CrossMark

# Uncertainty, Learning and International Environmental Agreements: The Role of Risk Aversion

Alistair Ulph[1] · Pedro Pintassilgo[2] · Michael Finus[3,4]

## Abstract

This paper analyses the formation of international environmental agreements (IEAs) under uncertainty, learning and risk aversion. It bridges two strands of the IEA literature: (i) the role of learning when countries are risk neutral; (ii) the role of risk aversion under no learning. Combining learning and risk aversion seems appropriate as the uncertainties surrounding many international environmental problems are large, often highly correlated (e.g. climate change), but are gradually reduced over time through learning. The paper analyses three scenarios of learning. A key finding is that risk aversion can change the ranking of these three scenarios of learning in terms of welfare and membership. In particular, the negative conclusion about the role of learning in a strategic context under risk neutrality is qualified. When countries are significantly risk averse, then it pays them to wait until uncertainties have been largely resolved before joining an IEA. This may suggest why it has been so difficult to reach an effective climate change agreement.

**Keywords** International environmental agreements · Uncertainty · Learning and risk aversion · Game theory

**JEL Classification** C72 · D62 · D80 · Q54

✉ Alistair Ulph
  alistair.ulph@manchester.ac.uk

  Pedro Pintassilgo
  ppintas@ualg.pt

  Michael Finus
  michael.finus@uni-graz.at

[1] Faculty of Humanities, School of Social Sciences, Sustainable Consumption Institute, The University of Manchester, Manchester M13 9PL, UK

[2] Faculty of Economics and CEFAGE, University of Algarve, Faro, Portugal

[3] Department of Economics, University of Graz, Graz, Austria

[4] Department of Economics, University of Bath, Bath BA2 7AY, UK

## 1 Introduction

Environmental issues such as climate change pose three key challenges for economic analysis: (i) there are considerable uncertainties about the likely future costs of environmental damages and abatement; (ii) our understanding of these uncertainties changes over time as a result of learning more about climate science, possible technological responses and behavioral responses by households, firms and governments; (iii) the problem is global, but since there is no single global agency to tackle climate change, policies need to be negotiated through international environmental agreements (IEAs).[1,2] Recently, these three issues have begun to be integrated in one framework. Two strands of literature can be distinguished.

The first strand of literature studies uncertainty and IEA formation with the focus on the role of *learning*, but under the assumption of risk neutrality. Ulph and Ulph (1996) and Ulph and Maddison (1997) compare the fully cooperative and the non-cooperative scenarios when countries face uncertainty about damage costs. They show that the value of learning about damage costs may be negative when countries act non-cooperatively and damage costs are correlated across countries. Na and Shin (1998), Ulph (2004), Kolstad (2007), Kolstad and Ulph (2008, 2011) have considered how the prospect of future resolution of uncertainty affects the incentives for countries to join an IEA. Kolstad and Ulph consider a model where countries face common uncertainty about the level of environmental damage costs.[3] Three scenarios of learning are considered: with *Full Learning*, uncertainty about damage costs is resolved before countries decide whether to join an IEA; with *Partial Learning*, uncertainty is resolved after countries decide whether to join an IEA, but before they choose their emissions levels; with *No Learning*, uncertainty is neither resolved before stage 1 nor stage 2. They showed that the prospect of learning, either full or partial, generally reduces the expected aggregate payoff in stable IEAs. In particular, Kolstad and Ulph (2008, 2011) showed that *Partial Learning* would yield the highest aggregate payoff for only a small proportion of parameter values. For a significant majority of parameter values, the highest expected aggregate payoff arose under *No Learning*. Hence, it is better to form an IEA before waiting for better information: removing the "veil of uncertainty" seems to be detrimental to the success of international environmental cooperation.

All these models have assumed that countries are risk neutral. However, in the climate context, risks are highly correlated and hence possibilities for risk sharing are limited so that the assumption of risk aversion may be quite relevant. Therefore, we extend the two-stage coalition formation setting by Kolstad and Ulph (2008) by departing from the assumption of risk neutrality. In this paper, we allow for countries to be risk averse, and show that if countries have a relatively high degree of risk aversion, then for a majority of parameter values *Full Learning* yields higher expected aggregate utility than *No Learning*. This may help to explain why it has taken such a long time between the start of the process of tackling climate change (the Kyoto Protocol) to reach a more substantial agreement in Paris—countries are risk averse

---

[1]  On the first two issues, see for instance Arrow and Fisher (1974), Epstein (1980), Kolstad (1996a, b), Ulph and Ulph (1997), Gollier et al. (2000) as well as Narain et al. (2007).

[2]  On the third issue, see for instance the classic papers by Carraro and Siniscalco (1993) and Barrett (1994). The most influential papers have been collected in a volume by Finus and Caparrós (2015) who also provide a survey.

[3]  By common uncertainty we mean that each country faces the same *ex-ante* distribution of possible damage costs, and when uncertainty is fully resolved they face the same *ex-post* level of damage costs, i.e. the risks they face are fully correlated across countries. Kolstad and Ulph (2011) extend this model to consider the case where the risks each country faces are uncorrelated. Uncorrelated uncertainty is also considered in a slightly different model in Finus and Pintassilgo (2013) and empirically investigated in a climate model with twelve world regions in Dellink and Finus (2012).

and so needed to wait until they had more information about the risk of climate change before committing to significant action to tackle climate change.

The second strand of literature studies uncertainty and IEA formation with the focus on the role of *risk aversion*, though under the assumption of No Learning. Endres and Ohl (2003) show in a simple two-player prisoners' dilemma, using the mean-standard deviation approach to capture risk aversion, that risk aversion can increase the prospects of cooperation once it reaches a certain threshold. The reason is that the benefits of mutual cooperation increase relative to the payoffs of unilateral cooperation and no cooperation because cooperation reduces the variance of payoffs. The more risk averse players are, the more attractive cooperation becomes compared to free-riding. In their model, there is a first threshold above which the prisoners' dilemma turns into a chicken game and a second threshold above which the game turns into an assurance game. Compared to their paper, we allow for an arbitrary number of players, model cooperation as a two-stage coalition formation game and consider explicitly the role of learning.

Bramoullé and Treich (2009) consider risk-averse players in a global emission model, in which all players behave non-cooperatively as singletons. They show that equilibrium emissions are lower under uncertainty than under certainty, as part of a hedging strategy, but the effect on global welfare is ambiguous. The authors also find that emissions decrease with the level of risk aversion. Unlike our paper, Bramoullé and Treich are not concerned with learning and coalition formation.

Boucher and Bramoullé (2010) consider the effects of risk aversion on coalition formation, but only with No Learning. They analyze the formation of an international environmental treaty using a similar coalition game and payoff function as adopted in this paper. Using an expected utility approach, their analysis focuses on the effect of uncertainty and risk aversion on signatories' efforts, the participation level in an agreement and expected aggregate utility. They show that if additional abatement reduces the variance of countries' payoffs, then, under risk aversion, an increase in uncertainty tends to increase abatement levels and may decrease equilibrium IEA membership while the reverse is true if additional abatement increases the variance of countries' payoffs.[4] In this paper, our model of No Learning satisfies the first condition, but we extend the analysis of Boucher and Bramoullé (2010) by considering also Partial Learning and Full Learning.

