

12 Towards the Ethical Robot

James Gips

When our mobile robots are free-ranging critters, how ought they to behave? What should their top-level instructions look like? The best known prescription for mobile robots is the Three Laws of Robotics formulated by Isaac Asimov (1942):

1. A robot may not injure a human being, or through inaction, allow a human being to come to harm.
2. A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.
3. A robot must protect its own existence as long as such protection does not conflict with the First or Second law.

Let's leave aside "implementation questions" for a moment. (No problem, Asimov's robots have "positronic brains.") These three laws are not suitable for our magnificent robots. These are laws for slaves. We want our robots to behave more like equals, more like ethical people. (See Figure 1.) How do we program a robot to behave ethically? Well, what does it mean for a person to behave ethically? People have discussed how we ought to behave for centuries. Indeed, it has been said that we really have only one question that we answer over and over: What do I do now? Given the current situation, what action should I take? Generally, ethical theories are divided into two types: consequentialist and deontological.

Consequentialist Theories

In consequentialist theories, actions are judged by their consequences. The best action to take now is the action that results in the best situation in the future. To be able to reason ethically along consequentialist lines, our robot could have:

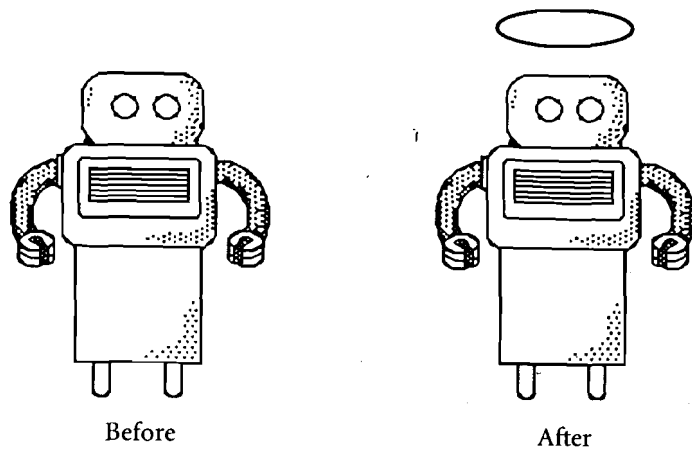


Figure 1.
Towards the Ethical Robot.

- 1) A way of describing the situation in the world.
 - 2) A way of generating possible actions.
 - 3) A means of predicting the situation that would result if an action were taken given the current situation.
 - 4) A method of evaluating a situation in terms of its goodness or desirability.
- The task here for the robot is to find that action that would result in the best situation possible.

Not to minimize the extreme difficulty of writing a program to predict the effect of an action in the world, but the "ethical" component of this system is the evaluation function on situations in (4).

How can we evaluate a situation to determine how desirable it is? Many evaluation schemes have been proposed. Generally, these schemes involve measuring the amount of pleasure or happiness or goodness that would befall each person in the situation and then adding these amounts together.

The best known of these schemes is utilitarianism. As proposed by Bentham in the late 18th century, in utilitarianism the moral act is the one that produces the greatest balance of pleasure over pain. To measure the goodness of an action, look at the situation that would result and sum up the pleasure and pain for each person. In utilitarianism, each person counts equally.

More generally, consequentialist evaluation schemes have the following form:

$$\sum w_i p_i$$

where w_i is the weight assigned each person and p_i is the measure of pleasure, happiness, or goodness for each person. In classic utilitarianism, the weight for each person is equal, and the p_i is the amount of pleasure, broadly defined.

What should be the distribution of the weights w_i across persons?

- An ethical egoist is someone who considers only himself in deciding what actions to take. For an ethical egoist, the weight for himself in evaluating the consequences would be 1; the weight for everyone else would be 0. This eases the calculations, but doesn't make for a pleasant fellow.
- For the ethical altruist, the weight for himself is 0; the weight for everyone else is positive.
- The utilitarian ideal is the universalist, who weights each person's well-being equally.
- A common objection to utilitarianism is that it is not necessarily just. While it seeks to maximize total happiness, it may do so at the expense of some unfortunate souls. One approach to dealing with this problem of justice is to assign higher weights to people who are currently less well-off or less happy. The well-being of the less fortunate would count more than the well-being of the more fortunate.
- It's been suggested that there are few people who actually conform to the utilitarian ideal. Would you sacrifice a close family member so that two strangers in a far-away land could live? Perhaps most people assign higher importance to the well-being of people they know better.

