

| | |
|--|---|
| Is it in our self-interest to being good word count: 1834] | Is it in our self-interest to being good word count: 1953] |
| #20387090 | #20387090 |
| July 27, 2015 | June 29, 2015 |
| It is in our self-interest to be good through two parts. First, being good rewards. Second, being good benefits individually. | We explore this problem from the perspective of a rational agent would act, and their problem so far. If we follow up, the diligent reader would find many different methods to solving a problem. Some methods would involve being good, and some would involve being bad. When it is possible to achieve the same goal through unjust means? In response to this, we can see an example of an unjust man versus a just man. It is easier to be good. |
| The main arguments are drawn from two textbooks, The Fundamentals of Ethics by Shafer-Landau, R. (2014). The Fundamentals of Ethics, 4th Edition, University Press of Kansas, and the textbook used in CS 486 in the University of Illinois at Urbana-Champaign, called Artificial Intelligence: A Modern Approach, Russell, S.J. & Norvig, P. (2010). Artificial Intelligence: A Modern Approach. Prentice Hall (referred as AIMA). There are similarities in the fields, different approaches, but same goals. | The main arguments are drawn from the textbook, The Fundamentals of Ethics by Shafer-Landau, R. (2014). The Fundamentals of Ethics, 4th Edition, University Press of Kansas, and the textbook used in CS 486 in the University of Illinois at Urbana-Champaign, called Artificial Intelligence: A Modern Approach, Russell, S.J. & Norvig, P. (2010). Artificial Intelligence: A Modern Approach. Prentice Hall (referred as AIMA). There are similarities in the fields, different approaches, but same goals. |
| As in the case with mathematical proofs, we need to define our terms. Definitions to the thesis. | Definition of our terms. |
| Definition of our terms | We define the term "our" in the question "Is it in our self-interest to be good?" to mean a ratio of good to bad. This is a fundamental assumption we make, in that intelligence is plagued with emotional and irrational decisions. Hence it is in our self-interest to act upon emotional decisions, and not rational. We also define good later on in the framework. |
| The term "our" refers to a rational agent that always acts to achieve the best outcome under uncertainty, the best expected outcome. Definition from page 4 of AIMA [1]. This simplifying assumption is made to focus on rational decisions. | |
| The philosophy text defines rationality (theory) slightly differently. Rationality is based on one's preference orderings. | From a Rational Agent's Perspective |
| Both are similar, but I highlight that the philosophy text includes acting under uncertainty. The philosophy text is more focused on the rational agent's perspective. | First, let's look at the reasons why a rational agent's self-interest to be good, in the future, and secondly, why it benefits from being good. |
| Definition of good | Definition of good. |
| Being good means having insight and understanding. Consequently, as I will explore, being good is always trying to maximize your own utility (which is always trying to maximize your own utility). | We're now ready to define what good means. "Our self-interest to be good": being good always maximizes the reward for your utility. We'll see that this is not necessarily picking the most seemingly good moment. |
| Part I: Long term rewards | Firstly, why it would be in a rational agent's self-interest to be good, due to future rewards. |
| The key insight into gaining more long-term utility is choosing the best action (being selfish) to maximize the best rewards. Some actions do not provide the best rewards, but they are beneficial in the long run, but not in the short run. | From our definition of what it means to be good, we shouldn't be bad (not good), which means picking the action that maximizes the reward. In other words, being bad is being always greedy. We'll see that this is not necessarily picking the most seemingly good moment. |
| Philosophy Textbook | |
| The first argument is from the course text, The Fundamentals of Ethics, which states that action isn't to be greedy, but to be rational. | |
| The previous sections talk about different theories, and this quotation below concludes about the differences from those different social theories [1]. I won't go in detail about what each theory is, but merely look at the commonalities: | page 215 of The Fundamentals of Ethics [1]: |
| The key to understanding [the social contract theory] is the idea that contractors are, above all, self-interested. | The key to understanding them [social contract theory], however, lies in the idea that the contractors are rational and self-interested. |
| Being self-interested is not the same as being selfish. Being self-interested is having a strategy to maximize your utility. | Being self-interested is not the same as being selfish. Being self-interested is having a strategy to maximize your utility. |

you are faring in life. Being selfish importance on your own well-being rel others.

While being selfish maximize your own theories have similar commonalities i isn't being selfish but rather being saying the best action isn't the same

This supports the idea that there are good, at a cost of gaining less benef selfish maximizes the present-moment future rewards.

Artificial Intelligence Textbook

The second argument is from the AIMA reinforcement learning, supporting th isn't always best for yourself.

