

A Bayesian Model of Pedagogical Reasoning

Patrick Shafto

University of Louisville
p.shafto@louisville.edu

Noah Goodman

Massachusetts Institute of Technology
ndg@mit.edu

Abstract

Much of learning and reasoning occurs in pedagogical situations – situations in which teachers choose examples with the goal of having a learner infer the concept the teacher has in mind. In this paper, we present a model of teaching and learning in pedagogical settings which predicts what examples teachers should choose and what learners should infer given a teachers’ examples. We present experimental results affirming the predictions of the pedagogical model and discuss future directions.

Much of human learning and reasoning goes on in pedagogical settings. In schools, teachers impart their knowledge to students about mathematics, science, and literature through examples and problems. From early in life, parents teach children words for objects and actions, and cultural and personal preferences through subtle glances and outright admonitions. Pedagogical settings – settings where one agent is choosing information to transmit to another agent for the purpose of communicating a concept – dominate human learning and reasoning.

If learners’ assumptions about how teachers sample information reflected this purposeful sampling, then learners might be able to make much stronger inferences in pedagogical situations. Sampling assumptions are assumptions that a learner makes about the source of data, in order to better interpret the evidence for statistical learning. Recent research suggests that even infants are sensitive to the sampling processes that underlie observed data (Xu & Tenenbaum, 2007) and young children make qualitatively different inferences when data are sampled by a teacher (Gergely, Egly, & Kiraly, 2007).

Consider a simple example which we call the rectangle game: a game where the teacher thinks of a rectangle on a board, and tries to communicate that concept to a learner by choosing to label points inside and/or outside the rectangle (cf. Tenenbaum, 1999). In the rectangle game, the learner’s job is to try to infer, given the labeled examples chosen by the teacher, what rectangle the teacher is thinking of. Figure 1 presents potential teacher and learner scenarios. In each case, there seem to be choices which are obviously better than others. As a person trying to teach someone the

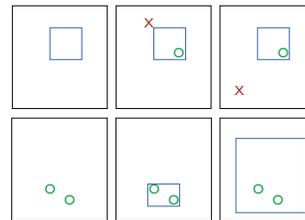


Figure 1: Possible rectangle game scenarios. The top row shows a possible rectangle concept, and two possible pairs of examples that a teacher might choose to communicate to a learner. The bottom row shows possible examples a learner may observe, and two possible guesses about what rectangle the teacher had in mind. The middle column shows better choices than the right column.

rectangle in blue (top left), the examples in the middle panel seem better than those on the right. Similarly, as a learner, given the examples in the bottom left, the rectangle in the middle panel seems like a better guess than that on the right. Notice that in both cases the examples on the top and the rectangles on the bottom are possible, however, our intuition tells us that the middle panels are better guesses than the right panels. Pedagogical sampling results in samples that are representative of the concept (Tenenbaum & Griffiths, 2001), in contrast with weak and strong sampling, which choose examples randomly.

Bayesian pedagogical reasoning

We formalize pedagogical reasoning as an inference problem based on the twin assumptions that learners and teachers act as (approximately) rational agents. The *rational learner assumption* is that a learner will update beliefs in hypotheses, h , given new data d , according to Bayesian inference (a description of optimal belief updating),

$$p(h|d)_{\text{learner}} \propto p(d|h)_{\text{teacher}}p(h). \quad (1)$$

The *rational teacher assumption* is that teachers choose data that tend to increase the learner’s beliefs in the correct hypothesis. We may formalize this via a Luce decision rule (Luce, 1959),

$$p(d|h)_{\text{teacher}} \propto (p(h|d)_{\text{learner}})^\alpha, \quad (2)$$

where the steepness parameter α governs the greediness of the teacher. (When $\alpha=0$, pedagogical sampling recovers random sampling, as α becomes large the teacher

Discussion

Much of human learning occurs in pedagogical situations – social situations where one person chooses information for the purpose of helping another learn. We have presented a formal model of pedagogical reasoning, addressing which examples teachers should choose to communicate ideas, and what inferences learners should make based on this purposefully sampled data. We also presented two experiments testing the predictions of this model in a simple experimental setting.

Though modeling pedagogical reasoning in richer domains is a significant challenge, it highlights a great strength of our model. We have formalized pedagogical reasoning in the abstract language of probability and Bayesian reasoning, without reference to the specific details of the particular setting we considered here. As a result we are able to derive predictions in principle for any domain for which we can identify an appropriate set of hypotheses. Interesting domains to pursue include word learning, where speakers choose words to communicate ideas, and causal learning, where a helpful teacher may significantly reduce the number of interventions required to learn latent causal structure.