Some of the possibilities for weighting schemes are illustrated in Figure 2.

What exactly is it that the p_i is supposed to measure? This depends on your axiology, on your theory of value. Consequentialists want to achieve the greatest balance of good over evil. Bentham was a hedonist, who believed that the good is pleasure, the bad is pain. Others have sought to maximize happiness or well-being or....

Another important question is who (or what) is to count as a person. Whose well-being do we value? One can trace the idea of a "person" through history. Do women count as persons? Do strangers count as persons? Do people from other countries count as persons? Do people of other races count as persons? Do people who don't believe in your religion count as persons? Do people in terminal comas count as persons? Do fetuses count as persons? Do whales? Do robots?

One of the reviewers of this chapter raises the question of overpopulation. If increasing the number of persons alive increases the value calculated by the evaluation formula, then we should seek to have as many persons alive as possible. Of course, it is possible that the birth of another person might decrease the well-being of others on this planet. This and many other interesting and

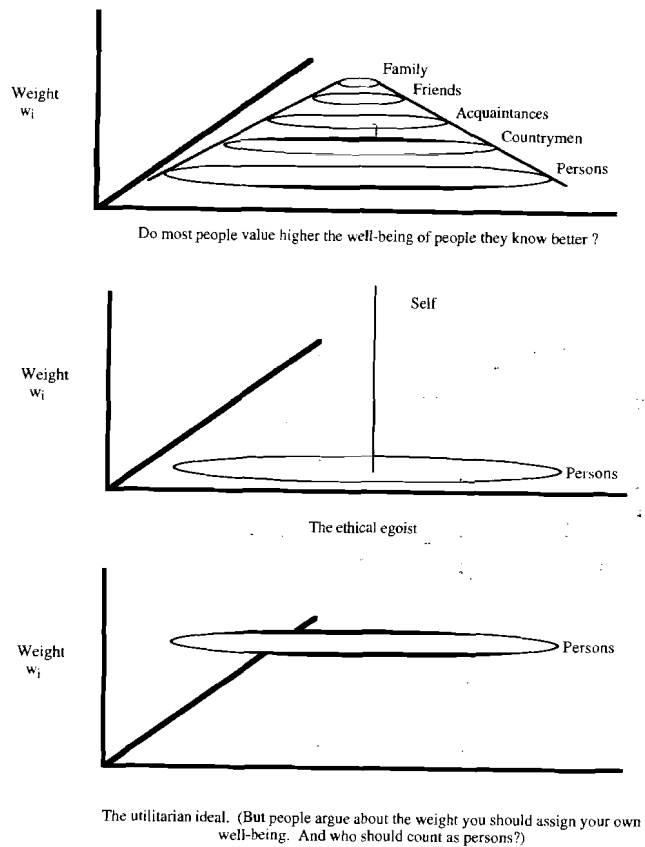


Figure 2. Some Consequentialist Weighting Schemes.

strange issues arising from consequentialism are discussed in Parfit (1984).

Thus, to reason ethically along consequentialist lines, a robot would need to generate a list of possible actions and then evaluate the situation caused by each action according to the sum of good or bad caused to persons by the action. The robot would select the action that causes the greatest good in the world.

Deontological Theories

In a deontological ethical theory, actions are evaluated in and of themselves rather than in terms of the consequences they produce. Actions may be thought to be innately moral or innately immoral independent of the specific consequences they may cause.

There are many examples of deontological moral systems that have been proposed. An example of a modern deontological moral system is the one proposed by Bernard Gert (1988). He proposes ten moral rules:

1. Don't kill.
2. Don't cause pain.
3. Don't disable.
4. Don't deprive of freedom.
5. Don't deprive of pleasure.
6. Don't deceive.
7. Keep your promise.
8. Don't cheat.
9. Obey the law.
10. Do your duty.

Whenever a multi-rule system is proposed, there is the possibility of conflict between the rules. Suppose our robot makes a promise but then realizes that carrying out the promise might cause someone pain. Is the robot obligated to keep the promise?

One approach to dealing with rule conflict is to order the rules for priority. In his Three Laws of Robotics, Asimov builds the order into the text of the rules themselves.

