

Policy-making for small-probability catastrophes*

Preliminary and incomplete

Michael Mandler[†]

Royal Holloway College, University of London

This version: November 2023

Abstract

I consider policy-making in the face of two types of small-probability catastrophes: environmental events that can lead to very low consumption and extermination events (such as large meteor strikes) that cut short an indefinitely long flow of future utility. A robustness requirement that a policy ranking should not be overturned by a small change in the distributions of outcomes will block any ranking of policies. The maxmin criterion for multiple priors can get around this impasse. In the environment setting, maxmin policymakers should minimize the likelihood of the tail event of very low consumption while in the extermination setting they should ignore the tail event where civilization survives until the very distant future. The former conclusion provides a rationale for Weitzman's Dismal Theorem while the latter conclusion validates conventional policy comparisons based on discounting.

JEL codes: D81, Q51, Q54

Keywords: climate change, dismal theorem, multiple priors, maxmin expected utility

*I thank John Conley, Prajit Dutta, Piotr Dworzak, Harvey Lederman, David Miller, Scott Page, Christian Tarsney, and Marcus Pivato for helpful advice and suggestions.

[†]Address: Department of Economics, Royal Holloway College, University of London, Egham, Surrey, TW20 0EX, UK. Email: m.mandler@rhul.ac.uk.

1 Introduction

Criticisms of conventional welfare analyses of climate change often combine arguments against expected utility theory with calls to avert global warming. Weitzman in his Dismal Theorem (2009, 2014) argued that the expected social loss from carbon emissions is infinitely great and that large sacrifices in current consumption to reduce emissions are therefore warranted. These positions lead to several policy and theoretical puzzles. When each policy option leads to an expected utility of $-\infty$, the options will be unranked, even when one of them leads to greater probabilities for the catastrophic outcomes. Although expected utilities that diverge underscore the magnitude of what is at stake, they fail to discriminate among policy choices. There is also an immiseration problem: if expected utility maximization implies that we should avoid catastrophic outcomes that bring unboundedly great losses, should we accept any sacrifice of current consumption that lessens the likelihood of such losses even slightly?¹

The expected utility delivered by a policy can equal $-\infty$ in Weitzman's work because society's utility function is not bounded below. Each policy generates a probability distribution over environmental outcomes and thus over future utility levels. If, as global warming increases and future utility levels fall, the probabilities that weight those outcomes diminish slowly then an expected utility of $-\infty$ can result. As a result, the rankings of policies become highly sensitive to the tails of the distributions that policies specify: small changes in the distribution of extreme outcomes can overturn policy judgments. So without precise knowledge of these tails, decision-making will be crippled. Expected utilities that diverge are not required for this conclusion: unbounded utilities are Weitzman's mathematical short-hand for policymakers' uncertainty about where the lower bounds lie.

To model imprecision about policies, the decision-maker in this paper will assign multiple priors – a set of probability distributions over outcomes – to each policy. Since nihilism will result if a policy ranking has to secure the backing of all of these distributions, a policymaker will evaluate options by the maxmin criterion: the value of a policy option will be the smallest expected utility the option can achieve as the distributions vary over the policymaker's priors

¹See Nordhaus (2009, 2011).

(Gilboa and Schmeidler (1989)). To eliminate the ties where many policies deliver a utility of $-\infty$, our policymaker will truncate the left tails (the worst outcomes) of the distributions assigned to each policy option, apply the maxmin criterion to the truncated options, and take the limit as the truncations become small. As long as there is a well-defined limit, a winning policy will emerge; policy discrimination can therefore be salvaged. With mild restrictions on priors, this method will recommend minimizing the likelihood of the most extreme negative outcomes, thus giving formal support for the guardrail policies backed by Stern and Stiglitz (2021) and many others. These conclusions hold even when every policy leads to infinitely negative utility levels.

The ‘fatness’ of the tails of the distributions of outcomes will drive the ranking of policies: when the tails are relevant, the options that can lead to distributions of outcomes with fatter left tails will be rejected. This conclusion echoes Weitzman, who put great emphasis on fat tails. But the absolute fatness of tails does not matter: it is *comparisons* of fatness that determine the winning policies.² The maxmin assumption shows that policy options can be ranked even when decision-makers exhibit a muscular aversion to ambiguity.

A second literature on catastrophic risks will clarify when the present approach deviates from conventional welfare analyses. ‘Long-termism’ in philosophy (see Bostrom (2003), Greaves and MacAskill (2021)) analyzes various civilization-ending calamities such as a large meteor collision or uncontrollable artificial intelligence.³ The reach of human civilization is assumed to be vast, encompassing quadrillions of lives when humanity does not perish prematurely, and the value of those lives, if they are lived, is undiscounted. In judging an expenditure that reduces the probability of an extermination event, the undiscounted utility of a quadrillion or more lives is functionally equivalent to an infinitely large benefit: virtually any expenditure that reduces the likelihood of calamity will be justified.

The two literatures consider the policy consequences of utility functions that are effectively unbounded: instead of global warming leading to an expected utility level that can be unboundedly negative, extermination events can lead to the loss of an expected utility level that can be unboundedly positive. But the environment and extermination cases are

²For related empirical work drawing on Weitzman, see Ackerman et al. (2010) and Dietz (2011).

³On AI risks, see Bostrom (2014), Ngo et al. (2023), Russell (2019), and Yudkowsky (2008).

dramatically different from the maxmin point of view. In the extermination model, each policy option will be assessed by one of the priors that leads to a *finite* expected value. When judging an option, a decision-maker will be guided by the distribution that assigns a small likelihood to the right tail of the distribution of outcomes, where civilization lasts into the far future; right tails therefore become unimportant. This conclusion backs the traditional economic objective functions that discount the far future into irrelevance, the very approach that the long-termist literature has aimed to overturn. So if, say, infectious diseases or bioterrorism exemplify the civilization-ending catastrophes of the proximate future and uncontrollable AI is a longer-term risk, the current analysis will recommend investing in protections against the first category.⁴

This paper will place the environment and extermination problems in a common mathematical framework where utility functions evaluate current consumption and a future benefit – the only difference is that in the environment model it is the utilities of the small values of the future benefit that diverge (to $-\infty$) while in the extermination model it is the utilities of the large values that diverge (to $+\infty$). In both cases the unbounded utilities will mean that small modeling differences can upend policy rankings. But to avoid duplication, I will emphasize the environment rather than the extermination model. Until we come to the maxmin criterion, every result in one model has a close parallel in the other setting.

Three results will lay out the difficulties with expected utilities that are either unbounded are where policymakers do not know where the bounds are. First, no pair of policies can be robustly ranked: arbitrarily small changes in the distributions of outcomes associated with policies can overturn any policy ranking.⁵ Second, if we model unknown bounds for utility functions by considering various truncations to the tail of extreme outcomes then policy recommendations become acutely sensitive to where the truncations are made. Third, optimal policies can reduce the current generation to extreme poverty: for any set of policy options and any level of current consumption $\varepsilon > 0$, no matter how small, there is a nearby set of policy options such that the optimal policy reduces current consumption to ε .

⁴See Jones (2016) for an economic analysis, with discounting, of extermination events. Discounting is of course prominent in most economic analyses of climate change. See Arrow et al. (1996), Stern (2007), and Nordhaus (2017, 2019).

⁵Environmental policy rankings can also fail to be robust with respect to other modeling changes besides changes in distributions. See Budolfson et al. (2017) and Gillingham et al. (2018).

The maxmin-cum-truncation proposal in this paper can rescue policy discrimination from this quagmire. But despite the mathematical similarities of the environment and extermination models, the terms of the rescue will depend on where the utility of the future benefit is unbounded. In the environment model, utility is unbounded ‘on the left’ where future consumption is small. The worst-case or minimum-utility evaluation of a policy option will then be determined by the distribution of the policymaker’s that has the fattest left tail – where low future consumption is assigned the greatest probability. Maxmin therefore endorses guardrail policies that avoid the extreme low-consumption outcomes of climate change: these policies make the minimum-utility evaluation as favorable as possible. In the extermination model, utility is unbounded ‘on the right’ and hence the maxmin focus on the worst case will evaluate a policy option using a distribution with a thin right tail that deems survival into the distant future, where utility is unboundedly great, to be unlikely. Little emphasis should therefore be placed on increasing the probability of very long survival times.

The main response from traditional decision theory to models where environmental policies can lead to expected utilities of negative infinity has been to object to unbounded utility functions, including those such as the CRRA utility that are common in applications. See, for example, Arrow and Priebisch (2014) or Pindyck (2011), which argue in different ways for bounds on utility. In the extermination model, the corresponding complaint would criticize utility or social welfare functions that place no upper bound on the value of the length of human civilization: in the orthodox view those functions should be bounded above.

Traditional decision theory moreover does not *assume* bounded utility functions; bounded utilities are deduced from axioms on preferences – the primitive of the subject. Specifically the classical von Neumann-Morgenstern assumption of continuity implies that any expected utility representation of an agent’s preferences must use a bounded utility function for outcomes. It is questionable however to take an agent’s or society’s preferences as given: advice on the environment requires policymakers to consider outcomes that no one has ever encountered. In that vein, I will pose an alternative continuity axiom that applies to individuals who are endowed with preferences over truncated domains of consumption and are trying to construct a preference over the full, nontruncated domain. This axiom, which is inconsistent

with classical continuity, will generate preferences represented by unbounded utilities.

When a single individual makes a decision that risks death, the agent implicitly assigns a finite utility level to death. But for environmental catastrophes, the worst-case outcomes are difficult to fathom: the 0 level of future consumption in this paper will represent not just many deaths but the end of society as we know it. Individuals therefore might not be able to identify an a priori lower bound on how bad they judge that outcome to be. Although agents might nevertheless make decisions that risk society-ending catastrophes, those decisions may be artifacts of being forced to choose rather than thought-through judgments.

Whether the case for unbounded utilities is right or wrong, no point in this paper turns on utility functions that are literally unbounded: comparably to Weitzman (2009), each result can be reposed using large but finite bounds. For example, the conclusion that orderings of policy options are not robust to small changes in the distributions of outcomes can be rewritten as ‘for any $\varepsilon > 0$ and any ranking of two policies, there exist finite bounds for utility functions such that the ranking can be overturned by some ε change in the distribution of outcomes that the policies induce.’ Unbounded utilities are a modeling convenience that allows for less convoluted statements.