I define some terms for readers in th learning agent has a fixed policy tha whereas an active agent must decide w utility function (also known as rewar the agent's performance. A state encl action. When actions and the search s utility function can be learned by va

The experiment setup is to finding an grid. But the grid is non-determinist take may or may not result in the des choses to move left, but the move hap

Take it without proof that in this se be extracted by one-step look-ahead t utility. After the 276^{th} value iteration and learning the optimal policy at each step, the poli sticks to using that policy, never le other states. We call this agent the chooses is optimal for each state (pl interestingly, the agent does not lea true optimal policy.

Repeated experiments show that the gr converges to the optimal policy[footn pg. 839 AIMA 3rd ed. j!]

How can it be that choosing the optim suboptimal results? The answer is tha the same as the true environment. An between exploitation (to maximize its maximize its long-term well-being). P getting stuck in a rut. Pure explorat knowledge is of no use if one never p practice.

The greedy agent's selfish strategy t rewarding action doesn't imply the be out on other opportunities. The optim of benefits from the future, through

Part I: Conclusion

The two textbooks shares the conclusi isn't necessarily one that maximizes one that factors in the future long t

Part II: Collective benefits

The main idea in this part is it is i because being good benefits the colle Rewards to a group of rational agents

I explore the readings from Collectiv the problem as a prisoner's dilemma p case of Nash equilibrium. Then applyi learned about the prisoner's dilemma.

Collective Action Problem

you are faring in life. Being selfish importance on your own well-being rel others.

What it means is that there are futur a cost of gaining less benefit at-the be self-interested instead of being s maximizes the reward we get now but a is also seen in the AI textbooks, thr function[footnote:

Also known as an utility function.], measuring how much reward was atta self-interest is to choose the action reward, i.e. the optimal action to ta that maximizes the benefit now, and i [r]epeated experiments show that the converges to the optimal policy"[foot

pg. 839 AIMA 3rd ed.]. Essentially, the AI textbooks also the philosophy textbook – the idea th and self-interested but not selfish.

Secondly, why it rewards rational a being good.

In other words, we'll look at how rew agents cooperate with each other and understand that when others have good have the trust of others, this opens choice of options to choose from. off collective group than being bad (gree happens when everyone is good, when e when everyone is bad except you, and bad. To the diligent reader, this ind Dilemma. And indeed we will look at N this, we'll look at why behaving badl form of a collective action problem.

Collective Action Problem

In the Collective Action Problem [footnote: p. 366 Constellations Volume 7, Number 3, 2015].
<http://homes.chass.utoronto.ca/~jheat>
 3, 2015
], people behave irrationally because but at a cost to others.

The reading suggests that we are lazy to be the one putting in the effort that is being good. And if everyone is being good, it is easier to not be good and reap the catch-22 situation.

This problem can be casted to a prisoner's well known case of Nash equilibrium, learnt from prisoner's dilemma. The case

1. Everyone being lazy and you are lazy both defect.
2. Everyone being good and you are lazy friend cooperates and you defect.
3. Everyone being good and you are good both cooperate.
4. Everyone being lazy and you are good friend defects and you cooperate.

Now I apply results learnt from analysis of dilemma.

Prisoner's Dilemma

A key assumption made is the prisoner reward or punish their partner and thus affect their reputation in the future.

This assumption is not realistic in that punishment is death. Similarly, this is the collective action problem that you are cooperating. Realistically there would be from the collective good on your behaviour, thus it benefits you to cooperate.

But what if the assumption holds?

Since betraying a partner offers a greater reward than cooperating with him, all purely rational agents would betray the other, and so the only possible outcome for them to betray each other. The time element is added to the classic iterated prisoner's dilemma, then cooperation is a rational outcome, as well explained in a version of prisoner's dilemma, in the chapter [footnote: p. 118 Multi-agent Interactions, in Artificial Intelligence Systems].

].

The game of the prisoner's dilemma is that each play is referred to as a round. That each agent can see what the opponent did in the previous round: player i can see whether j defected or not, and j can see whether i defected or not.

In the Collective Action Problem read p. 366 Constellations Volume 7, Number 3, 2015.
<http://homes.chass.utoronto.ca/~jheat>
 3, 2015

], people seem to behave irrationally best action, at a cost to others. This is the Prisoner's Dilemma problem, which we can look at from an agent perspective. This reading suggests that individuals who don't want to be the good one, unless everyone is doing good, then certainly it is going to be in your best interest to go back and reap the fruits of other's labour in the catch-22 problem.

While you can reap rewards by being good, so. Assuming that people are not fool (not always rational agents) is quite reasonable to assume others learn eventually from the whole system collapses. There would be no point to reap. Everyone being good is an unresolvable Prisoner's Dilemma.

Thus it is in our self-interest to be good, stable equilibrium and cooperation benefit is greater than individual greed. This is essentially the Nash equilibrium concept.