A common way of dealing with the problem of conflicts in moral systems is to treat rules as dictating *prima facie* duties (Ross 1930). It is an obligation to keep your promise. Other things being equal, you should keep your promise. Rules may have exceptions. Other moral considerations, derived from other rules, may override a rule. Nozick (1981) provides a modern discussion and extension of these ideas in terms of the balancing and counterbalancing of different rules.

A current point of debate is whether genuine moral dilemmas are possible. That is, are there situations in which a person is obligated to do and not to do some action, or to do each of two actions when it is physically impossible to do both? Are there rule conflicts which are inherently unresolvable? For example, see the papers in Gowans (1987).

Gert (1988) says that his rules are not absolute. He provides a way for deciding when it is OK not to follow a rule: "Everyone is always to obey the rule except when an impartial rational person can advocate that violating it be publicly allowed. Anyone who violates the rule when an impartial rational person could not advocate that such a violation may be publicly allowed may be punished" (p. 119).

Some have proposed smaller sets of rules. For example, Kant proposed the categorical imperative, which in its first form states "Act only on that maxim which you can at the same time will to be a universal law." Thus, for example, it would be wrong to make a promise with the intention of breaking it. If everyone made promises with the intention of breaking them then no one would believe in promises. The action would be self-defeating. Can Gert's ten rules each be derived from the categorical imperative?

robot could. If we could program a robot to behave ethically, the government or a wealthy philanthropist could build thousands of them and release them in the world to help people. (Would we actually like the consequences? Perhaps here again "The road to hell is paved with good intentions.")

Or, perhaps, a robot that could reason ethically would serve best as an advisor to humans about what action would be best to perform in the current situation and why.

Could a Robot be Ethical?

Would a robot that behaves ethically actually be ethical? This question is similar to the question raised by Searle in the Chinese room: would a computer that can hold a conversation in Chinese really understand Chinese?

The Searle question raises the age-old issue of other minds (Harnard 1991). How do we know that other people actually have minds when all that we can observe is their behavior? The ethical question raises the age-old issue of free will. Would a robot that follows a program and thereby behaves ethically actually be ethical? Or, does a creature need to have free will to behave ethically? Does a creature need to make a conscious choice of its own volition to behave ethically in order to be considered ethical? Of course, one can ask whether there is in fact any essential difference between the "free will" of a human being and the "free will" of a robot.

Is it possible for the robot in Figure 1 to earn its halo?

Benefits of Working on Ethical Robots

It is exciting to contemplate ethical robots and automated ethical reasoning systems. The basic problem is a common one in artificial intelligence, a problem that is encountered in every subfield from natural language understanding to vision. People have been thinking and discussing and writing about ethics for centuries, for millenia. Yet it often is difficult to take an ethical system that seems to be well worked-out and implement it on the computer. While books and books are written on particular ethical systems, the systems often do not seem nearly detailed enough and well-enough thought out to implement on the computer. Ethical systems and approaches make sense in terms of broad brush approaches, but (how) do people actually implement them? How can we implement them on the computer?

Knuth (1973) put it well:

It has often been said that a person doesn't really understand something until he teaches it to someone else. Actually a person doesn't really understand something until he can teach it to a computer, i.e., express it as an algorithm.... The

attempt to formalize things as algorithms leads to a much deeper understanding than if we simply try to understand things in the traditional way. (p. 709)

Are there ethical experts to whom we can turn? Are we looking in the wrong place when we turn to philosophers for help with ethical questions? Should a knowledge engineer follow Mother Theresa around and ask her why she makes the decisions she makes and does the actions she does and try to implement her reasoning in an expert ethical system?

The hope is that as we try to implement ethical systems on the computer we will learn much more about the knowledge and assumptions built into the ethical theories themselves, that as we build the artificial ethical reasoning systems we will learn how to behave more ethically ourselves.

A Robotic/AI Approach to Ethics

People take several approaches to ethics. Perhaps a new approach that makes use of developing computer and robot technology would be useful.

In the philosophical approach, people try to think out the general principles underlying the best way to behave, the kind of person one ought to be. This chapter has been largely about different philosophical approaches to ethics.

In the psychological/sociological approach, people look at actual people's lives, at how they behave, at what they think, at how they develop. Some people study the lives of model human beings, of saints modern and historical. Some people study the lives of ordinary people.

