


问题：这样的 AI 是真的吗？

题目描述：

我们设计的逻辑是, 如果狼抓到了羊, 就奖励 10 分. 如果狼撞到了障碍物, 就扣 1 分, 为了让狼尽快的抓, 所以每秒钟狼都会受到 0.1 分的惩罚

结束的时候, 分越高奖励就越好嘛

听着是不是特别科学



然后我们训练了 20W 次

确实

根据图表显示, 效果越来越差

原因是

狼发现大部分情况下他都吃不到羊.

然后追羊的过程中还会浪费很多时间, 最后还吃不到

草

那不如干脆一开始就一头撞死....

这样负的少一点

抑郁狼草啊

我和大佬研究了三天

笑死我了

最后发现这个原因的时候都他妈笑傻了

他妈 打不过就自杀

我们告诉了狼他面朝的方向, 他面前的东西是什么, 羊在哪里. 羊在不在视野里. 他的速度, 羊的速度. 之类一大堆的东西

但是因为狼想着自杀就完事了

所以羊的输入就等于没用了

他就判断面前有没有石头, 有就撞上去

草草草草草草草

所以我们当时设置了好多条件, 发现一点用没有

下午 7:53

这就是绝对理性 AI 我逐渐理解一切

就你也不知道什么条件神经网络用了 什么没用

一开始就不合理。

驱动意志体的基本逻辑是争取痛苦最小化，而不是争取收益最大化。

收益最大化只是痛苦最小化的一种策略，而根本不是可以相提并论的原则。

这也是为什么人类总是难以实现收益最大化的原因——因为当收益最大化与痛苦最小化冲突时，收益最大化常态性的被放弃。

而这是健康的。

人不可以常态性的靠着贪婪战胜痛苦，这既是错误的——因为这会极大概率（甚至是必然的）导致自己的毁灭，也是不义的——因为这更会造成 ta 人的灾难。

幸亏如此。

其实在 ta 们同意“这很合理”的时候，就已经双双错误了。

想要模拟人类，要去模拟人类的懦弱，而不是人类的贪婪。

懦弱自然会衍生出恰到好处的贪婪，其真实程度会让你吃惊。

但你如果反过来，得到的恰恰是这个荒谬的结果。

评论区：

Q: 搭配这篇，幸存靠的是休息，不是努力，应该会好理解一点。

<https://www.zhihu.com/answer/1413609374> (#懒散#)

Q: 有点不太明白“驱动意志体的基本逻辑是争取痛苦最小化，而不是争取收益最大化。”这句话，是因为意志是消耗品，所以基于无法避免的损耗转而选择更少消耗意志的方式吗？如果有看懂的朋友，希望能够解答一下，谢谢[开心]

B: 你好啊。我试着解释这句话本身，但是并不满意便一概删去了。于是我想换种方式回应你的疑问，其实答主在这句话下面做了同样具有经验提炼式的解释，但这可能仍然让人不满足：为什么贪婪无效于自己的痛苦还会造成他人的痛苦。所以我想最好理解这句话本身是亲身体会到追求更大收益来缓解痛苦的不可持续性，然而这总归不是一种好的体验。以上是我的粗浅解释，稍作参考。

Q: 非常感谢你的回复！（抱歉啦现在才看到[小情绪]）再结合以上几位层主的回复，感觉稍微有点眉目了，特别是狮子和恋爱的缘由那一篇，希望以后还能多多在评论区相遇~

B: 不必歉意。看到你的回复后，再看问题、答案和我当初语意不清的回复，是一种很不错的体验，所以也要谢谢你的回复呢。我现在是这样理解答主的思想的：【人在痛苦前的懦弱是优先于对享乐的追求】，现实中多少存在问题的人找不到问题的出路，于是选择了短暂的愉悦刺激来度过时间，就是最贴切的印证。谢谢阅读。

Q: 痛苦最小化是原则(目的)，而收益最大化是策略(手段)。这个准则在越大的群体尺度下，越具有刚性。联系之前关于“大规模组织体明知存在弊病，却不敢刮骨疗毒，反而等到最后破灭”，也是这个准则的应然逻辑。

而放在越小的尺度下，越容易存在例外。联系之前的“世间学问，无非真善美”那篇（<https://www.zhihu.com/answer/554109512>），超越现实的艺术流派实际上是放弃自身痛苦最小化，转而谋求人类痛苦最小化。并且“善、美”似乎也是在追求痛苦最小化，而它们需要接受“真”的检验，也是因为追求痛苦最小化。

不知道这样理解，对不对。[调皮]

Q: 读到模拟人类要模拟怯懦这一段我想起了前几天萦绕在我心里的一个问题。杰克伦敦小说中有一段描写，初读眼眶湿润，再读依旧如此：

不知疲倦的桨叶推着轮船日日夜夜空哐哐地前进。每一天看起来都和另一天都差不多，可巴克觉得出来，这天气可是越来越冷了。一天早晨，桨叶终于安静下来，“纳华”号上弥漫着一股不安分的气息。巴克和其他的狗都知道要有变故了。弗朗索斯牵着他们上了甲板。刚踏上冰冷的舱面，巴克的爪子就陷在软乎乎、好像泥一样的白东西里面。他打了个喷嚏，跳了开去。更多的白东西还在半空中往下落。他抖开一些，身上却又沾了很多。他好奇地闻闻，又伸出舌头舔了舔。这东西像火一样灼人，刚一入口，马上就没了。他有点纳闷，又试了试，还是一样。旁边的人哄笑起来，他不好意思，又不知道怎么回事，因为这是他第一次看到雪。

这个答案好像更能让我感触到这段文字为什么打动我，这力量来自何处。我在一只狗的身上体会到了人的情感，是细腻和深入人性的笔触，实现了人类情感的细致模拟，唤醒了移情与代入，文字也因此有了魂魄。

Q: 订阅。

请问如何理解“懦弱自然会衍生出恰到好处的贪婪，其真实程度会让你吃惊。”这句话？在最小化痛苦的前提下追求收益最大化，这不就是现实中的常态嘛？并不能让人很吃惊啊[好奇]能否展开说一下？

B: 我觉得最小化痛苦就是利益最大化了，不存在其他的利益最大化。所谓的其他利益最大，只是一种不切实际的理想，模拟割掉了其中必然发生的痛苦。

恰到好处的贪婪是由敬畏和恐惧带来的，它可能来源于懦弱，也可能来源于爱。没有畏，贪婪会让人因自恋缺少必要的敬，妨碍人去看到真实的客观世界和自然规律，也就会对掌控力有损了。

每个人的承受力不同，无法承受的点也不同。有时，很难说人的决策模型是恐惧失去还是追求幸福了，也很难说究竟是因为懦弱，还是因为爱。贪婪啊，没有人是不贪的。只是想着，能够最小痛苦的贪就好了。如果这在意的痛苦不是只与自己一个人有关，那就能算是爱了。

Q: 同意，痛苦的出现不会因为人的投机心理而改变。或者说，如果痛苦可以量化，那么可能无论采取何种行为痛苦的数值都是恒定的，“不确定”本身就是一种确定，“逃避痛苦”本身就是一种痛苦。重新过了一遍答案，原来这里的“真实程度”指的是“AI 运行后的结果与现实的契合度”的意思，第一遍硬是没读出来……

更新于 2023/10/14