Thus, taken together, in our paper, we generalize the analysis of Kolstad and Ulph (2008) by allowing for risk aversion, and the analysis of Boucher and Bramoullé (2010) and Endres and Ohl (2003) by considering the role of learning. The key findings are that as countries become more risk averse it is no longer the case that for most parameter values the scenario of No Learning yields the highest expected aggregate utility, but increasingly it is the scenario Full Learning. Moreover, the set of parameter values for which the scenario Partial Learning yields the highest expected aggregate utility, which is a small subset of such values when countries are risk neutral, becomes slightly larger as countries become more risk averse. Thus, we qualify the negative conclusion about the role of learning in a strategic context if players are sufficiently risk averse.

As in Kolstad and Ulph (2008), our model consists of a two-stage game of IEA formation. The damage cost from global emissions is uncertain and three learning scenarios are considered. Risk aversion is incorporated into the model, through risk preferences and the expected utility framework. In our model, emissions last for just one period, which may

---

[4] Hong and Karp (2014) show that it does not matter whether one analyses the provision of a public good or the amelioration of a public bad. What matters is whether players' actions increase or decrease the volatility of payoffs. In our model, as in Endres and Ohl (2003) and the emission game in Boucher and Bramoullé (2010), abatement (emissions) reduces (increase) the volatility of payoffs.

seem restrictive in the context of dynamic environmental problems such as climate change. However, it has been shown in the literature on IEA formation under uncertainty that one period models produce similar results to multi-period models. For instance, Kolstad and Ulph (2008), using a one-period model, and Ulph (2004), using a two-period model, found similar results regarding the implications of different scenarios of learning for the size of an IEA and expected welfare.[5]

Taken together, the tension we seek to capture in our modeling is between No Learning, where countries decide to join an IEA and base their decisions for ever on expected damage costs ignoring any later information, Full Learning where countries delay making any decision to join an agreement until (almost) all uncertainty about damage costs has been resolved, or Partial Learning where countries start the process of joining an IEA knowing that as they get better information they will be able to use that to adjust their emissions policies. This would seem to be particularly relevant to the kind of situation to which Weitzman (2009) has drawn attention—a small probability of catastrophic climate change—but also generally, given the continuous updates in recent years about expected damages from climate change.

The paper proceeds as follows. In Sect. 2, we set out the theoretical model and present our general results in Sect. 3. Section 4 presents some additional results based on simulations, while Sect. 5 summarizes our main conclusions and implications for future research.

## 2 The Model

### 2.1 No Uncertainty

To establish the basic framework, we set out the model with no uncertainty. There are $N$ identical countries, indexed $i = 1, \ldots, N$. Country $i$ produces emissions $x_i$ with aggregate emissions denoted by $X = \sum_{i=1}^{N} x_i$. Aggregate emissions cause global environmental damages. The cost of environmental damages per unit of global emissions is $\gamma$ and the benefit per unit of individual emissions is normalized to 1. (Thus, $\gamma$ essentially measures the cost–benefit ratio.) The payoff to country $i$, as a function of own and aggregate emissions, is given by

$$\pi_i(x_i, X) \equiv \pi_0 + x_i - \gamma X. \tag{1}$$

with $\pi_0$ a positive constant. In this simple model with a linear payoff function, following the literature, the (continuous) strategy space can be normalized to $x_i \in [0, 1]$. If benefits from emissions are lower than damage costs then the equilibrium emissions are zero. If benefits exceed costs, equilibrium emissions are set at their upper bound, which is normalised to 1. To make this model interesting, we make the following assumption:

**Assumption 1** (i) $1/N < \gamma < 1$; (ii) $\pi_0 > \gamma(N - 1)$.

The individual benefit exceeds the individual unit damage cost from pollution, i.e. $1 > \gamma$ (so countries pollute in the Nash equilibrium) but does not exceed the global unit damage cost, i.e. $1 < \gamma N$ (so countries abate in the social optimum); the second condition is a sufficient condition to ensure that $\pi_i(.) > 0$, which we will need when we consider expected utility.[6]

---

[5] In many papers with a dynamic payoff structure but fixed membership, results are qualitatively similar to the one period emission game (e.g. Rubio and Casino 2005). The extension to flexible membership would be more interesting but is technically very challenging. See Rubio and Ulph (2007).

[6] The lowest possible payoff arises when a country completely abates and all other countries emit at the maximum level, in which case $\pi_i = \pi_0 - \gamma(N - 1)$; so $\pi_0 > \gamma(N - 1) \Rightarrow \pi_i(.) > 0$.

In order to study coalition formation, we employ the widely used two-stage model of IEA formation (Carraro and Siniscalco [1993]; Barrett [1994]) which is solved backwards. In stage 2, the *emission game*, for any arbitrary number of IEA members $n$, $1 \leq n \leq N$, the members of the IEA and the remaining countries set their emission levels as the outcome of a Nash game between the coalition and the fringe countries.[7] That is, the coalition members together maximize the aggregate payoff to their coalition, whereas fringe countries maximize their own individual payoff. Hereafter, symbols $c$ and $f$ will be used to denote coalition countries and fringe countries, respectively. Given $1 > \gamma$, $x_f = 1$ follows; coalition members chose $x_c = 0$ if $1 \leq \gamma n$, and so $\pi_f(n) = \pi_0 + 1 - \gamma(N - n)$ and $\pi_c(n) = \pi_0 - \gamma(N - n)$; however if $1 > \gamma n$, then coalition members will also pollute, $x_c = 1$ and so $\pi_c(n) = \pi_f(n) = \pi_0 + 1 - \gamma N$.[8]

Knowing the payoffs to coalition and fringe countries for any arbitrary number of IEA members, we then determine the stable (Nash) equilibrium in stage 1, the *membership game*. No member should have an incentive to leave the coalition, that is, the coalition is internally stable, $\pi_c(n) \geq \pi_f(n-1)$, and no fringe country should have an incentive to join the coalition, that is, the coalition is externally stable, $\pi_f(n) > \pi_c(n+1)$.[9] It is now easy to show that the stable IEA is of size $n^*(\gamma) = I(1/\gamma)$, which is the smallest integer no less than $1/\gamma$. Consider internal stability and consider the non-trivial situation where $n$ members do not pollute because $1 \leq \gamma n$. If after one member left the coalition, the remaining coalition members continued not to pollute, that is $1 \leq \gamma(n-1)$ was satisfied, then the gain from leaving would be positive: the additional benefit is 1 and the additional damage is $\gamma$, with $1 > \gamma$ by Assumption 1. Thus, a coalition of $n$ members can only be stable if and only if $1 > \gamma(n-1)$ is true after one member left as the remaining coalition members would switch from $x_c(n) = 0$ to $x_c(n-1) = 1$. Then the additional benefit from pollution of 1 falls short of the additional damage $\gamma n$, as by assumption $1 \leq \gamma n$ in the initial situation with $n$ members. It is easily checked that such an equilibrium is also externally stable. The aggregate payoff, that is, the sum of the payoffs over all countries, when a stable coalition of size $n^*(\gamma)$ forms is given by:

$$\begin{aligned} \Pi(\gamma) &\equiv n^*(\gamma)\pi_c(n^*(\gamma)) + (N - n^*(\gamma))\pi_f(n^*(\gamma)) \\ &= N\pi_0 + (N - n^*(\gamma))(1 - \gamma N) \end{aligned}. \tag{2}$$

Thus, this simple model provides a relationship between the unit damage cost $\gamma$ and the equilibrium number of coalition members. The equilibrium is a knife-edge equilibrium with $n^*(\gamma)$ countries forming the coalition, which de facto dissolves once a member leaves the coalition as no country would abate anymore. The equilibrium coalition size weakly decreases in the cost–benefit ratio from emissions $\gamma$—the larger is $\gamma$ the smaller is the number of countries in a stable IEA.