2 The need for robustness: an example

The need for robust policy assessment when utility is unbounded can be obscured by the parametric details of a model. In the Weitzman (2009) model of the environment, there are two layers of uncertainty. Society’s future consumption c is governed by a family of log-normal distributions (or some other thin-tailed distribution) that differ with respect to the standard deviation s they assign to $\ln c$. That is, the conditional distribution of c given s is log-normal, but s is uncertain. The probability density function f for the distribution of s is assumed to obey a power law, $f(s) \propto s^{-k}$. The probabilities in the right tail of the distribution therefore diminish slowly: the tail is fat. Since the fat right tail of the distribution of s assigns relatively large probabilities to the large standard deviations, the probability of $\ln c$ diminishes slowly as $c \rightarrow 0$. Compared to the exponential speed of convergence of the probability of $\ln c$ to 0 when c is log-normally distributed with known

parameters, the left tail of the distribution of $\ln c$ is fat. When this distribution is combined with a CRRA utility, $u(c) = \frac{1}{1-\alpha}c^{1-\alpha}$ where $\alpha > 1$, utility can reach unboundedly negative values as $c \rightarrow 0$ and the expected utility of c can equal $-\infty$.⁶

Without the fatness assumption on the right tail of the distribution of s , the $-\infty$ expected utility conclusion would not obtain. Even putting aside the infinite expected values for s , most economists would be hard pressed to opine on the tail of the distribution of standard deviations, which suggests that policy advice ought not turn on such details. At a more technical level, assumptions on the exact form of the distributions of tails place emphasis on a secondary factor. The fatness of tails is neither necessary nor sufficient for expected utilities equal to $\pm\infty$: it is the unboundedness of utility functions that matters.

The significance of this example lies in the discontinuity between the distribution of outcomes and the expected utility of a policy option: a small change to the tail of the distribution of the parameter s can lead to a large utility change, from finite to infinite or vice versa. I will formalize this fragility in the next section. But the example suggests that unbounded utilities will make policy discrimination difficult, both because of the subtlety of the link between distributions and expected utility outcomes and because ties will become more likely (many options will deliver an expected utility of $-\infty$). Unbounded utilities thus do not lend support to guardrail policies that minimize climate change; they make all policy rankings problematic.

3 Policy paralysis

A policy will consist of current consumption $x_0 \geq 0$ and a probability distribution P for a nonnegative random x_1 that represents a future benefit. In a model of the environment, a policymaker chooses a level of, say, carbon abatement today that affects both x_0 and the distribution of consumption x_1 at some future date. In a model of extermination, the policymaker can invest in safeguards against a future civilization-ending calamity, thus reducing x_0 , and the benefit x_1 measures how long civilization survives or the consumption

⁶In Weitzman's exact formulation, the expected marginal utility of future consumption is infinitely great. See Nordhaus (2009).

stream of all who live until a calamity does arrive. A x will denote a pair (x_0, P) .

Let u_0 and u_1 be strictly increasing utility functions for current consumption and for the future benefit respectively. Neither function is assumed to be bounded above or below; utility levels can diverge as current consumption or the future benefit approaches 0 or ∞ . The utility of 0 current consumption or a 0 future benefit is allowed to equal $-\infty$.

Each policy is evaluated by the sum of the utility of its current consumption and the expected utility of its future benefit. So, for a policy $x = (x_0, P)$, define the overall utility of x by

$$u(x) = u_0(x_0) + \mathbb{E}_P u_1(x_1)$$

where expectations and hence $u(x)$ can equal ∞ or $-\infty$. I assume that, for every policy (x_0, P) , $P(0) = 0$, as for example when P is atomless. Appendix A will expand on the paper's formal concepts and lay out some technical assumptions, for example on when expectations and infinite utility levels are well-defined.

Definition 1 *The policy preference \succsim , a binary relation on policies, is defined by $x \succsim \hat{x}$ if and only if $u(x) \geq u(\hat{x})$.*

So when $u(x) = u(\hat{x}) = \infty$ or when $u(x) = u(\hat{x}) = -\infty$ holds, indifference between x and \hat{x} obtains. The strict policy preference \succ associated with \succsim is defined by $x \succ \hat{x}$ if and only if $x \succsim \hat{x}$ holds and $\hat{x} \not\succsim x$ does not hold.

Corresponding to the two main interpretations of x_1 as future consumption and as the length of civilization, two different assumptions can apply to u_1 . In the *environment model*, the utility of future consumption x_1 decreases without a lower bound as x_1 converges to 0 (as in Weitzman (2009, 2014)). In the *extermination model*, the utility of the survival date x_1 increases without bound as x_1 increases. Since u_1 is unbounded in both cases, the models violate the rules of expected utility theory. The extermination model notably cannot arise from a discounted sum of the utilities of future lives: that would bound u_1 .

Define a sequence of probability distributions $P(n)$ to *converge to* distribution P or $P(n) \rightarrow P$ if the maximum, over sets S of positive numbers, of the distance between $P(n)(S)$ and $P(S)$ converges to 0.⁷

⁷The reader should assume here and subsequently that any $S \subset \mathbb{R}_+$ to which a probability P is applied is measurable. See Appendix A.

Definition 2 A sequence of policies $\langle (x_0(n), P(n)) \rangle$ **converges to** policy (x_0, P) if $x_0(n) \rightarrow x_0$ and $P(n)$ converges to P .

When a sequence of policies $x(n)$ converges to x and n is sufficiently large, the difference between $x(n)$ and x is effectively unobservable.

If $P(n)$ converges to P then $\mathbb{E}_{P(n)}f \rightarrow \mathbb{E}_P f$ for any bounded and continuous real-valued function f . The distinctive feature of our setting is that the random variable of interest, $u_1(x_1)$, will not be bounded. So, even when $P(n) \rightarrow P$, $\mathbb{E}_{P(n)}u_1(x_1)$ might not converge to $\mathbb{E}_P u_1(x_1)$.⁸

Call the policy preference \succsim *robust* if its strict rankings are maintained at least weakly for all nearby policies. That is, for all policies x and \hat{x} with $x \succ \hat{x}$, if $x(n)$ and $\hat{x}(n)$ converge to x and \hat{x} respectively then $x(n) \succsim \hat{x}(n)$ for all n sufficiently large.⁹

The following example shows that \succsim fails to be robust: an arbitrarily small change in probability distributions can reverse a policy ranking.

Example 1 In the environment model, fix some $x_0 > 0$ and let the probability distributions P and \hat{P} satisfy $P(2) = \hat{P}(1) = 1$. Then $(x_0, P) \succ (x_0, \hat{P})$. For each positive integer n , let $P(n)$ satisfy $P(n)(2) = 1 - \frac{1}{n}$ and $P(n)(\varepsilon_n) = \frac{1}{n}$ where $0 < \varepsilon_n < 2$ is chosen so that $u_1(\varepsilon_n) < -n^2$, which the environment model permits. Then $P(n)$ converges to P . Setting $\hat{P}(n) = \hat{P}$ for each n , $\hat{P}(n)$ converges to \hat{P} . But since $u_0(x_0) + \mathbb{E}_{P(n)}u_1(x_1) \rightarrow -\infty$, we have $(x_0, \hat{P}(n)) \succ (x_0, P(n))$ for all n sufficiently large. ■

Theorem 1 Let x and \hat{x} be policies in either the environment or extermination model. Then there are sequences of policies $\langle x(n) \rangle$ and $\langle \hat{x}(n) \rangle$ that converge to x and \hat{x} respectively such that $x(n) \succ \hat{x}(n)$ for all n .

It follows from Theorem 1 that there are pairs of policies near x and \hat{x} where the policy near x is strictly preferred to the policy near \hat{x} and other pairs near x and \hat{x} where the policy near \hat{x} is strictly preferred to the policy near x .

⁸There are alternative definitions of convergence, provided in Appendix A, that would allow slight generalizations of Theorems 1 and 2 below.

⁹With this definition, a failure of robustness must pass a relatively demanding test: if \succsim fails to be robust it would also fail to meet the requirement that the strict rankings of \succsim remain strict for nearby policies.

Proofs are in Appendix B.

Theorem 1 shows that policy rankings are fragile with respect to vanishingly small changes in the distributions of x_1 : when utilities are unbounded, essentially any policy advice can be overturned by a slight and unobservable adjustment of the model. Utility comparisons therefore cannot tell policymakers how much current consumption should be invested in the future benefit.

If agents in the environment model hold that some risks of 0 future consumption to be acceptable or if agents in the extermination model discount the future sufficiently, utilities will be bounded. Theorem 1 then will not hold. But when utilities are bounded, policy rankings will depend delicately on where exactly the bounds lie.

To model this fragility, I will truncate the extreme values of the future benefit or its utility by cutting off the future benefit below or above a certain level. In the environment model, reset the domain of u_1 to consist of those x_1 greater than some cutoff $\varepsilon > 0$ and (to deal with one problem at a time) assume that u_1 is bounded above. The model would then appear to be well-behaved: a small change in policies cannot reverse a \succ ranking. Similarly, if in the extermination model we reset the domain to consist of those x_1 less than some $\delta > 0$ and assume that u_1 is bounded below then small changes again cannot overturn a \succ ranking of policies.

Truncations unfortunately address the symptom not the underlying problem. Although \succsim will robustly rank any pair of policies for a truncated model, the rankings that are generated will depend on the ε 's and δ 's that define the cutoffs: the ranking of policies can switch back and forth repeatedly as ε shrinks or δ grows. Since truncations will later play a key constructive role, the details of this conclusion will prove valuable.

Define a *sequence of truncations for policy* $x = (x_0, P)$ in the environment model by a sequence of numbers $\langle \varepsilon_n > 0 \rangle$ such that $\varepsilon_n \rightarrow 0$. Each ε_n is identified with the policy $x^n = (x_0, P^n)$ where P^n is the conditional distribution of P given $[\varepsilon_n, \infty)$. In the extermination model, a sequence of truncations for $x = (x_0, P)$ is defined by a $\langle \delta_n > 0 \rangle$ such that $\delta_n \rightarrow \infty$ and each δ_n is identified with $x^n = (x_0, P^n)$ where P^n is the conditional distribution of P given $[0, \delta_n]$.¹⁰

¹⁰In both models, these conditional distributions are unique for all large n . See Appendix A.

These truncations will serve as a stand-in for a wider family of potential truncations. For example, in the environment model we could instead declare that $u_1(x_1) = u_1(\varepsilon)$ for all future consumption levels x_1 less than ε . Theorem 2 below would then continue to hold exactly as stated.

Let the distance between two policies (x_0, P) and (x'_0, P') equal $|x_0 - x'_0| + d(P, P')$, where $d(P, P')$ is the distance between probability distributions P and P' . See Appendix A. This notion of distance defines ‘arbitrarily near’ in the result below.

Theorem 2 *Let x and \hat{x} be policies in either the environment or extermination model. Then there exist policies x^* and \hat{x}^* arbitrarily near to x and \hat{x} and truncations $\langle x^n \rangle$ and $\langle \hat{x}^n \rangle$ of x^* and \hat{x}^* such that $x^n \succ \hat{x}^n$ for infinitely many n and $\hat{x}^n \succ x^n$ for infinitely many n .*

Truncations thus convert the robustness failure of Theorem 1 into a sensitivity to the cutoffs of truncations.

Weitzman’s dismal theorem linked infinitely negative expected utilities with the policy recommendation to invest nearly all of x_0 in improvements in the distribution of x_1 : present-day immiseration can be rationalized as optimal. We can re-express this conclusion in a manner similar to the above paralysis results.

Suppose a set of policies X is assigned a set of indices I and let the superscript of a policy x^i in X denote the index. Define a sequence of sets of policies $\langle X(n) \rangle$ to converge to X if each $X(n)$ and X can be indexed by the same I and, for each $i \in I$, $x^i(n)$ converges to x^i . Given a set of policies X , $x \in X$ is \succsim -maximal in X if $x \succsim \hat{x}$ for all $\hat{x} \in X$.