We argued these points through analysis of the prisoner's dilemma. The keen reader will notice that Equilibrium is a dominant strategy. The focus of this section is (cleverly) to look at rational agents as a group. Hence it is well why should the rational agent cooperate? The answer to this is that if you are rational, meaning they know your strategy and the advantage of them.

Results.

The reader should see that it is in a self-interest to be good because it is above two reasons. First, making bad decisions decreases future expected rewards. Second, if a group forms a Nash Equilibrium, which is not being good destroys the cooperation of agents and results in overall less reward.

Archiving Same Goals Through Just v

This section explores the issue of doing the same goals, but the methods have different goodness.

Roughly translated, being just means up to no good. A reading of "The Ring of Gyges" by Plato [footnote: "The Ring of Gyges" by Plato - Philosophy Lander. <http://philosophy.lander.edu/intro/ar> June 2, 2015].

>], an example of an unjust man vs. a just man.

> In the reading, it explores a sly, unjust man and a just man. Both of whom have the same goals but through different methods. Injustice and deceit can live a life achieving the same goals through just means. Intuition tells us that archery through deceitful means is certainly better than through honorable means. And indeed Plato's course textbook summarizes it well [footnote: page 108 of The Fundamentals of Ethics].

Certainly many immoral people are deceitful. But others are able to sleep well at night for well done (assassination, theft, betrayal) within a network of like-minded associates. Sometimes they get away with it, having a clear conscience and never regret the harm they have caused.

see whether i
defected or not. Now, for the sake o
assume that the agents will continue
every round will be followed by anoth
assumptions, what is the rational thi

If you know that you will be meeting
rounds, the incentive to defect appea
diminished, for two reasons.

Reason 1: If you defect now, your opp
defecting. Punishment is not possible
dilemma.

Reason 2: If you 'test the water' by
receive the sucker's payoff on the fi
are playing the game indefinitely, th
util) can be 'amortized' over the fut

When taken into the context of an inf
long) run, then the loss of a single
represent a small percentage of the o
if you play the prisoner's dilemma ga
cooperation is a rational outcome.

The summary is that selfish rationali
possible result. Maximizing individua
possible reward. There are more benef
collectively, if they both cooperated
that being good benefits more as a wh
individually.

Interestingly, humans display a syste
cooperative behavior in this, which g
rationales[footnote:
Tversky, Amos; Shafir, Eldar (2004).
similarity: selected writings. (PDF).
Technology Press. Retrieved July 26,
]. It makes sense because the idea of
into play, because the assumption tha
to reward or punish their partner is

Part II: Conclusion

A rational agent's self-interest is t
rewarding. Rational agents in a group
Nash equilibrium, so being good means
more rewards. In the face of punishme
selfish decisions decreases future ex
other agents are less inclined to coo
defecting agent.

Conclusion

I conclude that it is in our self-int
analyzed in the above two parts. Part
reaps more long term rewards, and Par
benefits more as a whole.

References

[1]

[2]

[3]

[4]

[5]

Certainly based on pure accomplishmen
same, since both can accomplish their
or unjust means. But the primary diff
chance for the dishonest man to get c
penalized. And sure, certainly the ba
with it. A major assumption here is t
it.

Let's break the argument into two pie
sometimes in lying, and suppose he ne

He might fail at deceit.

This is the more realistic argument.
underlying assumption is unattainable
entire life from the moment they were
outrageous assumption, we might as we
everything he could ever want in life
goal but taking the unjust method to
more risk, but at a chance to gain hi

But here is the kicker: there is no d
had assumed the two man achieved the
means. Then obviously it makes no sen
All the deceitful man has done is sav
Which may be valuable, but at a heavy

He never fails at deceit.

Even if we make the (unlikely) assump
capable of maintaining his composure
out of character, it still doesn't ma
riskier path as it leads to the same
save some time and effort. Seems that
difficult to accomplish (if not impos
layman. Sure, there may exist those o
enough to carry this out, but they mo
argument. It is likely the assumption
cannot be proved.

Results.

A lot of the arguments explored here
solid in logic analysis as from the r
But nevertheless, the fundamental arg
show that it's much easier to maintai
composure than a deceitful and unjust
it difficult to maintain a deceitful
risks. Consequently, if we act just,
self-interest to be good.

Conclusion

We've now explored this problem from
rational agent would act, and their p
in a rational agent's self-interest t
rewarding, because making bad (not go
decreases future expected rewards, an
the cooperation with other rational a
less reward.

And we've followed up on what do we d
different methods to solving a proble
good and some not. It seems that the
reasonably invalid, and even if it wa
be gained.

| The common result from the analysis s
> optimal choice.