---

[7] A sequential Stackelberg game in the second stage, as an alternative assumption (e.g. Barrett [1994]), would make no difference here as players have dominant strategies. This also applies to Boucher and Bramoullé ([2010]).

[8] It is now evident why we need Assumption 1: it avoids trivial outcomes where all countries either abate or pollute in the Nash equilibrium and the social optimum (and hence also for partial cooperation).

[9] Without loss of generality, the strong inequality for external stability could be replaced by a weak inequality sign. Our assumption avoids knife-edge cases where a fringe country is indifferent between staying outside and joining a coalition.

## 2.2 Uncertainty, Risk Aversion and Learning

Now assume that the unit damage cost of global emissions is uncertain and equal for all countries, both ex-ante and ex-post. We denote the value by $\gamma_s$ in the state of the world $s$ and hence (1) becomes:

$$\pi_{i,s}(x_i, X) \equiv \pi_0 + x_i - \gamma_s X. \tag{3}$$

Following Kolstad and Ulph (2008) and Boucher and Bramoullé (2010), we assume for simplicity that $\gamma_s$ can take one of two values: low damage costs, $\gamma_l$, with probability $p$, and high damage costs, $\gamma_h$, with probability $1 - p$, where $\gamma_l < \gamma_h$ and $0 < p < 1$. We denote by $\bar{\gamma} \equiv p\gamma_l + (1 - p)\gamma_h$ the expected value of unit damage costs and by $\bar{\bar{\gamma}} = 0.5(\gamma_l + \gamma_h)$ the median value of damage costs. We define $n_l \equiv I(1/\gamma_l)$, $n_h \equiv I(1/\gamma_h)$, $\bar{n} \equiv I(1/\bar{\gamma})$, $\bar{\bar{n}} \equiv I(1/\bar{\bar{\gamma}})$ and make the following assumptions:

**Assumption 2** (i) $1/N < \gamma_l < \gamma_h < 1$; (ii) $\pi_0 > \gamma_h(N - 1)$; (iii) $2 \leq n_h < \bar{\bar{n}} - 1 < \bar{\bar{n}} + 1 < n_l \leq N$.

Assumptions 2(i) and 2(ii) are just the analogues to Assumption 1(i) and 1(ii) when there is uncertainty. Assumption 2(iii) means that uncertainty matters, in the sense that it implies significant differences in the size of the stable IEAs that would arise if we knew for certain which state of the world prevailed.

To allow for risk aversion, we assume that each country has an identical utility function over payoffs: $u(\pi_i)$, $u'(\pi_i) > 0$, $u''(\pi_i) < 0$. While ex-ante countries face uncertainty about the true value of unit damage costs, we want to allow for the possibility that countries may learn information during the course of the game which changes the risk they face. We shall follow Kolstad and Ulph (2008) in considering three very simple scenarios of learning, denoted by $m$, $m \in \{NL, PL, FL\}$. With No Learning ($m = NL$) countries make their decisions about membership and emissions with uncertainty about the true value of unit damage costs. With Full Learning ($m = FL$) countries learn the true value of unit damage costs before they have to take their decisions on membership (stage 1) and emissions (stage 2). With Partial Learning ($m = PL$) countries learn the true value of damage costs at the end of stage 1, that is after they have made their membership decisions but before they make their emission decisions (stage 2). Thus, in this simple analysis, learning takes the form of revealing perfect information.[10] We will compare the outcomes of the three scenarios of learning in terms of the expected size of IEAs and expected aggregate utility from an ex-ante perspective, i.e. before stage 1.

## 3 Analytical Results

In this section, we set out the equilibrium of the IEA model for each of the three models of learning with risk aversion, generalizing the results of Kolstad and Ulph (2008) who assumed risk neutrality. The proofs are provided in "Appendix A1".

---

[10] We use the term "Partial Learning" in line with the IEA-literature, although this may be misleading; partial learning is usually referred to as delayed but perfect learning, so not partial learning in the sense of Bayesian updating.

### 3.1 Full Learning

We start with the benchmark scenario of Full Learning (FL). Players know the realization of the damage parameter $\gamma$ at the outset of the coalition formation game, i.e. before stage 1. Thus, the results follow directly from what we know from the game with certainty in Sect. 2.1 above.

**Proposition 1: Full Learning** *If state $s = l, \ h$ has been revealed before stage 1, then the subsequent membership and emission decisions are as in the model with certainty with damage cost parameter $\gamma_s$. In each state $s = l, \ h$ the size of the stable IEA is $n_s = I(1/\gamma_s)$ and the expected membership is $E(n^{FL}) = pn_l + (1-p)n_h$. The expected aggregate utility is given by:*

$$E\left(U^{FL}\right) = pn_l u\left(\pi_{c,l}^{FL}\right) + p(N - n_l)u\left(\pi_{f,l}^{FL}\right) + (1-p)n_h u\left(\pi_{c,h}^{FL}\right) + (1-p)(N - n_h)u\left(\pi_{f,h}^{FL}\right)$$

*with $\pi_{c,l}^{FL} = \pi_0 - \gamma_l(N - n_l), \ \pi_{f,l}^{FL} = \pi_0 + 1 - \gamma_l(N - n_l), \ \pi_{c,h}^{FL} = \pi_0 - \gamma_h(N - n_h)$ and $\pi_{f,h}^{FL} = \pi_0 + 1 - \gamma_h(N - n_h)$.*

Note that with Full Learning, while the degree of risk aversion does not affect the expected size of the IEA it will affect expected utility. Importantly, the expected size and expected aggregate utility are computed from an ex-ante perspective to make a comparison with the other models of learning meaningful.