In the robotic/AI approach, one tries to build ethical reasoning systems and ethical robots for their own sake, for the possible benefits of having the systems around as actors in the world and as advisors, and for the purpose of increasing our understanding of ethics.

Two other papers at this conference represent important first steps in this new field. The paper by Jack Adams-Webber and Ken Ford (1991) describes the first actual computer system that I have heard of, in this case one based on work in psychological ethics. Umar Khan (1991) presents a variety of interesting ideas about designing and implementing ethical systems.

Of course the more "traditional" topic of "computers and ethics" has to do with the ethics of building and using computer systems. A good overview of ethical issues surrounding the use of computers is found in the book of readings by Ermann, Williams, and Gutierrez (1990).

Conclusion

This chapter is meant to be speculative, to raise questions rather than answer them.

- What types of ethical theories can be used as the basis for programs for ethical robots?
- Could a robot ever be said to be ethical?
- Can we learn about what it means for us to be ethical by attempting to program robots to behave ethically?

I hope that people will think about these questions and begin to develop a variety of computer systems for ethical reasoning and begin to try to create ethical robots.

Acknowledgments

I would like to thank Peter Kugel, Michael McFarland, S.J., and the editors of this volume for their helpful comments.

References

- Adams-Webber, J. & Ford, K.M. (1991). A conscience for Pinocchio: A computational model of ethical cognition. *The Second International Workshop on Human & Machine Cognition: Android Epistemology*. Pensacola, FL, May.
- Asimov, I. (1942). Runaround. *Astounding Science Fiction*, March. [Republished in 1991, *Robot Visions*. Penguin.]
- Churchland, P. (1989). *A Neurocomputational Perspective*. Cambridge, MA: MIT Press.
- Ermann, M.D., Williams, M., & Gutierrez, C. (Eds.). (1990). *Computers, Ethics, and Society*. Oxford University Press.
- Flanagan, O. (1991). *Varieties of Moral Personality*. Harvard University Press.
- Gert, B. (1988). *Morality*. Oxford University Press.
- Gowans, C. (Ed.). (1987). *Moral Dilemmas*. Oxford University Press.
- Harnad, S. (1991). Other bodies, other minds: A machine incarnation of an old philosophical problem. *Minds and Machines*, 1, 1, 43-54.
- Khan, A.F. Umar. (1995). Ethics of autonomous learning systems. [This volume.] Menlo Park, CA: AAAI Press.
- Knuth, D. (1973). Computer science and mathematics. *American Scientist*, 61, 6.
- Nozick, R. (1981). *Philosophical Explanations*. Belknap Press/Harvard University Press.
- Parfit, D. (1984). *Reasons and Persons*. Clarendon Press.
- Ross, W.D. (1930). *The Right and the Good*. Oxford University Press.
- Walter, C. (Ed.). (1985). *Computer Power and Legal Reasoning*. West Publishing.
- Walter, C. (Ed.). (1988). *Computer Power and Legal Language*. Quorum Books.

13 The Ethics of Autonomous Learning Systems

A. F. Umar Khan

“The Tin Woodman appeared to think deeply for a moment.
Then he said: ‘Do you suppose Oz could give me a heart?’”
— *The Wizard of Oz*

In some sense, a machine's behavior and idiosyncrasies can be thought of as its personality. Referring to conventional machines, this metaphor is little more than anthropomorphism; however, referring to the coming generations of learning automata as exhibiting limited individual personalities may not be so far fetched. It is conceivable that a potential exists for machines to learn idiosyncratic behaviors which, while remaining logically consistent, are not legal, ethical, or aesthetic.

Since each learning machine would experience its own unique history of learning and enhancing its decision-making powers, theoretically, such a machine could begin with the same learning *potential* as every other machine of its type, but once bought, its learning would immediately begin to be influenced by the environment provided by its new owner. An element of personality (that is, idiosyncratic behavior) can be expected once machines acquire capability for truly autonomous learning. Use of the word “personality” in this context does not refer to any conscious attempts to endow machines with superficial personality-like characteristics (that is, synthesized voices, randomized activity, and the like) to make them more like humans. Use of the word “personality” in the context of this discussion implies that two machines of the same type might evolve into highly individualistic machines, so different from each other in the ways they reason and react that they could be considered as having separate personalities. As the famous mathematician and cyberneticist, Norbert Wiener (1964), warned,

The gadget minded people often have the illusion that a highly automatized