The following Observation is essentially a corollary of Theorem 1.

Observation. Let the set of policies X be indexed by I . In either the environment or extermination model, for any $i \in I$ there exists a sequence of sets of policies $\langle X(n) \rangle$ that converges to X such that $x^i(n)$ is the unique \succsim -maximal policy in $X(n)$ for all n .

Since the choice of $i \in I$ in the Observation is arbitrary, any policy (x_0, P) in a feasible choice set of policies X is optimal, no matter how small its current consumption x_0 , when a policymaker chooses from some approximation of X .

4 A multiple priors solution

Rankings of policies will fail to be robust when utilities are unbounded and truncation will not by itself eliminate the fragility. To escape these difficulties, I will formalize a policymaker's inability to assign a single distribution to each policy option by assuming that each option is associated with a set of distributions or priors.¹¹ For each (x_0, P) , the policymaker believes that P' might also govern the benefit x_1 for at least some of the distributions P' near to P . Requiring a policy ranking to be approved by all of a policymaker's priors would return us to the paralysis of the previous section. I will therefore follow the maxmin rule of Gilboa and Schmeidler (1989): our policymaker will act as if the true distribution is the prior that minimizes the expected utility of the policy. This rule does not solve all of the policymaker's problems. In the environment model, many or even all policy options might deliver the same maxmin expected utility of $-\infty$. It is here that truncation comes to the rescue. When the tails of distributions are truncated expected utilities are finite and hence one policy can defeat all of its competitors. That victory moreover can persist in the limit as we decrease the size of the tails that are truncated.

Formally, the discrimination procedure has three steps. First, truncate the model to ensure that expected utility calculations are finite. Second, evaluate a policy $x = (x_0, P)$ based on the minimum utility delivered by the policies $x' = (x_0, P')$ where P' is drawn from an appropriate set of distributions near to P . For each pair of policy options, one of the options will then be at least weakly superior to the other. Third, let the interval of x_1 's that the truncations cut off decrease in size. If x defeats \hat{x} in the limit then x is the maxmin winner.

A policymaking option will now specify a level of current consumption x_0 and a set of distributions \mathcal{P} over the future benefit and can therefore be represented as a set of policies of the form $X = \{(x_0, P) : P \in \mathcal{P}\}$. I will call such a set of policies a *decision*.

Given a decision X and either a cutoff $\varepsilon_n > 0$ for the environment model or a cutoff $\delta_n > 0$ for the extermination model, define the decision X^n by applying the cutoff to each policy in X . Let $x^n = (x_0, P^n)$ be a representative policy in X^n , where P^n is a conditional

¹¹See Millner et al. (2013) and Lemoine and Traeger (2016) for applications of multiple priors to climate change.

distribution of P given $[\varepsilon_n, \infty)$ or $[0, \delta_n]$. A sequence of cutoffs $\langle \varepsilon_n \rangle$ with $\varepsilon_n \rightarrow 0$ or a sequence $\langle \delta_n \rangle$ with $\delta_n \rightarrow \infty$ thus defines a *sequence of truncated decisions* $\langle X^n \rangle$ for X .

Since we can view a decision X as a set of probability distributions, the maxmin criterion can be applied. Define the *min utility of decision* X for a truncation ε^n or δ^n by

$$u_{\min}^n(X) = \min_{x^n \in X^n} u(x^n).$$

(When X^n contains infinitely many policies, the min above should be replaced by inf but I will use the notation min throughout the text.)

Definition 3 *In either the environment or extermination model, decision X is **maxmin superior** to decision \hat{X} if for all sequences of truncated decisions for X and \hat{X} we have $u_{\min}^n(X) > u_{\min}^n(\hat{X})$ for all n sufficiently large.¹²*

Beyond the specific advice that maxmin superiority gives for environmental and extermination risks, which will come in the next two sections, the criterion enjoys two structural advantages.

First, maxmin superiority can rank policies that deliver the same infinite expected utility level. The following example illustrates.

Example 2 *In the environment model, let $x = (x_0, P)$ and $\hat{x} = (\hat{x}_0, \hat{P})$ be two policies (or equivalently singleton decisions) where $u(x)$ is finite and $u(\hat{x}) = -\infty$. Let x_α be the mixture of x and \hat{x} given by $(\alpha x_0 + (1-\alpha)\hat{x}_0, \alpha P + (1-\alpha)\hat{P})$ where $0 < \alpha < 1$.¹³ Although x_α delivers the better policy x with positive probability and would therefore appear to be superior to \hat{x} , policies x_α and \hat{x} share the same utility level $-\infty$. But for any positive sequence $\varepsilon_n \rightarrow 0$ we have $u(x_\alpha^n) > u(\hat{x}^n)$ for all n sufficiently large.¹⁴ A limit of truncations thus delivers the reasonable judgment that x_α is superior to \hat{x} . ■*

Second, maxmin superiority is robustness-proof: if we begin with two policies that are strictly ranked by \succsim then the decisions that contain all policies with nearby distributions

¹²This test shares some similarities to overtaking criteria in growth theory. See Brock (1970), Gale (1967), Koopmans (1963), and von Weizsäcker (1965).

¹³The distribution $\alpha P + (1-\alpha)\hat{P}$ assigns probability $\alpha P(S) + (1-\alpha)\hat{P}(S)$ to any $S \subset \mathbb{R}_+$.

¹⁴Given that $P(0) = \hat{P}(0) = 0$, a monotone convergence argument shows that $u(x_\alpha^n) \rightarrow u(x)$ and $u(\hat{x}^n) \rightarrow u(\hat{x})$ as $n \rightarrow \infty$.

will be maxmin ranked the same way. For each truncation of any policy, the minimum utility of all nearby policies will be achieved by the distribution that maximally increases the probability of the smallest future benefit level the truncation allows, ε_n in the environment model. Since these maxmin adjustments will have the same impact on all policies, the \succsim ranking of policies will not change. A policy thus cannot be selectively poisoned by adding small probabilities of arbitrarily small future benefit levels as in Example 1; due to the truncation the poison will be ignored. The imprecision of a policymaker's information about the tails of distributions therefore no longer undercuts decision-making.

To formalize this second point, I define the set of policies near to some $x = (x_0, P)$. Given some $d > 0$, a decision X is a *d-neighborhood of policy x* if it consists of all policies with current consumption x_0 and a probability distribution within distance d of P . See Section 3 and Appendix A for more precision and for a formulation of the result below using a variant definition of the distance between probability distributions.

Call the extermination model *pure* if u_1 is bounded below.¹⁵

Theorem 3 *Suppose, in either the environment model or the pure extermination model, that $x \succ \hat{x}$ for two policies x and \hat{x} . Then each sufficiently small neighborhood of x is maxmin superior to each sufficiently small neighborhood of \hat{x} : there is a $\bar{d} > 0$ such that the d -neighborhood of x is maxmin superior to the d -neighborhood of \hat{x} for all $0 < d < \bar{d}$.*

There is one knife-edge case where maxmin superiority is suspect. Suppose we allowed policies to specify probability distributions that assign strictly positive probability to $x_1 = 0$. For a decision X in the environment model that contains such a policy, no sequence of truncated decisions for X will pick up on this fact: the min utility levels given by $u_{\min}^n(X)$ will therefore be inappropriately large. Such a X might maxmin defeat a decision \hat{X} that contains only policies with $\hat{P}(0) = 0$ but the victory will be unconvincing.

¹⁵In any extermination model that is not pure, each neighborhood of any policy will contain a policy with an expected utility of $-\infty$. See Theorem 1.

5 A fat-tail theorem for the environment model

For an arbitrary pair of decisions X and \widehat{X} , there is no guarantee that the min utility of one of the decisions will dominate the min utility of the other for all sufficiently small truncations. Instead of $u_{\min}^n(X) > u_{\min}^n(\widehat{X})$ or the reverse inequality holding for all n sufficiently large, the lead might switch back and forth as n increases. But there will be a winner in the limit, a maxmin superior decision, if in the environment model the tails of the distributions given by X and \widehat{X} are suitably ordered by their likelihoods.

In the environment model, let the distribution P of some policy x in decision X have a fatter left tail than the distributions for policies in decision \widehat{X} . For climate decision-making, \widehat{X} is the guardrail option: it assigns smaller probabilities than x to all of the sufficiently bad outcomes. All else being equal, the fat left tail will count against x in comparisons with the policies in \widehat{X} . But that effect will be substantial only if x delivers utility $-\infty$. If instead x has finite utility, P 's fat left tail is unlikely to be decisive: the utility contribution of the tail will also be finite and hence the negative effect of its fatness will become arbitrarily small as we proceed further out the left tail. These two assumptions about x , the fat left tail of P and its $-\infty$ utility, are enough to conclude that the guardrail decision \widehat{X} is maxmin superior to X .

The maxmin criterion will be indispensable for this conclusion: there may well be policies in X that lead to greater utility than some policies in \widehat{X} but under maxmin a policymaker will ignore the high-utility policies in X . Maxmin also lets a policymaker remain undecided about whether the policies in a decision can deliver unbounded expected utility levels. The maxmin superiority of \widehat{X} over X requires only that *one* of the policies associated with X delivers an expected utility of $-\infty$.

Let $X' = \{(x'_0, P') : P' \in \mathcal{P}'\}$ be a decision. Define policy (x_0, P) to have a *fatter left tail than* policies in X' if P assigns larger probability than the distributions in X' to all sufficiently small left tails, that is, there exist $\varepsilon > 0$ and $k < 1$ such that $kP(S) > P'(S)$ for all $P' \in \mathcal{P}'$ and $S \subset [0, \varepsilon]$.¹⁶

Theorem 4 *Let X and \widehat{X} be two decisions in the environment model with positive current*

¹⁶See Rojo (1992) for related orderings.

consumption. If there is a policy $x \in X$ with a fatter left tail than the policies in \widehat{X} and $u(x) = -\infty$ then \widehat{X} is maxmin superior to X .

Theorem 4 permits every policy in both decisions X and \widehat{X} to deliver an expected utility of $-\infty$: policymakers can therefore make recommendations even when all decisions lead to negatively infinite utilities. Indeed, as mentioned, policymakers can stay agnostic about whether X and \widehat{X} will deliver an expected utility of $-\infty$. The most pressing objections to Weitzman (2009) therefore do not apply to the current approach. Weitzman’s emphasis on tails moreover is validated, though their importance lies in policy discrimination rather than in showing that expected utilities can equal $-\infty$.

The proof of Theorem 4 argues that since x delivers an expected utility of $-\infty$, the min utility of X must converge to $-\infty$ as the truncations of the left tail shrink. The result then follows immediately if each policy in \widehat{X} has finite or greater expected utility. So suppose there is a policy \widehat{x} in \widehat{X} that, like x , delivers an expected utility of $-\infty$. The race between \widehat{x} and x will then be decided by their left tails (which by themselves deliver utility $-\infty$) and, due to x ’s fatter left tail, \widehat{x} will emerge as the victor. Since x probability-dominates \widehat{x} in the left tail and u_1 is not bounded below, the tail advantage of \widehat{x} over x grows ever larger as the tail shrinks and will overtake any nontail advantage that x might enjoy. Notice that x might not determine the min utility of X : there might be policies in X even worse than x but if \widehat{x} defeats x then it must also defeat the smallest-utility policy in X .