### 3.2 No Learning

In this section, we address the scenario of No Learning in which players take their membership (stage 1) and emission (stage 2) decisions under uncertainty.[11] We begin by solving for optimal emissions of countries for any number of IEA members $n$. Since the benefit of one unit of emissions exceeds the damage cost in both states of the world, it is straightforward to see that fringe countries will always pollute. To solve for the optimal emissions for a coalition member for any $n$, which we denote by $x_c(n)$, we need to introduce some notation. We define:

$$E(u_c(x_c(n), n))) \equiv pu(\pi_{c,l}(x_c(n), n)) + (1 - p)u(\pi_{c,h}(x_c(n), n)) \tag{4a}$$

where

$$\pi_{c,s}(x_c(n), n) = \pi_0 - \gamma_s(N - n) + x_c(n)(1 - \gamma_s n), \ s = l, h. \tag{4b}$$

$E(u_c(x_c(n), n)))$ is the expected utility to an IEA country when there are $n$ IEA members who set emissions $x_c$ and all fringe countries set emissions equal to 1. Then:

$$\frac{\partial E(u_c)}{\partial x_c} = p(1 - \gamma_l n)u'[\pi_{c,l}(x_c, n)] + (1 - p)(1 - \gamma_h n)u'[\pi_{c,h}(x_c, n)], \tag{5a}$$

$$\frac{\partial^2 E(u_c)}{\partial (x_c)^2} = p(1 - \gamma_l n)^2 u''[\pi_{c,l}(x_c(n), n)] + (1 - p)(1 - \gamma_h n)^2 u''[\pi_{c,h}(x_c(n), n)] < 0. \tag{5b}$$

From (5a) it is clear that if $\gamma_l n \geq 1$, then $\frac{\partial E(u_c)}{\partial x_c} < 0$ and hence it is optimal for IEA countries to completely abate, while if $\gamma_h n \leq 1$, then $\frac{\partial E(u_c)}{\partial x_c} > 0$ and hence it is optimal for IEA countries to completely pollute. To get tighter bounds on when IEA countries completely pollute or abate, we define $\tilde{n}$ as the largest value of $n$ such that:

---

[11] Our analysis of No Learning is similar to the analysis provided by Boucher and Bramoullé (2010).

$$p(1 - \gamma_l n)u'[\pi_{c,l}(1, n)] + (1 - p)(1 - \gamma_h n)u'[\pi_{c,h}(1, n)] > 0 \quad \forall n < \tilde{n}. \qquad (6a)$$

and $\tilde{\tilde{n}}$ as the smallest value of $n$ such that:

$$p(1 - \gamma_l n)u'[\pi_{c,l}(0, n)] + (1 - p)(1 - \gamma_h n)u'[\pi_{c,h}(0, n)] \leq 0 \quad \forall n \geq \tilde{\tilde{n}}. \qquad (6b)$$

We summarise the results on emissions in the following Lemma:

**Lemma 1: Emission Decisions with No Learning**

   (i)   $x_f(n) = 1 \, \forall n, \, 2 \leq n \leq N$;
  (ii)   $n_h \leq \tilde{n} \leq \tilde{\tilde{n}} \leq \bar{n}$;
 (iii)   $p \to 0 \Rightarrow \tilde{n}, \tilde{\tilde{n}}, \bar{n} \to n_h$;   $p \to 1 \Rightarrow \tilde{n}, \tilde{\tilde{n}}, \bar{n} \to n_l$;
  (iv)   $n < \tilde{n} \Rightarrow x_c(n) = 1$;
   (v)   $n = \tilde{n} \Rightarrow 0 \leq x^c(n) \leq 1$;
  (vi)   $\tilde{n} < n < \tilde{\tilde{n}} \Rightarrow 0 \leq x^c(n) < 1$;
 (vii)   $n \geq \tilde{\tilde{n}} \Rightarrow x_c(n) = 0$.

As already noted, for any size of an IEA, $n$, fringe countries always choose the upper limit of emissions. For IEA members, there is a critical range of values for $n$, $[\tilde{n}, \tilde{\tilde{n}}]$, which lies between $n_h$ and $\bar{n}$ such that IEA members will completely abate if $n \geq \tilde{\tilde{n}}$, and choose the upper limit of emissions for $n < \tilde{n}$; but if there are values of $n$ which lie within the range $]\tilde{n}, \tilde{\tilde{n}}[$, then coalition members choose a level of emissions below the upper limit. Note that, as in Boucher and Bramoullé (2010), with both risk neutrality and risk aversion $x_c(n) = 0$ for $n \geq \bar{n}$ and $x_c(n) = 1$ for $n < \tilde{n} \leq \bar{n}$. With risk neutrality for both $\tilde{n} \leq n < \tilde{\tilde{n}}$ and $\tilde{\tilde{n}} \leq n < \bar{n}$ the emissions are $x_c(n) = 1$; while with risk aversion $0 \leq x_c(n) \leq 1$ for $\tilde{n} \leq n < \tilde{\tilde{n}}$ and $x_c(n) = 0$ for $\tilde{\tilde{n}} \leq n < \bar{n}$. So for $\tilde{\tilde{n}} \leq n < \bar{n}$ aggregate emissions are lower with risk aversion than with risk neutrality.

**Proposition 2: No Learning** *With No Learning, for all parameter values, there exists a stable IEA with membership $n^{NL}$, which is the same in both states of the world with $n^{NL} \in [\tilde{n}, \tilde{\tilde{n}}]$. This is also expected membership, i.e. $E(n^{NL}) = n^{NL} \in [\tilde{n}, \tilde{\tilde{n}}]$ which is (weakly) lower than under risk neutrality with $E(n^{NL}) = \bar{n}$ as $\tilde{\tilde{n}} \leq \bar{n}$. Emissions of fringe and IEA members are given in Lemma* 1. *Expected aggregate utility is given by:*

$$E\left(U^{NL}\right) = pn^{NL}u\left(\pi_{c,l}^{NL}\right) + p\left(N - n^{NL}\right)u\left(\pi_{f,l}^{NL}\right)$$
$$+ (1 - p)n^{NL}u\left(\pi_{c,h}^{NL}\right) + (1 - p)\left(N - n^{NL}\right)u\left(\pi_{f,h}^{NL}\right)$$

with $\pi_{c,l}^{NL} = \pi_0 - \gamma_l(N - n^{NL}) + (1 - n^{NL}\gamma_l)x^c(n^{NL})$, $\pi_{f,l}^{NL} = \pi_0 - \gamma_l(N - n^{NL}) + 1 - n^{NL}\gamma_l x^c(n^{NL})$, $\pi_{c,h}^{NL} = \pi_0 - \gamma_h(N - n^{NL}) + (1 - n^{NL}\gamma_h)x^c(n^{NL})$ and $\pi_{f,h}^{NL} = \pi_0 - \gamma_h(N - n^{NL}) + 1 - n^{NL}\gamma_h x^c(n^{NL})$.

So the expected equilibrium coalition size is (weakly) smaller under risk aversion than risk neutrality. With uncertainty, countries are unsure about the state of the world. With risk aversion, and hence concave utility, countries shy away from the commitment to be a member in an IEA. This is in line with the findings in Boucher and Bramoullé (2010).