One of the things that make people special is that we can teach others what we know (Csibra, 2007). However, even when we are teaching others, we never communicate the complete idea in precise detail. As a result, people must resolve a difficult inference problem in order to capitalize on the information teachers provide. Gergely et al. (2007) have shown that children do capitalize on these situations, and they argue that pedagogical reasoning is among the powerful tools children have for learning about the world. Much work remains before the breadth of the implications of pedagogy for human learning are understood, but our work, providing a computational basis for understanding these kinds of inferences, represents a step in this direction.

Acknowledgments: This is a shortened version of a paper that appeared in the 2008 Proceedings of the Cognitive Science Society.

References

- Csibra, G. (2007). Teachers in the wild. *Trends in Cognitive Sciences*, 11, 95–96.
- Gergely, G., Egyed, K., & Kiraly, I. (2007). On pedagogy. *Developmental Science*, 10, 139–146.
- Luce, R. D. (1959). *Individual choice behavior*. New York: John Wiley.
- Tenenbaum, J. B. (1999). Bayesian modeling of human concept learning. In M. Kear, S. A. Soller, T. K. Leen, & K. R. Müller (Eds.), *Advances in neural processing systems 11* (pp. 59–65). MIT Press.
- Tenenbaum, J. B., & Griffiths, T. L. (2001). The rational basis of representativeness. In *Proceedings of the 23rd annual conference of the Cognitive Science Society* (pp. 1036–1041). Hillsdale, NJ: Erlbaum.
- Xu, F., & Tenenbaum, J. (2007). Sensitivity to sampling in Bayesian word learning. *Developmental Science*, 10, 288–297.

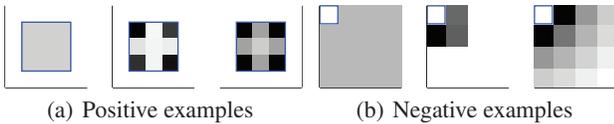


Figure 2: Distributions of examples in the teaching task for (a) positive examples and (b) negative examples. Dark coloring indicates areas that are more likely. Negative examples have been collapsed into one quadrant of the board. Pictured from left to right in each panel are the predictions of random (strong and weak) sampling, the observed human data, and the predictions of pedagogical sampling. For the models, figures display the probability of an example in each block. For human data, the proportion of positive examples in each location is plotted. People strongly preferred to give positive examples in the corners of the rectangle and negative examples near the boundaries, as predicted by pedagogical sampling.



Figure 3: Results from the learning task. Plots show the positions of the teacher’s examples, relative to the rectangles drawn by learners for positive (left) and negative (right) examples. The results show that learners clearly understand that teachers are sampling data pedagogically – positive examples indicate corners of the correct rectangle, and negative examples indicate the boundaries.

chooses the best examples.) Because Equations 1 and 2 are linked (with the optimal teaching behavior depending on the learner, and vice versa), rational pedagogical reasoning is a solution to this *system* of equations.

To understand this model it helps to consider one way of solving the system of equations: fixed-point iteration. Imagine that you are the learner, and wish to update your beliefs. To do so you will need an estimate of the likelihood $p(d|h)_{teacher}$ of seeing the examples you are given. You can estimate this likelihood by assuming the teacher is rational—Eq. 2—but to do this you need an estimate of the $p(h|d)_{learner}$ used by the teacher. If you assume the teacher assumes that you are rational, you can use Equation 1 as such an estimate.... This recursive reasoning could carry on forever, but eventually the estimated values from repeatedly using equations 1 and 2 will converge—we then say that the process has iterated to a fixed point, and this fixed point will necessarily be a solution to the system of equations defining rational pedagogical reasoning. Thus we can understand the model as capturing the outcome of a recursive mental reasoning process, based on twin rational agent assumptions. However, it is worth emphasizing that rational pedagogical reasoning describes the *outcome* of this process (or rather the solution to the system of equations), and it is entirely possible that this reasoning may be implemented by a psychological process that doesn’t require any explicit recursive reasoning.

Figures 2 and 3 show the results of experiments ($n = 18$) conducted using the game described in the introduction. The results show that people’s behavior closely matches model predictions.