Theorem 4 invokes a moderately weak fat tail assumption that, for some $k < 1$, $\widehat{P}(S) < kP(S)$ for all events S far enough out on the left tail and all \widehat{P} ’s for the policies in \widehat{X} , but it is not the weakest possible. If we had assumed the milder, $\widehat{P}(S) < P(S)$ for all events S far enough out on the left tail, then a nontail advantage of X could outweigh the tail advantage of \widehat{X} .

To apply maxmin superiority via Theorem 4, a policymaker will have to judge the tails of distributions. To see if our fat tail assumption holds, a policymaker will need to weigh claims about the tails of various distributions and determine if decision A makes extreme climate outcomes more likely than decision B. But the demands of decision-making are more modest than in Weitzman (2009), which first put the spotlight on fat tails. Truncations allow policymakers to avoid debates over whether the policies/distributions in a decision

necessarily lead to an expected utility of $-\infty$ or whether tails meet some target level of fatness. Whether a policy delivers $-\infty$ expected utility can turn on minute modeling details, as we saw in section 2, and there is no critical level of fatness that necessarily leads to $-\infty$ expected utility. Policymakers fortunately face an easier task: they just need to *compare* the likelihoods of tails for different decisions. Theorem 4's assumption that the x in X with the fatter tail delivers a utility of $-\infty$ could be satisfied simply because every policy in X has a fatter tail than the policies in \widehat{X} and the policymaker is open-minded about their utility levels.

The maxmin proposal makes little progress on the immiseration problem: it can be optimal to invest nearly all of x_0 on tail improvements if by so doing the tail of the distribution of x_1 becomes less fat. Immiseration conclusions however present more of a predicament for social welfare maximization than for efficiency. The variables x_0 and x_1 are the consumptions of different generations; and perhaps social welfare would be enhanced if the present generation suffered severe hardship for the greater good. An efficiency test instead asks how, given the present generation's consumption x_0 , future welfare would be most improved. Since x_0 can be fixed across decisions, the maxmin approach can answer that question.

6 Tail irrelevance in the extermination model

From the maxmin point of view, the environment and extermination models differ dramatically. In the environment model, the utility of x_1 is unbounded 'on the left' which renders the left tail of the distribution of x_1 all important. If one of the policies in a decision, due to its left tail, leads to a utility level of $-\infty$ then the minimum-utility policy in the decision will also deliver a utility of $-\infty$. In a contest with other decisions, tail comparisons will then be paramount. In the extermination model, the utility of x_1 is unbounded 'on the right' and hence the utility of some policies can diverge to ∞ . But the minimum expected utility levels delivered by a decision will usually be finite or smaller and the right tail will then be irrelevant. In fact if we narrow our view to one problem at a time by assuming in the extermination model that u_1 is bounded below then the left tail will be irrelevant too: the expected utility of the smallest intervals of x_1 will converge to 0 as the intervals shrink. The

centers of the distributions of the future benefit will then be in the driver's seat and the min utilities of decisions will be finite. Policy comparisons can therefore proceed conventionally without any need for truncation.

Each decision X in this section again denotes a set of policies that differ only with respect to their distributions of x_1 . Call an extermination model *pure* if u_1 is bounded below and let us say that decision X has the *finite min utility* if $\min_{x \in X} u(x)$ is finite. When X contains infinitely many policies, min should be replaced by inf but I will use the min notation below.

Decisions in the pure extermination model will normally have finite min utilities and all truncations will lead to these utilities. The truncations themselves are therefore superfluous and policymaking can proceed without them: a decision-maker will face a straightforward choice between options with finite utility.

Theorem 5 *Let decision X in the pure extermination model contain a policy with finite expected utility. Then X has a finite min utility $\min_{x \in X} u(x)$ and the min utility of X for a truncation will converge to $\min_{x \in X} u(x)$ as the truncations decrease ($\delta_n \rightarrow \infty$). A decision X will contain a policy with finite expected utility if X contains every policy sufficiently near to some policy x .*

So under mild conditions decisions in the pure extermination model will have finite expected utility.

To derive clear-cut policy advice from Theorem 5, we need only eliminate flukes where the finite min utilities happen to tie. I will continue to use a limit-of-truncations definition of an optimal policy, in part to provide a result that covers both the environment and extermination models. But as we have now seen truncations are usually unnecessary in the extermination model.

Let a policymaker choose from a finite feasible set of decisions \mathcal{X} . Let $X \in \mathcal{X}$ be a *maxmin optimum* if X is maxmin superior to all other policies in \mathcal{X} .¹⁷

The set of decisions \mathcal{X} is *generic* if the (possibly non-finite) min utilities of decisions in \mathcal{X} differ, that is, $\min_{x \in X} u(x) \neq \min_{x \in X'} u(x)$ when $X \neq X'$ and both are in \mathcal{X} . Genericity

¹⁷So if X is a maxmin optimum then, for all sequences of truncated decisions for the decisions in \mathcal{X} , the min utility of X is greater than the min utility of X' for truncation ε^n or δ^n , that is, $u_{\min}^n(X) > u_{\min}^n(X')$, for all decisions X' in \mathcal{X} besides X and for all n sufficiently large.

is best-suited to the extermination model. Under the assumption, at most one $X \in \mathcal{X}$ can have $\min_{x \in X} u(x) = -\infty$ and thus a leading case of the environment model, where many decisions have a min utility of $-\infty$, is excluded. In the pure extermination model, in contrast, Theorem 5 reports that decisions will normally have finite min utilities; genericity then just requires these finite utility levels to differ across decisions. However genericity comes to be satisfied, it implies that a finite set of decisions will have a unique maxmin optimum.

Theorem 6 *If the feasible set of decisions is finite and generic then in either the environment or extermination model there is a unique maxmin optimum.*

The upshot of Theorems 5 and 6 is that in the pure extermination model there will normally be a feasible X with finite utility that maxmin dominates all other feasible options and X 's superiority holds even if we ignore the right tail. The classical advice to ignore the distant future is valid.

The proof of Theorem 6 argues that the min utility of decision X for a truncation, $u_{\min}^n(X)$, converges to $\min_{x \in X} u(x)$ as the truncations shrink in size. Genericity and finiteness furthermore imply that only one decision X^* in the feasible set can maximize $\min_{x \in X} u(x)$. The min utility of X^* for all truncations sufficiently small must then outstrip the min utility of the other feasible decisions and is therefore the maxmin optimum.

7 Discussion: asymptotic discounting

The environment and extermination models come to different conclusions about the importance of the tails of the distribution of the future benefit, as Theorems 4 and 5 have shown.

In the environment model, utility diverges to $-\infty$ as future consumption x_1 decreases to 0; the left tails of the distribution of x_1 can therefore lead the expected utilities of policies to diverge to $-\infty$ as well. When this fact is combined with the maximin criterion, the left tails will drive the ranking of decisions. Since the likelihoods of the lowest future consumption levels will determine the min utilities, a policymaker should minimize the likelihood of the

consumption outcomes on the left tail – the guardrail advice that Weitzman (2009) backed informally.

In the pure extermination model, utility diverges to ∞ as the length of civilization x_1 increases. The distributions that govern the min utility will therefore assign high likelihood to an early end to civilization (small values of x_1) and thus deliver finite expected utility levels. A policymaker should therefore ignore the right tail and pay little attention to the dangers that are likely to end civilization only in the far future. In fact, in the pure extermination model the left tail is also unimportant: when the utility of x_1 is bounded below, the utility impact of the left tail will become vanishingly small as the lower bound that defines the left tail goes to 0.

Returning to the right tail, an asymptotic form of discounting or risk aversion appears in both the environment and extermination models. Just as the principle of diminishing marginal utility tells us to pay more for the high marginal utility gains that arise at small levels of consumption, maxmin optimization tells us to assign greater weight to small consumption levels and to ignore large civilization survival times. This advice moreover does not require u_0 or u_1 to be concave.

The tail asymmetry in our results – give priority to fat tails in the environment model but not in the extermination model – is driven by the pessimism of min utility evaluations. With a max or sup definition of the utility of a decision, the advice would be reversed: give priority to fat tails in the extermination model but not in the environment model. Dual to Theorem 4, if in the extermination model a policy x in decision X has a fatter right tail than the policies in \hat{X} and if $u(x) = \infty$ then X will be maxmax superior to \hat{X} .

8 Conclusion: unbounded utilities and classical decision theory

Weitzman (2009) and Stern and Stiglitz (2021) argue that expected utility theory does not apply when distributions are fat-tailed: fat tails combined with some standard utility functions lead to expected utilities that diverge. From the point of view of classical decision

theory, this criticism is hard to understand. In the von Neumann-Morgenstern setting where agents choose between known probability distributions (lotteries), axioms are placed on preference relations over probability distributions not on utility functions. Moreover, when preferences satisfy standard axioms it is a theorem that the utility functions that represent those preferences are bounded above and below – which is presumably why decision theory has mostly ignored Weitzman’s dismal theorem.¹⁸

According to the classical agenda, critics of expected utility theory should not posit unbounded utility functions; they should identify the axiom of choice under uncertainty they wish to challenge. It is clear which axiom the preferences in this paper (or the implicit preferences in Weitzman (2009)) violate: they fail to obey the continuity axiom. If we fix current consumption x_0 then in the language of this paper continuity states: if $(x_0, P) \succ (x_0, Q) \succ (x_0, R)$ then there exist $\alpha, \beta \in (0, 1)$ such that

$$(x_0, \alpha P + (1 - \alpha)R) \succ (x_0, Q) \succ (x_0, \beta P + (1 - \beta)R). \quad (\text{vNM continuity})$$

To see that the policy preference \succsim violates vNM continuity, suppose that

$$\infty > \mathbb{E}_P u_1 > \mathbb{E}_Q u_1 > \mathbb{E}_R u_1 = -\infty.$$

Then for every $\alpha \in (0, 1)$ we have $\mathbb{E}_Q u_1 > \mathbb{E}_{\alpha P + (1-\alpha)R} u_1$ and hence $(x_0, Q) \succ (x_0, \alpha P + (1 - \alpha)R)$, thus violating the first \succ in vNM continuity.

Instead of taking a position on whether or not vNM continuity is reasonable, I will show that it violates a different and equally compelling continuity axiom. Fix u_0 and u_1 and, for concreteness, a sequence of truncations in the environment model defined by $\varepsilon_n \rightarrow 0$. Consider an agent who for each n has a preference \succsim^n over the ε_n truncations defined by

$$(x_0, P^n) \succsim^n (\hat{x}_0, \hat{P}^n) \text{ if and only if } u_0(x_0) + \mathbb{E}_{P^n} u_1 \geq u_0(\hat{x}_0) + \mathbb{E}_{\hat{P}^n} u_1,$$

where P^n and \hat{P}^n are defined as in section 3. The \succsim^n fit together well: if $(x_0, P^n) \succsim^n (\hat{x}_0, \hat{P}^n)$ then, for any integer $i > 0$, $(x_0, P^n) \succsim^{n+i} (\hat{x}_0, \hat{P}^n)$. That is, a later preference in the sequence

¹⁸See Arrow and Priebsch (2014), however, which argues that sensible utilities should be bounded.

agrees with any earlier preference over the truncated policies that lie in the earlier preference's domain.