## 3.3 Partial Learning

In the scenario of Partial Learning, countries have to make their decision on whether to join an IEA without knowing the true damage cost of emissions, but can make their subsequent emission decisions based on that knowledge. It follows that the emission decisions of countries do not depend on risk aversion and so are the same as in Kolstad and Ulph (2008). Since for one unit of emissions the benefit exceeds damage costs in each state of the world, a fringe country will optimally set $x_{f,s} = 1, \ s = l, h$; for an IEA member optimal emissions depend on the size of the IEA $n$; so $n \geq n_l \ \Rightarrow \ x_{c,s}(n) = 0, \ s = l, h; n_h \leq n < n_l \ \Rightarrow \ x_{c,l}(n) = 1$ and $x_{c,h}(n) = 0; n < n_h \ \Rightarrow \ x_{c,s}(n) = 1, s = l, h$. That is, fringe countries always pollute; if there are at least $n_l$ IEA members, then IEA members always abate; if there are less than $n_h$ IEA members, then IEA members always pollute; otherwise IEA members pollute in the low damage cost state and abate in the high damage cost state.

As in Kolstad and Ulph (2008), for certain values of $p$ there may be more than one stable IEA with Partial Learning. In our model a second stable IEA exists iff $\tilde{p} \leq p \leq 1$ where $\tilde{p} \equiv \frac{\chi}{1+\chi}$ with $\chi \equiv \frac{u[\pi_0 - (N-n_l)\gamma_h + (1-\gamma_h)] - u[\pi_0 - (N-n_l)\gamma_h]}{u[\pi_0 - (N-n_l)\gamma_l] - u[\pi_0 - N\gamma_l + 1]} > 0$. It is straightforward to show that when countries are risk neutral $\tilde{p} = \frac{1-\gamma_h}{n_l\gamma_l - \gamma_h}$, as in Kolstad and Ulph (2008). Irrespective of the degree of risk aversion, it is straightforward to see that $\gamma_h \to 1 \Rightarrow \tilde{p} \to 0$. However, we have not been able to determine analytically how $\tilde{p}$ varies with the degree of risk aversion, as this depends on the exact form of the utility function. In Sect. 4 we report our findings on this from our simulation results.

**Proposition 3: Partial Learning** *With Partial Learning, for all parameter values there always exists a stable IEA with $n^{PL_1} = n_h$ members. All countries pollute in the low damage cost state, while in the high damage cost state coalition members abate and fringe countries pollute. Expected aggregate utility is given by:*

$$E\left(U^{PL_1}\right) = pn_h u\left(\pi_{c,l}^{PL_1}\right) + p(N - n_h)u\left(\pi_{f,l}^{PL_1}\right)$$
$$+ (1-p)n_h u\left(\pi_{c,h}^{PL_1}\right) + (1-p)(N - n_h)u\left(\pi_{f,h}^{PL_1}\right).$$

*with* $\pi_{c,l}^{PL_1} = \pi_0 + 1 - \gamma_l N, \ \pi_{f,l}^{PL_1} = \pi_0 + 1 - \gamma_l N, \ \pi_{c,h}^{PL_1} = \pi_0 - \gamma_h(N - n_h)$ *and* $\pi_{f,h}^{PL_1} = \pi_0 + 1 - \gamma_h(N - n_h)$.

*If $\tilde{p} \leq p \leq 1$, then there is a second stable IEA with $n^{PL_2} = n_l$ members. In both states of the world, coalition members abate and fringe countries pollute. Expected aggregate utility is given by:*

$$E\left(U^{PL_2}\right) = pn_l u\left(\pi_{c,l}^{PL_2}\right) + p(N - n_l)u\left(\pi_{f,l}^{PL_2}\right)$$
$$+ (1-p)n_l u\left(\pi_{c,h}^{PL_2}\right) + (1-p)(N - n_l)u\left(\pi_{f,h}^{PL_2}\right)$$

*with* $\pi_{c,l}^{PL_2} = \pi_0 - \gamma_l(N - n_l), \ \pi_{f,l}^{PL_2} = \pi_0 + 1 - \gamma_l(N - n_l), \ \pi_{c,h}^{PL_2} = \pi_0 - \gamma_h(N - n_l)$ *and* $\pi_{f,h}^{PL_2} = \pi_0 + 1 - \gamma_h(N - n_l)$.

Since the second equilibrium Pareto-dominates the first equilibrium if it exists, expected membership is either $E(n^{PL_1}) = n_h$ if $p < \tilde{p}$ or $E(n^{PL_2}) = n_l$ if $\tilde{p} \leq p \leq 1$.

As the degree of risk aversion affects $\tilde{p}$ it has an effect on the likelihood of a second coalition with higher membership $n_l$ being stable. This effect is further explored in Sect. 4.

### 3.4 Comparison Across the Three Scenarios of Learning

This section presents analytical results on the ranking of the learning scenarios in terms of expected membership and expected utility. General results are exposed as well as results on special cases with limit parameter values.

#### 3.4.1 General Results

In this sub-section, we investigate what we can say generically about expected IEA membership, payoffs and expected utility across the four possible equilibria, FL, NL, $PL_1$ and $PL_2$.

In terms of expected membership of an IEA it is clear that since $E(n^{PL_1}) = n_h$, this equilibrium has the lowest expected membership while since $E(n^{PL_2}) = n_l$, this equilibrium has the highest expected membership. Note also that:

$$E\left(n^{FL}\right) \geq p\left(\frac{1}{\gamma_l}\right) + (1-p)\left(\frac{1}{\gamma_h}\right) = \frac{p\gamma_h + (1-p)\gamma_l}{\gamma_l \gamma_h}.$$

Moreover it is straightforward to show that:

$$\frac{p\gamma_h + (1-p)\gamma_l}{\gamma_l \gamma_h} \geq \frac{1}{p\gamma_l + (1-p)\gamma_h} \Leftrightarrow p(1-p)(\gamma_h - \gamma_l)^2 \geq 0.$$

Hence:

$$E(n^{FL}) > \frac{1}{p\gamma_l + (1-p)\gamma_h}. \tag{7a}$$

From Proposition 2 and Lemma 1 we have that:

$$E\left(n^{NL}\right) \leq \tilde{\tilde{n}} \leq \bar{n} = I\left(\frac{1}{p\gamma_l + (1-p)\gamma_h}\right). \tag{7b}$$

Because $E(n^{FL}) = pI(\frac{1}{\gamma_l}) + (1-p)I(\frac{1}{\gamma_h})$, (7a) and (7b) are not sufficient to ensure that $E(n^{NL}) \leq E(n^{FL})$. Moreover, as $0 < p < 1$, then $E(n^{PL_1}) < E(n^{FL}) < E(n^{PL_2})$. Thus, as we shall see, the following two rankings can occur, according to the parameter values:

$$E\left(n^{PL_1}\right) \leq E\left(n^{NL}\right) \leq E\left(n^{FL}\right) < E\left(n^{PL_2}\right). \tag{8a}$$

$$E\left(n^{PL_1}\right) < E\left(n^{FL}\right) < E\left(n^{NL}\right) \leq E\left(n^{PL_2}\right). \tag{8b}$$

