanism in the class defined above must fail in some equilibrium and for some distribution. Consider a situation where all agents are informed and play truthfully: the mechanism succeeds if and only if it identifies all agents as informed and truthful. Now consider a new situation, one in which the equilibrium where all agents play the prejudice occurs, and the distribution of the prejudice in this situation is identical to the distribution of the truth in the previous situation. Since the play observed in both games is identical from the perspective of the mechanism, it must designate all agents as informed and truthful in the second situation and fail, or it must have identified some agents as prejudiced in the first situation and fail.

Using a gold standard

A mechanism that has access to sufficient gold standard data (a sample from the informed distribution) is enough for the impossibility result to no longer apply. In the situation where agents play the prejudice with high probability, the labels reported by those who are playing prejudice will contradict the labels the mechanism has access to and thus can be identified. The mechanism needs access to enough labels from the informative distribution to identify either an agent that is playing prejudice or one that is playing truthfully. It can then extrapolate to agents who have labeled points which it did not originally have labels for based on their agreement or disagreement with those agents it previously identified. This new set of players that has then been identified can be used to extrapolate the the players who labelled points in common with them, and this procedure can be repeated until all players have been identified (this requires that there is sufficient overlap between the data points labeled by agents). The same logic can be adapted to algorithms that do not have players label points in common such as (Dekel & Shamir 2009b) but the overlap then applies to how features map to labels instead of how labels map to points.

Interestingly, access to the prejudice is also sufficient to for the mechanism to succeed in that situation, as this can also be used to identify those who are playing prejudice with high probability when they overlap with the points to which the mechanism has access to the prejudice. The same overlapping procedure can then be used to reveal the strategies of the other players.

Since randomising provides the highest entropy to the sequence of play and thus is the hardest to distinguish from the true labels that a uninformed player can generate the mechanism can identify players who play the prejudice in this situation faster than those that are randomizing. Thus, there is no equilibrium where the mechanism uses the gold data where agents play prejudice. The gold data also guarantees that informed agents have a dominant strategy to be truthful. This implies the only equilibrium has informed agents playing truthfully and uninformed agents randomizing.

CONCLUSION

We consider algorithms that attempt to distinguish between informed and uninformed workers without using gold standard data, and show that when these algorithms are analyzed as mechanisms they can lead to equilibria where no agents truthfully reveal their private information about the label if they have access to it but rather report labels that are uninformative of the true label, but on which they can coordinate with other agents. In future research experimental work to identify whether the equilibria identified in the theoretical model occur, and to test theories of equilibrium selection if they do, would be extremely interesting.

ACKNOWLEDGMENTS

We would like to gratefully acknowledge Maureen Evans for help with editing the paper, and an anonymous reviewer for their helpful detailed comments. The idea that algorithms proposed for learning without gold standard data can be seen as models for herding behaviour come from a blog post by Paul Mineiro. (Mineiro 2011)

CITATIONS AND REFERENCES**REFERENCES**

- Banerjee, A. (1992), ‘A simple model of herd behavior’, *The Quarterly Journal of Economics* **107**(3), 797.
- Bikhchandani, S., Hirshleifer, D. & Welch, I. (1992), ‘A theory of fads, fashion, custom, and cultural change as informational cascades’, *Journal of political Economy* pp. 992–1026.
- Cipriani, M. & Guarino, A. (2005), ‘Herd behavior in a laboratory financial market’, *American Economic Review* pp. 1427–1443.
- Dekel, O., Fischer, F. et al. (2008), Incentive compatible regression learning, in ‘In The ACM-SIAM Symposium on Discrete Algorithms (SODA)’.
- Dekel, O. & Shamir, O. (2009a), Good learners for evil teachers, in ‘Proceedings of the 26th Annual International Conference on Machine Learning (ICML)’, ACM, New York, NY, USA, pp. 233–240.
- Dekel, O. & Shamir, O. (2009b), Vox populi: Collecting high-quality labels from a crowd, in ‘In Proceedings of the 22nd Annual Conference on Learning Theory (COLT)’.
- Drehmann, M., Oechssler, J. & Roeder, A. (2005), ‘Herding and contrarian behavior in financial markets: An internet experiment’, *The American Economic Review* **95**(5), 1403–1426.
- Duffy, J. & Fisher, E. (2005), ‘Sunspots in the laboratory’, *American Economic Review* pp. 510–529.
- Kumar, A. & Lease, M. (2011), Modeling annotator accuracies for supervised learning, in ‘Workshop on Crowdsourcing for Search and Data Mining (WSDM)’, pp. 19–22.
- Meir, R., Procaccia, A. & Rosenschein, J. (2010), On the limits of dictatorial classification, in ‘Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems’, Vol. 1, pp. 609–616.
- Miller, N., P. R. & Zeckhauser, R. (2008), ‘Eliciting informative feedback: The peer-prediction method’, *Management Science* **51**(9), 1359–1373.
- Mineiro, P. (2011), ‘Formalized herd mentality’.
URL:
<http://www.machinedlearnings.com/2011/10/formalized-herd-mentality.html>
- Prelec, D. (2004), ‘A Bayesian truth serum for subjective data’, *Science* **306**(5695), 462.
- Raafat, R., Chater, N. & Frith, C. (2009), ‘Herding in humans’, *Trends in cognitive sciences* **13**(10), 420–428.
- Raykar, V., Yu, S., Zhao, L., Jerebko, A., Florin, C., Hermosillo Valadez, G., Bogoni, L. & May, L. (2009), Supervised learning from multiple experts: whom to trust when everyone lies a bit, in ‘Proceedings of the 26th Annual International Conference on Machine Learning (ICML)’, pp. 889–896.
- Raykar, V., Yu, S., Zhao, L., Valadez, G., Florin, C., Bogoni, L. & Moy, L. (2010), ‘Learning from crowds’, *The Journal of Machine Learning Research* **99**, 1297–1322.
- Shell, K. (2008), Sunspot equilibrium, in ‘The New Palgrave: A Dictionary of Economics’, 2nd edn, Palgrave Macmillan.
- Snow, R., O’Connor, B., Jurafsky, D. & Ng, A. (2008), Cheap and fast—but is it good?: evaluating non-expert annotations for natural language tasks, in ‘Proceedings of the Conference on Empirical Methods in Natural Language Processing’, Association for Computational Linguistics, pp. 254–263.
- Wauthier, F. & Jordan, M. (2011), Bayesian bias mitigation for crowdsourcing, in ‘Proceedings of Advances in Neural Information Processing Systems (NIPS)’.
- Witkowski, J. & Parkes, D. C. (2011), Peer prediction with private beliefs, in ‘ACM EC Workshop on Social Computing’, San Jose, CA.
- Yan, Y., Rosales, R., Fung, G., Schmidt, M., Hermosillo, G., Bogoni, L., Mouy, L. & Dy, J. (2010), Modeling annotator expertise: learning when everybody knows a bit of something, in ‘Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS)’, pp. 932–993.