Now suppose that the agent wishes to extend his or her preferences to a \succsim that applies to nontruncated policies. Continuity considerations recommend the following axiom:

$$\text{if } x^n \succ^n \hat{x}^n \text{ for all } n \text{ sufficiently large then } x \succsim \hat{x}, \quad (\text{tail continuity})$$

where as usual $x^n = (x_0, P^n)$ and $\hat{x}^n = (\hat{x}_0, \hat{P}^n)$. But if we assume tail continuity then \succsim will violate vNM continuity. Fix some $\alpha \in (0, 1)$ and x_0 and let P , Q , and R be defined as before. Then $(x_0, Q^n) \succ^n (x_0, \alpha P^n + (1 - \alpha)R^n)$ for all n sufficiently large and hence, by tail continuity, $(x_0, Q) \succ (x_0, \alpha P + (1 - \alpha)R)$. Since this holds for all $\alpha \in (0, 1)$, we have a violation of vNM continuity. So even though our agent has nicely behaved preferences over truncated policies, the vNM continuity axiom will block their natural extension.

Truncations suggest a role for unbounded utility functions: they can summarize the preferences of agents who can judge the harm of diminishing consumption but are unsure how to evaluate very small levels of consumption. When in social decision-making 0 consumption represents the end of civilization, a judgment that the loss incurred has no obvious lower bound does not seem unreasonable.

A Appendix: technicalities

Further assumptions and conventions

I assume throughout the paper that, for each policy (x_0, P) , (1) $u_0(x_0) = -\infty$ and $E_P u_1(x_1) = \infty$ do not hold simultaneously and (2) $\mathbb{E}_P u_1(x_1)$ exists (with $\pm\infty$ allowed). As mentioned in the text, $P(0) = 0$ also holds for any (x_0, P) .

Due to (2), for any policy (x_0, P) the integrals with respect to P of the positive and negative parts of u_1 do not both equal ∞ .

Expectations will designate Lebesgue rather than Riemann integrals.

When a sequence of policies $x(n)$ converges to policy x , I will sometimes write $x(n) \rightarrow x$.

Convergence of probability distributions

The distance d between two measures P and P' on \mathbb{R}_+ is defined by $d(P, P') = \sup_S |P(S) - P'(S)|$ where the supremum is taken over all measurable $S \subset \mathbb{R}_+$. The measures do not have to be probability distributions. The sequence of measures $P(n)$ converges to measure P if $d(P(n), P) \rightarrow 0$. Weak convergence provides an alternative definition of convergence for probability measures and can be defined by the requirement that $\mathbb{E}_{P(n)}f \rightarrow \mathbb{E}_P f$ for any bounded and continuous $f : \mathbb{R}_+ \rightarrow \mathbb{R}$. Convergence as defined by d implies weak convergence.

Section 3

The two paralysis results, Theorems 1 and 2, will hold under slightly weaker conditions if they are restated using weak convergence. If ‘ $P(n)$ converges weakly to P ’ replaces ‘ $P(n)$ converges to P ’ in Definition 2 then Theorem 1 remains valid even if we no longer require a policy (x_0, P) to assign probability 0 to $x_1 = 0$, $P(0) = 0$. Theorem 2 will also continue to hold without the same requirement if a weak-convergence definition of the distance between probability distributions P and \hat{P} replaces the d given in A1. Specifically, distance could instead be given by the Lévy-Prokhorov metric d^* defined by: if F and \hat{F} are the distributions functions for P and \hat{P} then

$$d^*(P, \hat{P}) = \inf\{\gamma \in \mathbb{R}_+ : F(x_1 - \gamma) - \gamma \leq \hat{F}(x_1) \leq F(x_1 + \gamma) + \gamma \text{ for all } x_1 \in \mathbb{R}_+\}.$$

See Billingsley (1999).

In both the environment and extermination model, there is a unique conditional distribution of P given $[\varepsilon_n, \infty)$ and $[0, \delta_n]$ respectively for all n sufficiently large given by

$$P^n(S) = \frac{1}{P([\varepsilon_n, \infty))} P(S \cap [\varepsilon_n, \infty)) \quad \text{and} \quad P^n(S) = \frac{1}{P([0, \delta_n])} P(S \cap [0, \delta_n])$$

for each measurable $S \subset \mathbb{R}_+$.

In Theorem 2, a more precise statement of ‘arbitrarily near’ is that, for any distance $b > 0$, there exist policies x^* and \hat{x}^* such that the distance between x and x^* and the distance between \hat{x} and \hat{x}^* are both less than b .

Section 4

Using the definition of the distance d in A1 and letting $b > 0$, decision X is the b -neighborhood of (x_0, P) if $X = \{(x_0, P') : d(P', P) \leq b\}$.

If we use a metric for weak convergence to define b -neighborhoods then Theorem 3 will hold if we additionally assume that u_1 is Lipschitzian (see Gibbs and Su (2003) on the Kantorovich metric).

Section 6

In Theorem 5, a more precise statement of ‘ X contains every policy sufficiently near to some policy x' ’ is that there is a $b > 0$ and a policy x such that $x' \in X$ if the distance between x and x' is less than b .

B Appendix: proofs

Proof of Theorem 1. The proofs for the two models are similar. Consider the environment model and let $x = (x_0, P)$ and $\hat{x} = (\hat{x}_0, \hat{P})$. To construct $\langle x(n) \rangle$, let $\langle \varepsilon_n \rangle$ be a sequence of positive numbers with $\varepsilon_n \rightarrow 0$. For each n and measurable S , define $P(n)(S) = P(S \cap (\varepsilon_n, \infty)) + p$ where p equals $P([0, \varepsilon_n])$ if $\varepsilon_n \in S$ and 0 otherwise. Since, by continuity from below, $P(S \cap (\varepsilon_n, \infty)) \rightarrow P(S \setminus \{0\})$ and since $P(0) = 0$, $P(S \cap (\varepsilon_n, \infty)) \rightarrow P(S)$. Moreover, since $P(0) = 0$, continuity from above implies $P([0, \varepsilon_n]) \rightarrow 0$. Hence $P(n)(S) \rightarrow P(S)$ and $P(n) \rightarrow P$. Setting $x(n) = (x_0 + \varepsilon_n, P(n))$, we have $x(n) \rightarrow x$ and $u(x(n)) > -\infty$ for all n .

Let $\langle \hat{\delta}_n \rangle$ be a sequence of positive numbers such that $\hat{\delta}_n \rightarrow \infty$. Given some n and $0 < \hat{\varepsilon}_n < \hat{\delta}_n$, set $\hat{x}_0(n) = \hat{x}_0 + \hat{\varepsilon}_n$ and define $\hat{P}(n)$ by setting, for any measurable S , $\hat{P}(n)(S) = (1 - \frac{1}{n+1})\hat{P}(S \cap (\hat{\varepsilon}_n, \hat{\delta}_n)) + p + q$ where p equals $(1 - \frac{1}{n+1})\hat{P}([0, \hat{\varepsilon}_n]) + \frac{1}{n+1}$ if $\hat{\varepsilon}_n \in S$ and 0 otherwise and q equals $(1 - \frac{1}{n+1})\hat{P}([\hat{\delta}_n, \infty))$ if $\hat{\delta}_n \in S$ and 0 otherwise. Similarly to the previous paragraph, $\hat{P}(S \cap (\hat{\varepsilon}_n, \hat{\delta}_n)) \rightarrow \hat{P}(S \setminus \{0\})$ and, since $\hat{P}(0) = 0$, $\hat{P}(S \cap (\hat{\varepsilon}_n, \hat{\delta}_n)) \rightarrow \hat{P}(S)$. Moreover $\hat{P}([0, \hat{\varepsilon}_n]) \rightarrow 0$ and $\hat{P}([\hat{\delta}_n, \infty)) \rightarrow 0$. Hence $\hat{P}(n)(S) \rightarrow \hat{P}(S)$ and $\hat{P}(n) \rightarrow \hat{P}$. Since $\hat{x}_0(n) \rightarrow \hat{x}_0$, we have $\hat{x}(n) \rightarrow \hat{x}$.

Set $\langle \hat{\varepsilon}_n \rangle$ recursively by letting each $\hat{\varepsilon}_n$ satisfy $\hat{\varepsilon}_n \leq \frac{1}{2}\hat{\varepsilon}_{n-1}$ for $n > 1$, $0 < \hat{\varepsilon}_n < \hat{\delta}_n$, and

$$u_0(\hat{x}_0(n)) + \mathbb{E}_{\hat{P}(n)} u_1(x_1) = u_0(\hat{x}_0(n)) + \hat{P}(n)(\hat{\varepsilon}_n) u_1(\hat{\varepsilon}_n) + \int_{(\hat{\varepsilon}_n, \hat{\delta}_n]} u_1(x_1) d\hat{P}(n) < u(x(n)).$$

The inequality can be satisfied since (1) $\int_{(\hat{\varepsilon}_n, \hat{\delta}_n]} u_1(x_1) d\hat{P}(n)$ and $u_0(\hat{x}_0(n))$, each seen as a

function of $\widehat{\varepsilon}_n$, are bounded above by $u_1(\widehat{\delta}_n)$ and $u_0(\widehat{x}_0 + \widehat{\varepsilon}_0)$ respectively, (2) $\widehat{P}(n)(\widehat{\varepsilon}_n) \geq \frac{1}{n+1} > 0$ for all $\widehat{\varepsilon}_n$, (3) $u_1(\widehat{\varepsilon}_n) \rightarrow -\infty$ as $\widehat{\varepsilon}_n \rightarrow 0$, and (4) $u(x(n)) > -\infty$. So for $\widehat{x}(n) = (\widehat{x}_0(n), \widehat{P}(n))$, we have $u(x(n)) > u(\widehat{x}(n))$ for all n . ■

Proof of Theorem 2. As the proofs for the two models are again similar, assume the environment model. Fix a number $b > 0$ which will serve as a bound for the distance between x^* and x and between \widehat{x}^* and \widehat{x} . I will show that there exists a $\varepsilon_0 > 0$, a sequence of positive numbers $\langle \varepsilon_n \rangle$, and sequences of measures $\langle Q^n \rangle$ and $\langle \widehat{Q}^n \rangle$ such that:

$$u_0(\widehat{x}_0 + \varepsilon_0) + \frac{1}{\widehat{Q}^n(\mathbb{R}_+)} \int_{[\varepsilon_n, \infty)} u_1(x_1) d\widehat{Q}^n < u_0(x_0 + \varepsilon_0) + \frac{1}{Q^n(\mathbb{R}_+)} \int_{[\varepsilon_n, \infty)} u_1(x_1) dQ^n \quad (\text{Odd})$$

for all odd $n > 0$,

$$u_0(\widehat{x}_0 + \varepsilon_0) + \frac{1}{\widehat{Q}^n(\mathbb{R}_+)} \int_{[\varepsilon_n, \infty)} u_1(x_1) d\widehat{Q}^n > u_0(x_0 + \varepsilon_0) + \frac{1}{Q^n(\mathbb{R}_+)} \int_{[\varepsilon_n, \infty)} u_1(x_1) dQ^n \quad (\text{Even})$$

for all even $n > 0$. There will furthermore be probability distributions P^* and \widehat{P}^* such that (A) for each n , $\frac{1}{Q^n(\mathbb{R}_+)} Q^n$ and $\frac{1}{\widehat{Q}^n(\mathbb{R}_+)} \widehat{Q}^n$ are conditional distributions of P^* and \widehat{P}^* given $[\varepsilon_n, \infty)$ and (B) $d(P^*, P) + \varepsilon_0 < b$, and $d(\widehat{P}^*, \widehat{P}) + \varepsilon_0 < b$.