In terms of payoffs across the four equilibria, it is straightforward to see from Propositions 1, 2 and 3 that:

$$\pi_{c,l}^{PL_2} = \pi_{c,l}^{FL} \geq \pi_{c,l}^{PL_1}; \quad \pi_{f,l}^{PL_2} = \pi_{f,l}^{FL} > \pi_{f,l}^{PL_1}; \quad \pi_{c,h}^{PL_2} > \pi_{c,h}^{FL} = \pi_{c,h}^{PL_1};$$
$$\pi_{f,h}^{PL_2} > \pi_{f,h}^{FL} = \pi_{f,h}^{PL_1}. \tag{9a}$$

For NL, in the low damage cost state of the world, the highest payoff to coalition members is when $x_c = 1$, which is less than or equal to the payoff to coalition members in $PL_2$ since $n_l \gamma_l \geq 1$; in the high damage cost state of the world, the highest payoff to coalition members is when $x_c = 0$, which is less than the payoff to members in $PL_2$ since $E(n^{NL})\gamma_h < n_l \gamma_h$. So it must be the case that:

$$\pi_{c,l}^{PL_2} \geq \pi_{c,l}^{NL}; \quad \pi_{f,l}^{PL_2} > \pi_{f,l}^{NL}; \quad \pi_{c,h}^{PL_2} > \pi_{c,h}^{NL}; \quad \pi_{f,h}^{PL_2} > \pi_{f,h}^{NL}. \tag{9b}$$

The results in (9a) and (9b) allow us to rank a number of the payoffs across the four possible equilibria for both members and fringe countries in the high and low damage cost states of the world. However, this is not sufficient to allow us to rank expected aggregate utility across different models of learning at an analytical and general level because (i) the weights of the different utilities in the aggregation depend on the equilibrium coalition size, which differs across the learning scenarios; and (ii) how differences in payoffs translate into differences in expected aggregate utility depend on the exact form of the utility function. The next section reports the simulations we have carried out to compare expected IEA membership and expected welfare across the different models of learning.

The analytic results obtained allow us, however, to get some insights about the role of risk aversion on expected IEA membership and welfare. Comparing to the case of risk neutrality explored by Kolstad and Ulph (2008), we conclude that risk aversion does not affect the expected coalition under FL but (weakly) decreases it under NL. Regarding PL it does not affect the two possible coalition sizes under PL, $n^{PL_1} = n_h$ and $n^{PL_2} = n_l$, although it affects the likelihood of the larger equilibrium. Hence, if we restrict our analysis to the extreme cases of FL and NL, we can conclude that the prospects of learning being conducive to larger IEAs are higher under risk aversion. Thus, the result found by Kolstad and Ulph (2008) that expected welfare is higher under NL compared to FL is less likely under risk aversion. Adding PL to this picture, the novelty is that the likelihood of the larger equilibrium depends on the level of risk aversion.

### 3.4.2 Special Cases

In this subsection, the comparison of learning scenarios in terms of expected IEA membership and expected utility is undertaken for special parameter values. Following Karp (2012), we use two limit cases regarding the probability of low damages: $p = \varepsilon \approx 0$ and $p = 1 - \varepsilon \approx 1$, where $\varepsilon$ denotes an infinitesimal. Limit values for the damage cost from pollution were also considered: $\gamma_l = 1/N + \varepsilon \approx 1/N$ and $\gamma_h = 1 - \varepsilon \approx 1$.

Lemma 2 shows the results that could be obtained analytically on the ranking of learning scenarios, in terms of expected membership.

### Lemma 2: Expected Membership

(i)   If $p \approx 0$ then $E(n^{FL}) > E(n^{PL}) = E(n^{NL})$;
(ii)  If $p \approx 1$ then $E(n^{PL}) > E(n^{NL}) = E(n^{FL})$;
(iii) If $\gamma_l \approx 1/N$ then $E(n^{FL}) > n^{PL} = n_h$;
(iv)  If $\gamma_h \approx 1$ then $n^{PL} = n_l > E(n^{FL})$.

When high damages are very likely, $p \approx 0$, then FL leads to larger expected IEA membership than PL and NL, irrespective of the level of risk aversion. Hence, in this context, learning is conducive to the formation of larger agreements. The opposite holds when low damages are very likely, $p \approx 1$, with PL and NL leading to larger expected IEA membership.

In the context of risk neutrality, Karp (2012) found that $p \approx 0$ implies $E(n^{FL}) > E(n^{NL})$ and $p \approx 1$ implies $E(n^{NL}) > E(n^{FL})$. Thus, we show that Karp's result also holds under risk aversion and we have extended it to the scenario of PL.

If the low damage cost from pollution is close to its lower bound, $\gamma_l \approx 1/N$, then under PL the smaller equilibrium coalition size forms, $n^{PL_1} = n_h$, and hence the expected size of an IEA under FL is higher than under PL. If the high damage cost is at its upper bound, $\gamma_h \approx 1$, then the larger IEA forms under PL, $n^{PL_2} = n_l$, which exceeds the expected coalition size

under FL. For these two limit values of the damage cost, analytical results could not be found on the rankings between FL and NL, as well PL and NL.

The results of the special cases in terms of expected utility are shown in Lemma 3.

### Lemma 3: Expected Utility

  (i)  *If $p \approx 0$ then $E(U^{FL}) > E(U^{PL})$;*
 (ii)  *If $p \approx 1$ then $E(U^{PL}) = E(U^{NL})$;*
(iii)  *If $\gamma_l \approx 1/N$ then $E(U^{FL}) > E(U^{PL})$;*
 (iv)  *If $\gamma_h \approx 1$ then $E(U^{PL}) > E(U^{FL})$ and $E(U^{PL}) > E(U^{NL})$.*

Regarding expected utility, for each special case we could obtain only pairwise rankings of the three learning scenarios, and not a complete ranking. When high damages are very likely, $p \approx 0$, then FL yields higher expected utility than PL. When these damages have a low likelihood, $p \approx 1$, then PL and NL lead to same expected utility. When low damage cost from pollution is close to its lower bound, $\gamma_l \approx 1/N$, FL provides a larger expected utility than PL. When $\gamma_h \approx 1$, then as we showed in Sect. 3.3, $\tilde{p} \approx 0$, and PL$_2$ will be the selected equilibrium for Partial Learning, so expected utility be higher than with either FL or NL.

Combining the results of expected membership size and expected utility a few conclusions can be made. First, if $p \approx 0$ then learning leads to a better outcome in terms of membership, $E(n^{FL}) > E(n^{PL}) = E(n^{NL})$. Regarding utility, the only analytical result obtained also points in that direction, $E(U^{FL}) > E(U^{PL})$. Second, for $p \approx 1$, learning leads to worse results in terms of membership, $E(n^{PL}) = E(n^{NL}) > E(n^{FL})$. The result obtained on utility, $E(U^{NL}) = E(U^{PL})$, points to a neutrality of learning. Third, for $\gamma_l \approx 1/N$ the results obtained indicate an advantage of learning in terms of membership, $E(n^{FL}) > n^{PL} = n_h$, and utility $E(U^{FL}) > E(U^{PL})$. The opposite occurs for $\gamma_h \approx 1$. Using these four limit cases, we can conclude that under risk aversion the role of learning in terms of membership and utility depends on the distribution of the damage cost, namely parameters, $\gamma_l$, $\gamma_h$, $p$. This is in line with the results obtained by Karp (2012) under risk neutrality. In the next section we comment on the implications for these limiting results arising from our simulations.