Given $0 < \varepsilon_0 < \delta$, define Q^0 by letting, for any measurable S , $Q^0(S) = P(S \cap (\varepsilon_0, \delta))$. For the same $0 < \varepsilon_0 < \delta$, define \widehat{Q}^0 by $\widehat{Q}^0(S) = \widehat{P}(S \cap (\varepsilon_0, \delta))$. Comparably to the proof of Theorem 1, set $\varepsilon_0 > 0$ small enough and δ large enough that $d(Q^0, P) + 1 - Q^0(\mathbb{R}_+) + \varepsilon_0 < b$ and $d(\widehat{Q}^0, \widehat{P}) + 1 - \widehat{Q}^0(\mathbb{R}_+) + \varepsilon_0 < b$.

Proceeding by induction, suppose that ε_n , Δ_n , $\widehat{\Delta}_n$, Q^n , and \widehat{Q}^n for $1 \leq n \leq m$ satisfy: (i) Odd and Even, (ii) $d(Q^n, P) < d(Q^0, P) + 1 - Q^0(\mathbb{R}_+)$ and $d(\widehat{Q}^n, \widehat{P}) < d(\widehat{Q}^0, \widehat{P}) + 1 - \widehat{Q}^0(\mathbb{R}_+)$, (iii) $0 < 1 - Q^n(\mathbb{R}_+) < \frac{1}{n}$ if n is odd and $0 < 1 - \widehat{Q}^n(\mathbb{R}_+) < \frac{1}{n}$ if n is even, (iv) $Q^n = Q^{n-1}$ if n is even and $\widehat{Q}^n = \widehat{Q}^{n-1}$ if n is odd, (v) $d(Q^n, Q^{n-1}) = \Delta_n$ and $Q^n(\mathbb{R}_+) = Q^{n-1}(\mathbb{R}_+) + \Delta_n$ if n is odd and $d(\widehat{Q}^n, \widehat{Q}^{n-1}) = \widehat{\Delta}_n$ and $\widehat{Q}^n(\mathbb{R}_+) = \widehat{Q}^{n-1}(\mathbb{R}_+) + \widehat{\Delta}_n$ if n is even, (vi) $\varepsilon_n < \varepsilon_{n-1}$, and (vii) the supports of Q^n and \widehat{Q}^n are contained in $[\varepsilon_n, \delta]$. Due to vii, $\int_{[\varepsilon_m, \infty)} u_1(x_1) dQ^m$ and $\int_{[\varepsilon_m, \infty)} u_1(x_1) d\widehat{Q}^m$ are finite. For concreteness, let m be even. Set $Q^{m+1} = Q^m$ and $\Delta_{m+1} = 0$. Since (1) $u_1(\varepsilon_{m+1}) \rightarrow -\infty$ as $\varepsilon_{m+1} \rightarrow 0$, (2) $\int_{[\varepsilon_m, \infty)} u_1(x_1) d\widehat{Q}^m$, $\int_{[\varepsilon_m, \infty)} u_1(x_1) dQ^m$, $u_0(\widehat{x}_0 + \varepsilon_0)$ and $u_0(x_0 + \varepsilon_0)$ are finite, and (3) $\int_{[\varepsilon_{m+1}, \infty)} u_1(x_1) dQ^{m+1} = \int_{[\varepsilon_m, \infty)} u_1(x_1) dQ^m$ for

any $\varepsilon_{m+1} \leq \varepsilon_m$, we can for any $\widehat{\Delta}_{m+1} > 0$ choose $0 < \varepsilon_{m+1}(\widehat{\Delta}_{m+1}) < \varepsilon_m$ so that

$$\begin{aligned} u_0(\widehat{x}_0 + \varepsilon_0) + \frac{1}{\widehat{Q}^m(\mathbb{R}_+) + \widehat{\Delta}_{m+1}} \widehat{\Delta}_{m+1} u_1(\varepsilon_{m+1}(\widehat{\Delta}_{m+1})) + \frac{1}{\widehat{Q}^m(\mathbb{R}_+) + \widehat{\Delta}_{m+1}} \int_{[\varepsilon_m, \infty)} u_1(x_1) d\widehat{Q}^m \\ < u_0(x_0 + \varepsilon_0) + \frac{1}{Q^{m+1}(\mathbb{R}_+)} \int_{[\varepsilon_{m+1}, \infty)} u_1(x_1) dQ^{m+1}. \end{aligned}$$

Fix $\widehat{\Delta}_{m+1}$ so that $0 < 1 - \widehat{Q}^m(\mathbb{R}_+) + \widehat{\Delta}_{m+1} < \frac{1}{m+1}$ and set $\varepsilon_{m+1} = \varepsilon_{m+1}(\widehat{\Delta}_{m+1})$. Define \widehat{Q}^{m+1} by $\widehat{Q}^{m+1}(S) = \widehat{Q}^m(S) + \widehat{\Delta}_{m+1}$ for any measurable S such that $\varepsilon_{m+1} \in S$ and $\widehat{Q}^{m+1}(S) = \widehat{Q}^m(S)$ for any measurable S such that $\varepsilon_{m+1} \notin S$, thus satisfying iii and v.

We then have

$$\begin{aligned} d(\widehat{Q}^{m+1}, \widehat{P}) &\leq d(\widehat{Q}^m, \widehat{P}) + d(\widehat{Q}^{m+1}, \widehat{Q}^m) = d(\widehat{Q}^m, \widehat{P}) + \widehat{\Delta}_{m+1} = d(\widehat{Q}^0, \widehat{P}) + \sum_{n=1}^{m+1} \widehat{\Delta}_n \\ &= d(\widehat{Q}^0, \widehat{P}) + \widehat{Q}^{m+1}(\mathbb{R}_+) - \widehat{Q}^0(\mathbb{R}_+) \leq d(\widehat{Q}^0, \widehat{P}) + 1 - \widehat{Q}^0(\mathbb{R}_+), \end{aligned}$$

where the inequalities are due to the triangle inequality and iii, the first equality is due to the definition of \widehat{Q}^{m+1} , and the final two equalities are due to v. Condition ii is therefore satisfied for $m+1$. Since i, iv, vi, and vii are evidently satisfied as well, the induction is complete.

Define P^* and \widehat{P}^* by $P^*(S) = Q^0(S \cap [\varepsilon_0, \delta]) + \sum_{n=1}^{\infty} Q^n(S \cap \{\varepsilon_n\})$ and $\widehat{P}^*(S) = \widehat{Q}^0(S \cap [\varepsilon_0, \delta]) + \sum_{n=1}^{\infty} \widehat{Q}^n(S \cap \{\varepsilon_n\})$ for each measurable S . Due to iii, $P^*(\mathbb{R}_+) = 1$ and $\widehat{P}^*(\mathbb{R}_+) = 1$ and so P^* and \widehat{P}^* are probability distributions. Since $d(Q^n, P^*) \rightarrow 0$ and $d(\widehat{Q}^n, \widehat{P}^*) \rightarrow 0$, taking the limit of ii gives $d(P^*, P) \leq d(Q^0, P) + 1 - Q^0(\mathbb{R}_+)$ and $d(\widehat{P}^*, \widehat{P}) \leq d(\widehat{Q}^0, \widehat{P}) + 1 - \widehat{Q}^0(\mathbb{R}_+)$. Given that $d(Q^0, P) + 1 - Q^0(\mathbb{R}_+) + \varepsilon_0 < b$ and $d(\widehat{Q}^0, \widehat{P}) + 1 - \widehat{Q}^0(\mathbb{R}_+) + \varepsilon_0 < b$, we have B. Define $x^* = (x_0 + \varepsilon_0, P^*)$ and $\widehat{x}^* = (\widehat{x}_0 + \varepsilon_0, \widehat{P}^*)$ and let the truncations $\langle x^n \rangle$ and $\langle \widehat{x}^n \rangle$ of x^* and \widehat{x}^* be defined by the $\langle \varepsilon_n \rangle$ given in the induction argument. Since $Q^n|_{[\varepsilon_n, \infty)} = P^*|_{[\varepsilon_n, \infty)}$ and $\widehat{Q}^n|_{[\varepsilon_n, \infty)} = \widehat{P}^*|_{[\varepsilon_n, \infty)}$, we have A. Odd therefore shows that $x^n \succ \widehat{x}^n$ for each odd n and Even shows that $\widehat{x}^n \succ x^n$ for each even n . ■

Proof of Theorem 3. I again restrict the proof to the environment model. Let $x^* = (x_0^*, P^*)$ and $\widehat{x} = (\widehat{x}_0, \widehat{P})$ satisfy $x^* \succ \widehat{x}$ and let X^* and \widehat{X} be the b -neighborhoods of x^* and \widehat{x} . (In the proof, I use b rather than d to denote a particular distance.) Since

$x^* \succ \hat{x}$, $u(x^*) > -\infty$ and therefore $x_0^* > 0$. Fix some positive sequence $\varepsilon_n \rightarrow 0$.

Let $(x_0, P) \in \{x^*, \hat{x}\}$ and $0 < b < 1$. Then $P([0, \varepsilon_n]) < 1 - b$ for all n sufficiently large. Fix n at such a value, allowing the construction of a probability distribution P' with $P'(\varepsilon_n) = P(\varepsilon_n) + b$. Let $X = \{(x_0, P') : d(P, P') \leq b\}$. Let $\delta > \varepsilon_n$ equal $\max\{\delta' : P([\delta', \infty)) \geq b\}$, where the max exists due to continuity from above. Then $\inf_{x^n \in X^n} u(x^n)$ is achieved by the P' defined by, for any measurable S , $P'(S) = P(S \cap [0, \delta)) + p + q$ where $p = b$ if $\varepsilon_n \in S$ and 0 otherwise and $q = P([\delta, \infty)) - b$ if $\delta \in S$ and 0 otherwise.

Fixing some n ,

$$u(x_0, P'^n) = u_0(x_0) + \frac{1}{1 - P'([0, \varepsilon_n])} \left(P(\varepsilon_n)u_1(\varepsilon_n) + bu_1(\varepsilon_n) + \int_{(\varepsilon_n, \delta)} u_1(x_1)dP + (P([\delta, \infty)) - b)u_1(\delta) \right).$$

For $x = x^*$ define $P^{*'} = P'$ and $\delta^* = \delta$, and for $x = \hat{x}$ define $\hat{P}' = P'$ and $\hat{\delta} = \delta$. Then, using the fact that

$$\frac{1}{1 - P^{*'}([0, \varepsilon_n])} - \frac{1}{1 - \hat{P}'([0, \varepsilon_n])} = \frac{P^{*'}([0, \varepsilon_n]) - \hat{P}'([0, \varepsilon_n])}{1 + P^{*'}([0, \varepsilon_n]) - \hat{P}'([0, \varepsilon_n]) + P^{*'}([0, \varepsilon_n])\hat{P}'([0, \varepsilon_n])},$$

we have

$$u(x_0^*, P^{*n}) - u(\hat{x}_0, \hat{P}^n) = u_0(x_0^*) - u_0(\hat{x}_0) + \frac{1}{1 - P^{*'}([0, \varepsilon_n])} P^{*'}(\varepsilon_n)u_1(\varepsilon_n) - \frac{1}{1 - \hat{P}'([0, \varepsilon_n])} \hat{P}'(\varepsilon_n)u_1(\varepsilon_n) \quad (1a)$$

$$+ \left(\frac{P^{*'}([0, \varepsilon_n]) - \hat{P}'([0, \varepsilon_n])}{1 + P^{*'}([0, \varepsilon_n]) - \hat{P}'([0, \varepsilon_n]) + P^{*'}([0, \varepsilon_n])\hat{P}'([0, \varepsilon_n])} \right) bu_1(\varepsilon_n) \quad (1b)$$

$$+ \frac{1}{1 - P^{*'}([0, \varepsilon_n])} \int_{(\varepsilon_n, \delta^*)} u_1(x_1)dP^{*'} - \frac{1}{1 - \hat{P}'([0, \varepsilon_n])} \int_{(\varepsilon_n, \hat{\delta})} u_1(x_1)d\hat{P}' \quad (1c)$$

$$+ \frac{1}{1 - P^{*'}([0, \varepsilon_n])} (P([\delta^*, \infty)) - b)u_1(\delta^*) - \frac{1}{1 - \hat{P}'([0, \varepsilon_n])} (P([\hat{\delta}, \infty)) - b)u_1(\hat{\delta}). \quad (1d)$$

Suppose first that $u(x_0^*, P^*)$ and $u(\widehat{x}_0, \widehat{P})$ are both finite. Then, for $P \in \{P^*, \widehat{P}\}$, as $n \rightarrow \infty$ we have $P(\varepsilon_n)u_1(\varepsilon_n) \rightarrow 0$, $P'([0, \varepsilon_n)) \rightarrow 0$, and $P(\mathbb{R}_+ \setminus \{\varepsilon_n\}) \rightarrow 1$. For $x \in \{x^*, \widehat{x}\}$, the finiteness of $u(x_0, P)$ implies that $\int_{[0, \varepsilon_n)} u_1(x_1)dP' \rightarrow 0$ as $n \rightarrow \infty$. For all n sufficiently large, $|u_1(x_1)| > |u_1(\varepsilon_n)|$ for all $x_1 < \varepsilon_n$. Hence $P^{*'}([0, \varepsilon_n))u_1(\varepsilon_n) \rightarrow 0$ and $\widehat{P}'([0, \varepsilon_n))u_1(\varepsilon_n)$ as $n \rightarrow \infty$. Therefore 1a and 1b converge to 0 as $n \rightarrow \infty$. For 1c and 1d, observe that $\delta \rightarrow \infty$ as $b \rightarrow 0$ and, since the finiteness of $u(x_0, P)$ implies that $\int_{[\delta, \infty)} u_1(x_1)dP \rightarrow 0$ as $\delta \rightarrow \infty$, $(P([\delta, \infty)) - b)u_1(\delta) \rightarrow 0$ as $b \rightarrow 0$. Hence 1d converges to 0 as $b \rightarrow 0$. Since the finiteness of $u(x_0, P)$ also implies that $\int_{(\varepsilon_n, \delta)} u_1(x_1)dP \rightarrow \int_{(\varepsilon_n, \infty)} u_1(x_1)dP$ as $\delta \rightarrow \infty$ and $\int_{(\varepsilon_n, \delta)} u_1(x_1)dP \rightarrow \int_{(\varepsilon_n, \delta)} u_1(x_1)dP$ as $n \rightarrow \infty$, for any $\varepsilon > 0$ there exist $\bar{\varepsilon} > 0$ and $\bar{\delta} > 0$ such that, for all $\varepsilon_n < \bar{\varepsilon}$ and $\delta > \bar{\delta}$, $\left| \int_{(\varepsilon_n, \delta)} u_1(x_1)dP - \int_{\mathbb{R}_+} u_1(x_1)dP \right| < \varepsilon$ and therefore $\left| u_0(x_0) + \int_{(\varepsilon_n, \delta)} u_1(x_1)dP - u(x) \right| < \varepsilon$. Given that $P^*(\mathbb{R}_+ \setminus \{\varepsilon_n\}) \rightarrow 1$ and $\widehat{P}(\mathbb{R}_+ \setminus \{\varepsilon_n\}) \rightarrow 1$ as $n \rightarrow \infty$, there is a $\bar{b} > 0$ such that, for all $0 < b < \bar{b}$, $u(x_0^*, P^{*n}) - u(\widehat{x}_0, \widehat{P}^n)$ and $u(x^*) - u(\widehat{x})$ will have the same sign for all n sufficiently large. So, for these values of b , $u(x_0^*, P^{*n}) > u(\widehat{x}_0, \widehat{P}^n)$ for all n sufficiently large.

Next suppose that $u(x_0^*, P^*)$ is finite and $u(\widehat{x}_0, \widehat{P}) = -\infty$. For 1a, as before $P^*(\varepsilon_n)u_1(\varepsilon_n) \rightarrow 0$ as $n \rightarrow \infty$ while $\widehat{P}'(\varepsilon_n)u_1(\varepsilon_n) < 0$ for all n sufficiently large. Similarly for 1b, $P^{*'}([0, \varepsilon_n))u_1(\varepsilon_n) \rightarrow 0$ and $\widehat{P}'([0, \varepsilon_n))u_1(\varepsilon_n) < 0$ for all large n . For 1d and either $P \in \{P^*, \widehat{P}\}$, $(P([\delta, \infty)) - b)u_1(\delta) \rightarrow 0$ as $b \rightarrow 0$ as before, again letting δ be a function of b . For 1c, given the finiteness of $u(x_0^*, P^*)$, there again exist $\bar{\varepsilon} > 0$ and $\bar{\delta} > 0$ such that, for all $\varepsilon_n < \bar{\varepsilon}$ and $\widehat{\delta} > \bar{\delta}$, $\left| u_0(x_0^*) + \int_{(\varepsilon_n, \widehat{\delta}^*)} u_1(x_1)dP^* - u(x^*) \right| < \varepsilon$. So, since $\int_{(\varepsilon_n, \widehat{\delta}^*)} u_1(x_1)dP^* \rightarrow -\infty$ as $\varepsilon_n \rightarrow 0$ for all $\widehat{\delta}^*$, there is a $\bar{b} > 0$ such that, for $0 < b < \bar{b}$, $u(x_0^*, P^{*n}) > u(\widehat{x}_0, \widehat{P}^n)$ for all n sufficiently large.

Finally suppose that $u(x_0^*, P^*) = \infty$ and $u(\widehat{x}_0, \widehat{P})$ is finite. For $P \in \{P^*, \widehat{P}\}$ and any $\varepsilon > 0$, $\int_{[0, \varepsilon]} u_1(x_1)dP^*$ must then be finite. Hence 1a and 1b converge to 0 as in the all-finite case. For 1d, $(\widehat{P}'([\widehat{\delta}, \infty)) - b)u_1(\widehat{\delta}) \rightarrow 0$ as $b \rightarrow 0$ and $(P^*([\delta^*, \infty)) - b)u_1(\delta^*) \geq 0$ for all b sufficiently small. For 1c, given the finiteness of $u(\widehat{x}_0, \widehat{P})$, there again exist $\bar{\varepsilon} > 0$ and $\bar{\delta} > 0$ such that, for all $\varepsilon_n < \bar{\varepsilon}$ and $\widehat{\delta} > \bar{\delta}$, $\left| u_0(\widehat{x}_0) + \int_{(\varepsilon_n, \widehat{\delta})} u_1(x_1)d\widehat{P} - u(\widehat{x}) \right| < \varepsilon$. Since $\int_{(\varepsilon_n, \widehat{\delta}^*)} u_1(x_1)dP^* \rightarrow \infty$ as $\widehat{\delta}^* \rightarrow \infty$, uniformly for all ε_n , there is a $\bar{b} > 0$ such that, for $0 < b < \bar{b}$, $u(x_0^*, P^{*n}) > u(\widehat{x}_0, \widehat{P}^n)$ for all n sufficiently large. ■

Proof of Theorem 4. Let $x^* = (x_0^*, P^*)$ be the posited policy with the fatter left tail. Given a policy $\bar{x} = (\bar{x}_0, \bar{P})$ and a measurable $S \subset \mathbb{R}_+$, define $u_S(\bar{x}) = u_0(\bar{x}_0) + \int_S u_1(x_1) d\bar{P}$. Since x^* has a fatter left tail than the policies in \hat{X} , there is a $\bar{\varepsilon} > 0$ and $\bar{\gamma} > 0$ such that $\frac{\int_{[\gamma, \bar{\varepsilon}]} u_1(\hat{x}_1) d\hat{P}}{\int_{[\gamma, \bar{\varepsilon}]} u_1(x_1^*) dP^*} < k$ for all $(\hat{x}_0, \hat{P}) \in \hat{X}$ and all $0 < \gamma \leq \bar{\gamma}$. Since $u(x^*) = -\infty$, $P^*(0) = 0$, and $u_0(x_0^*)$ is finite, $\int_{[\gamma, \bar{\varepsilon}]} u_1(x_1^*) dP^* \rightarrow -\infty$ as $\gamma \rightarrow 0$. In combination with the fact that $u_0(\hat{x}_0)$ is finite for $(\hat{x}_0, \hat{P}) \in \hat{X}$, it follows that there is a $\varepsilon' > 0$ and $\gamma' > 0$ such that $\frac{u_{[\gamma, \varepsilon']}(\hat{x})}{u_{[\gamma, \varepsilon']}(x^*)} < k$ for all $\hat{x} \in \hat{X}$ and all $0 < \gamma \leq \gamma'$. Since $u(x^*)$ is a Lebesgue integral (plus a finite constant), the assumption that $u(x^*) = -\infty$ implies that $u_{(\tilde{\varepsilon}, \infty)}(x^*)$ is finite for each $\tilde{\varepsilon} > 0$. Since in addition $\lim_{\gamma \rightarrow 0} u_{[\gamma, \tilde{\varepsilon}]}(x^*) = -\infty$ for any $\tilde{\varepsilon} > 0$, there is a $\varepsilon'' > 0$ and $\gamma'' > 0$ such that $\frac{u_{[\gamma, \varepsilon'']}(\hat{x})}{u_{[\gamma, \infty]}(x^*)} < k$ for all $\hat{x} \in \hat{X}$ and $0 < \gamma \leq \gamma''$. Partition \hat{X} into $\hat{X}_1 = \{\hat{x} \in \hat{X} : u_{(\tilde{\varepsilon}, \infty)}(\hat{x}) \geq 0 \text{ for all } \tilde{\varepsilon} > 0\}$ and $\hat{X}_2 = \hat{X} \setminus \hat{X}_1$. If $\hat{x} \in \hat{X}_1$ then for any pair of sequences of truncated decisions $\langle X^n \rangle$ and $\langle \hat{X}^n \rangle$ there is a n''' such that $u(\hat{x}^n) > u(x^{*n})$ for all $n > n'''$. For any $\hat{x} \in \hat{X}_2$, $u_{(\varepsilon'', \infty)}(\hat{x})$ must lie in the bounded interval $[u_0(\hat{x}_0) + u_1(\varepsilon''), 0]$. Since in addition $\lim_{\gamma \rightarrow 0} u_{[\gamma, \infty]}(x^*) \rightarrow -\infty$, for any $\delta > 0$ there exists a $\tilde{\gamma} > 0$ such that $\frac{u_{(\varepsilon'', \infty)}(\hat{x})}{u_{[\gamma, \infty]}(x^*)} < \delta$ for all $\hat{x} \in \hat{X}_2$ and $0 < \gamma \leq \tilde{\gamma}$. Since