## 4 Results from Simulations

There are three sets of issues we were unable to resolve analytically in Sect. 3 and which we explore using numerical simulations. (i) What is the expected size of the IEA in the case of No Learning, $E(n^{NL})$, in relation to the theoretical limits $\tilde{n}$ and $\tilde{\tilde{n}}$ and, more importantly, to the key parameters of our model, $n_h$ and $\bar{n}$ and how does this vary across different degrees of risk aversion? (ii) In the case of Partial Learning what is the critical value of the likelihood of low damage state of the world $\tilde{p}$ such that for $\tilde{p} \leq p \leq 1$ there is second stable IEA ($n^{PL_2} = n_l$) and how does $\tilde{p}$ vary across different degrees of risk aversion? (iii) Most importantly, how does the expected size of IEA and expected aggregate utility compare across the three different models of learning, Full Learning (FL), No Learning (NL), and Partial Learning (PL$_1$, PL$_2$) and how does this comparison depend on the degree of risk aversion and other parameters of the model?

We now discuss our choice of parameter values for the simulations. To guarantee non-negative payoffs we set:

$$\pi_0 = B\gamma_h(N-1), \quad B > 1 \tag{10}$$

and to ensure that payoffs are sensitive to countries' abatement decisions we chose $B = 1.1$.[12] For the parameter of relative risk aversion, $\rho$, in the CRRA utility function $u(\pi_i) = [1/(1-\rho)]\pi_i^{1-\rho}$ we use $\rho = 0$ (risk-neutral) as a benchmark case and then choose 7 values of $\rho = 0.05, 0.5, 0.99, 2.5, 5.0, 10.0$, and $20.0$ to capture what we believe to be a reasonable range of values for country-level risk aversion.[13]

For the remaining key parameters ($N, p, \gamma_l, \gamma_h$), we report results of 2 sets of simulations. In Sect. 4.1 we present results from a set of simulations using a small number of specific values of the parameters ($N, p, \gamma_l, \gamma_h$) to get some insights into what drives results. Then in Sect. 4.2 we present results of a more general set of 500,000 simulations where the parameters ($N, p, \gamma_l, \gamma_h$) are chosen randomly within a specified range. The details of how these parameters are chosen are set out in the "Appendix A2".

In addition to reporting the membership of the IEA and the expected utility for each of our three models of learning, $m$, we also report expected damage costs as a percentage of GDP,[14] which we define as $D^m$. The reason for doing this is to indicate that the environmental issue has significant implications.

Finally, in Sect. 4.3, we will test the robustness of our simulation results in three respects.[15] First, we check whether the results in Sects. 4.1 and 4.2, derived for $B = 1.1$, hold for higher values of $B$. Second, we will undertake simulations using the limiting conditions adopted in Lemmas 2 and 3, in Sect. 3.4, to check the ranking of the learning scenarios in terms of expected IEA membership and expected utility. Third, in Sects. 4.1 and 4.2 we assume that each country has a constant relative risk aversion (CRRA) utility function, and we will test whether our key numerical results also hold if we use a constant absolute risk aversion (CARA) utility function.

## 4.1 Results from Simple Simulations

In our simple simulations we use $B = 1.1$ and $N = 20$. The results[16] are shown in Table 1. A total of 13 simulations were used. The first part of the table shows the parameter values we used, while the next three sections show the outcomes for NL, FL and PL respectively. In the first 7 columns we use different values of risk aversion $\rho$, but keep constant the values for $p, \gamma_l, \gamma_h$. In the remaining 6 columns we set $\rho = 2.5$, the midpoint of the values we use, and vary in turn $p$, (columns 8 and 9), and then $\gamma_l, \gamma_h$ by choosing higher and lower values of $\bar{\bar{\gamma}}$ (columns 10 and 11) and then tighter and wider spreads of $\gamma_l, \gamma_h$ around $\bar{\bar{\gamma}}$. For these parameter values, expected damage costs as a percentage of GDP ranges between 0.34 and 3.40.

---

[12] We are grateful to a referee for this suggestion.

[13] Meyer and Meyer (2006) note that the CRRA utility function is widely used in empirical studies of risk aversion, and that empirical estimates of $\rho$ vary between 0 and 100. They note that such estimates depend on the variable that enters the utility function, and for the three most commonly used variables—wealth, income and profits—the appropriate empirical estimate increases as one moves from wealth to profits. In our one-period model the relevant variable is income, though there is no distinction between wealth and income. Hence, we have chosen a range of values for $\rho$ at the lower end of the range noted by Meyer and Meyer, namely between 0 and 20.

[14] A referee suggested we calculate this measure to indicate that the environmental problem is a significant one. We define GDP as: $\pi_0 N + (N - n) + nx^c$; each coalition country produces output $\pi_0 + x^c$; each fringe country produces output $\pi_0 + 1$.

[15] We are grateful to two referees for these suggestions.

[16] The key results are not sensitive to the choice of $N$.

**Table 1** Results of simple simulations: $B = 1.1$; $N = 20$

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\rho$ | 00.05 | 00.50 | 00.99 | 02.50 | 05.00 | 10.00 | 20.00 | 2.50 | 2.50 | 2.50 | 2.50 | 2.50 | 2.50 |
| $p$ | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.05 | 0.95 | 0.5 | 0.5 | 0.5 | 0.5 |
| $\pi_0$ | 2.69 | 2.69 | 2.69 | 2.69 | 2.69 | 2.69 | 2.69 | 2.69 | 2.69 | 6.36 | 1.55 | 2.40 | 2.90 |
| $\gamma_l$ | .0615 | .0615 | .0615 | .0615 | .0615 | .0615 | .0615 | .0615 | .0615 | .0732 | .0531 | .0755 | .0519 |
| $\gamma_h$ | .1289 | .1289 | .1289 | .1289 | .1289 | .1289 | .1289 | .1289 | .1289 | .3041 | .0743 | .1150 | .1385 |
| $\bar{\bar{\gamma}}$ | .0952 | .0952 | .0952 | .0952 | .0952 | .0952 | .0952 | .0952 | .0952 | .1887 | .0657 | .0952 | .0952 |
| $\bar{\gamma}$ | .0952 | .0952 | .0952 | .0952 | .0952 | .0952 | .0952 | .1256 | .0649 | .1887 | .0657 | .0952 | .0952 |
| $n_l$ | 17 | 17 | 17 | 17 | 17 | 17 | 17 | 17 | 17 | 14 | 19 | 14 | 20 |
| $n_h$ | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 4 | 14 | 9 | 8 |
| $\bar{\bar{n}}$ | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 6 | 16 | 11 | 11 |
| $\bar{n}$ | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 8 | 15 | 6 | 16 | 11 | 11 |
| *No learning* | | | | | | | | | | | | | |
| $\tilde{n}$ | 11 | 10 | 10 | 9 | 8 | 8 | 8 | 8 | 13 | 4 | 15 | 10 | 8 |
| $\tilde{\tilde{n}}$ | 11 | 11 | 10 | 9 | 9 | 8 | 8 | 8 | 16 | 4 | 16 | 10 | 9 |
| $x_c$ | 0.00 | 0.45 | 0.00 | 0.00 | 0.08 | 0.00 | 0.00 | 0.00 | .846 | 0.00 | .569 | 0.00 | 0.19 |
| $E(\pi^{NL})$ | 11 | 10 | 10 | 9 | 8 | 8 | 8 | 8 | 13 | 4 | 15 | 10 | 8 |
| $u^{NL}_{cl}$ | 42.71 | 53.55 | 53.03 | 70.26 | 90.20 | 99.18 | 99.99 | 67.74 | 83.16 | 93.99 | 42.50 | 56.06 | 77.03 |
| $u^{NL}_{ch}$ | 16.79 | 7.23 | 18.05 | 18.10 | 0.00 | 0.00 | 0.00 | 0.00 | 10.77 | 0.00 | 1.69 | 17.68 | 0.10 |
| $u^{NL}_{fl}$ | 84.69 | 75.13 | 88.23 | 94.15 | 98.75 | 99.98 | 99.99 | 93.25 | 86.65 | 97.97 | 75.35 | 94.46 | 92.71 |
| $u^{NL}_{fh}$ | 59.32 | 34.86 | 66.07 | 78.98 | 91.67 | 99.65 | 99.99 | 74.95 | 29.06 | 59.68 | 55.96 | 81.29 | 64.99 |
| $D^{NL}$ (%) | 1.85 | 2.73 | 2.02 | 2.18 | 2.45 | 2.35 | 2.35 | 2.35 | 3.23 | 3.40 | 2.26 | 1.98 | 2.62 |
| $E(U^{NL}_t)$ | 48.77 | 42.69 | 56.35 | 67.51 | 75.15 | 79.73 | 79.97 | 46.87 | 81.02 | 72.46 | 32.99 | 63.12 | 62.72 |