$$\frac{u_{[\gamma, \infty]}(\hat{x})}{u_{[\gamma, \infty]}(x^*)} = \frac{u_{[\gamma, \varepsilon'']}(\hat{x})}{u_{[\gamma, \infty]}(x^*)} + \frac{u_{(\varepsilon'', \infty)}(\hat{x})}{u_{[\gamma, \infty]}(x^*)},$$

the fact that $\frac{u_{[\gamma, \varepsilon'']}(\hat{x})}{u_{[\gamma, \infty]}(x^*)} < k$ for all $\hat{x} \in \hat{X}_2$ and $0 < \gamma \leq \gamma''$ implies there is a $\tilde{\gamma} > 0$ such that $\frac{u_{[\gamma, \infty]}(\hat{x})}{u_{[\gamma, \infty]}(x^*)} < k$ for all $\hat{x} \in \hat{X}_2$ and $0 < \gamma \leq \tilde{\gamma}$. Hence for any pair of sequences of truncated decisions $\langle X^n \rangle$ and $\langle \hat{X}^n \rangle$ and any $\hat{x} \in \hat{X}$ there exists n^* such that $u(\hat{x}^n) > u^n(x^{*n})$ for all $n \geq n^*$. Given that $\inf_{x^n \in X^n} u(x^n) \leq u(x^{*n})$ for each n , \hat{X} is maxmin superior to X . ■

Proof of Theorem 5. Letting $\hat{x} = (x_0, \hat{P})$ be the posited policy with finite utility, $u_0(x_0)$ is finite. Since the extermination model is pure, $u_1(0)$ is finite and therefore $u(x') \geq u_0(x_0) + u_1(0) > -\infty$ for each $x' \in X$. Hence $\inf_{x' \in X} u(x')$ is finite or greater. Since $\inf_{x' \in X} u(x') \leq u(\hat{x})$, $\inf_{x' \in X} u(x')$ is finite. Given any positive sequence $\delta_n \rightarrow \infty$, $\lim_{n \rightarrow \infty} u_{\min}^n(X) = \inf_{x' \in X} u(x')$ and hence the min utility of X for $\langle \delta^n \rangle$ is $\inf_{x' \in X} u(x')$.

Let $x = (x_0, P)$ be the policy given in the final sentence of the Theorem. Since $P(0) = 0$, for $\varepsilon > 0$ sufficiently small and δ sufficiently large the policy $\tilde{x} = (x_0 + \varepsilon, \tilde{P})$, where \tilde{P} is the conditional distribution of P given $[\varepsilon, \delta]$, is within distance d of x . So $\tilde{x} \in X$. Since

$u_0(x_0 + \varepsilon)$ and $\int_{[\varepsilon, \delta]} u_1(x_1) d\tilde{P}$ are both finite, so is $u(\tilde{x})$. ■

Proof of Theorem 6. Due to the maintained assumption that each policy (x_0, P) has $P(0) = 0$, for any $\varepsilon_n \rightarrow 0$ or $\delta_n \rightarrow \infty$, $\lim_{n \rightarrow \infty} u_{\min}^n(X) = \inf_{x \in X} u(x)$ for each $X \in \mathcal{X}$. Genericity and finiteness then imply there is a $X^* \in \mathcal{X}$ such that $\lim_{n \rightarrow \infty} u_{\min}^n(X^*) > \lim_{n \rightarrow \infty} u_{\min}^n(X)$ for all $X \in \mathcal{X} \setminus \{X^*\}$. Hence for any $\varepsilon_n \rightarrow 0$ or $\delta_n \rightarrow \infty$ there is a \bar{n} such that $u_{\min}^n(X^*) > u_{\min}^n(X)$ for all $X \in \mathcal{X} \setminus \{X^*\}$ and $n > \bar{n}$. ■

References

- [1] Ackerman, F., Stanton, E., and Bueno, R., 2010. ‘Fat tails, exponents, extreme uncertainty: simulating catastrophe in DICE.’ *Ecological Economics* 69: 1657-1665.
- [2] Arrow, K., Cline, W., Maler, K., Munasinghe, M., Squitieri, R., and Stiglitz, J., 1996. ‘Intertemporal equity, discounting, and economic efficiency.’ *Climate Change 1995: Economic and Social Dimensions of Climate Change*, J. Bruce, H. Lee, and E. Haites, eds, Cambridge: Cambridge University Press.
- [3] Arrow, K. and Priebsch, M., 2014. ‘Bliss, catastrophe, and rational policy.’ *Environmental and Resource Economics* 58: 491-509.
- [4] Billingsley, P., 1999. *Convergence of Probability Measures, Second Edition*. New York: Wiley.
- [5] Bostrom, N., 2003. ‘Astronomical waste: the opportunity cost of delayed technological development.’ *Utilitas* 15: 308-314.
- [6] Bostrom, N., 2014. *Superintelligence: Path, Dangers, Strategies*. Oxford: Oxford University Press.
- [7] Brock, W., 1970. ‘An axiomatic basis for the Ramsey-Weizsäcker overtaking criterion.’ *Econometrica* 38 927-929.
- [8] Budolfson, M., Dennig, F., Fleurbaey, M., Siebert, A., and Socolow, R., 2017. ‘The comparative importance for optimal climate policy of discounting, inequalities and catastrophes.’ *Climatic Change* 145: 481-494.
- [9] Dietz, S., 2011. ‘High impact, low probability? An empirical analysis of risk in the economics of climate change.’ *Climatic Change* 108: 519-541.
- [10] Gale, D., 1967. ‘On optimal development in a multisector economy.’ *Review of Economic Studies* 34: 1-18.
- [11] Gibbs, A. and Su, F., 2002. ‘On choosing and bounding probability metrics.’ *International Statistical Review* 70: 419-435.

- [12] Gilboa, I. and Schmeidler, D., 1989. ‘Maxmin expected utility with non-unique prior.’ *Journal of Mathematical Economics* 18: 141-153.
- [13] Gillingham, K., Nordhaus, W., Anthoff, D., Blanford, G., Bosetti, V., Christensen, P., McJeon, H., and Reilly, J., 2018. ‘Modeling uncertainty in integrated assessment of climate change: A multimodel comparison.’ *Journal of the Association of Environmental and Resource Economists* 5: 791-826.
- [14] Greaves, H. and MacAskill, W., 2021. ‘The case for strong longtermism.’ *GPI Working Paper No. 5-2021*, Global Priorities Institute.
- [15] Jones, C., 2016. ‘Life and Growth.’ *Journal of Political Economy* 124: 539-578.
- [16] Koopmans, T., 1963. ‘On the concept of optimal economic growth.’ *Cowles Foundation Discussion Paper 392*, Yale University.
- [17] Lemoine, D. and Traeger, C., 2016. ‘Ambiguous tipping points.’ *Journal of Economic Behavior and Organization* 132: 5-18.
- [18] Millner, A., Dietz, S., and Heal, G., 2013. ‘Scientific ambiguity and climate policy.’ *Environmental and Resource Economics* 55: 21-46.
- [19] Ngo, R., Chan, L., and Mindermann, S., 2023. ‘The alignment problem from a deep learning perspective.’ arXiv:2209.00626.
- [20] Nordhaus, W., 2009. ‘An analysis of the dismal theorem.’ *Cowles Foundation Discussion Paper 1686*, Yale University.
- [21] Nordhaus, W., 2011. ‘The economics of tail events with an application to climate change.’ *Review of Environmental Economics and Policy* 5: 240-257.
- [22] Nordhaus, W., 2017. ‘Revisiting the social cost of carbon.’ *Proceedings of the National Academy of Sciences* 114: 1518-1523.
- [23] Nordhaus, W., 2019. ‘Climate change: the ultimate challenge for economics.’ *American Economic Review* 109: 1991-2014.
- [24] Pindyck, R., 2011. ‘Fat tails, thin tails, and climate change policy.’ *Review of Environmental Economics and Policy* 5: 258-274.
- [25] Rojo, J., 1992. ‘A pure-tail ordering based on the ratio of the quantile functions.’ *Annals of Statistics* 20: 570-579.
- [26] Russell, S., 2019. *Human Compatible: Artificial Intelligence and the Problem of Control*. New York: Viking.
- [27] Stern, N., 2007. *The Economics of Climate Change: The Stern Review*. Cambridge: Cambridge University Press.

- [28] Stern, N. and Stiglitz, J., 2021. ‘The social cost of carbon, risk, distribution, market failures: an alternative approach.’ NBER Working Paper Series, Working Paper 28472, NBER.
- [29] Stern, N., Stiglitz, J., and Taylor, C., 2022. ‘The economics of immense risk, urgent action and radical change: towards new approaches to the economics of climate change.’ *Journal of Economic Methodology* 29: 181-216.
- [30] von Weizsäcker, C., 1965. ‘Existence of optimal programs of accumulation for an infinite time horizon.’ *Review of Economic Studies* 32: 85-104.
- [31] Weitzman, M., 2009. ‘On modeling and interpreting the economics of catastrophic climate change.’ *Review of Economics and Statistics* 91: 1-19.
- [32] Weitzman, M., 2014. ‘Fat tails and the social cost of carbon.’ *American Economic Review Papers and Proceedings* 104: 544-546.
- [33] Yudkowsky, E., 2008. ‘Artificial intelligence as a positive and negative factor in global risk’ in *Global Catastrophic Risks*, N. Bostrom and M. Čirković, eds. New York: Oxford University Press.