**Table 1** continued

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Full learning* | | | | | | | | | | | | | |
| $E(n^{FL})$ | 12.5 | 12.5 | 12.5 | 12.5 | 12.5 | 12.5 | 12.5 | 8.4 | 16.5 | 9 | 16.5 | 11.5 | 14 |
| $u_{cl}^{FL}$ | 58.31 | 63.96 | 69.89 | 84.98 | 96.71 | 99.92 | 99.99 | 84.97 | 84.97 | 97.07 | 51.88 | 72.86 | 88.22 |
| $u_{ch}^{FL}$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.87 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 2.96 |
| $u_{fl}^{FL}$ | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |
| $u_{fh}^{FL}$ | 42.99 | 49.13 | 55.91 | 74.96 | 92.97 | 99.65 | 99.99 | 74.95 | 74.95 | 59.68 | 87.78 | 80.43 | 72.58 |
| $D^{FL}(\%)$ | 1.38 | 1.38 | 1.38 | 1.38 | 1.38 | 1.38 | 1.38 | 2.25 | 0.42 | 1.87 | 0.69 | 1.49 | 1.19 |
| $E(U^{FL})$ | 45.18 | 49.42 | 53.98 | 66.11 | 76.69 | 79.86 | 79.99 | 47.08 | 85.11 | 72.84 | 40.37 | 62.63 | 66.47 |
| *Partial learning* | | | | | | | | | | | | | |
| $\tilde{p}$ | .949 | .947 | .945 | .939 | .929 | .917 | .921 | .939 | .939 | .978 | .985 | .926 | .940 |
| PL equil. selected | PL1 | PL1 | PL1 | PL1 | PL1 | PL1 | PL1 | PL1 | PL2 | PL1 | PL1 | PL1 | PL1 |
| $E(n^{PL})$ | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 17 | 8 | 8 | 8 | 8 |
| $u_{cl}^{PL}$ | 56.36 | 62.13 | 68.23 | 83.91 | 96.43 | 99.90 | 99.99 | 83.90 | 84.97 | 96.97 | 57.76 | 70.73 | 87.54 |
| $u_{ch}^{PL}$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.87 | 0.00 | 0.00 | 0.00 | 79.87 | 0.00 | 0.00 | 0.00 | 2.96 |
| $u_{fl}^{PL}$ | 56.36 | 62.13 | 68.23 | 83.91 | 96.43 | 99.90 | 99.99 | 83.90 | 100.0 | 96.97 | 57.76 | 70.23 | 87.54 |
| $u_{fh}^{PL}$ | 42.99 | 49.13 | 55.91 | 74.96 | 92.97 | 99.65 | 99.99 | 74.95 | 97.86 | 59.68 | 87.78 | 80.43 | 72.58 |
| $D^{PL}(\%)$ | 2.01 | 2.01 | 2.01 | 2.01 | 2.01 | 2.01 | 2.01 | 2.32 | 0.34 | 1.70 | 1.64 | 2.18 | 1.86 |
| $E(U^{PL})$ | 41.08 | 45.81 | 50.85 | 64.98 | 76.28 | 79.85 | 79.98 | 46.92 | 86.99 | 72.36 | 38.75 | 57.23 | 66.14 |
| *Ranking of learning models in terms of welfare* | | | | | | | | | | | | | |
| Rank | NFP | FPN | NFP | NFP | FPN | FPN | FPN | FPN | PFN | FNP | FPN | NFP | FPN |

We now turn to the implications for the three sets of issues we want to examine using our simulation results.

### 4.1.1 Implications for $x_c$ and $n^{NL}$

From the first two rows, simulations 1–7, in the No Learning part, we see first that for both low and high values of $\rho$ the gap between $\tilde{n}$ and $\tilde{\tilde{n}}$ tends to zero. The reason is that, for small values of $\rho$, $\tilde{n}$ and $\tilde{\tilde{n}}$ are bunched closely to the upper limit $\bar{n}$ while for large values of $\rho$ they are bunched closely to the lower limit $n_h$. Positive gaps between them only occur for intermediate values of $\rho$. So, as can be seen from rows 3 and 4, it is more likely we get interior solutions for $x_c$ for intermediate values of $\rho$. So these results show that as $\rho$ increases, $n^{NL}$ steadily decreases from $\bar{n}$ to $n_h$.

As we saw from Lemma 1 as $p$ tends to 0, $\tilde{n}, \tilde{\tilde{n}}, \bar{n} \to n_h$, and hence the gap between $\tilde{n}$ and $\tilde{\tilde{n}}$ vanishes. This result was obtained in simulation 8 using $p = 0.05$. From Lemma 1 we also know that when $p$ tends to 1, the gap between $\tilde{n}$ and $\tilde{\tilde{n}}$ also converges to zero as $\tilde{n}, \tilde{\tilde{n}}, \bar{n} \to n_l$. In simulation 9 we set $p = 0.95$, which is not sufficiently high for the convergence between $\tilde{n}$ and $\tilde{\tilde{n}}$, and hence we get an interior solution. Similarly, from columns 10, 11, 12 and 13 we see that for lower values of $\bar{\bar{\gamma}}$ and larger spreads of $\gamma_s$ around $\bar{\bar{\gamma}}$ we are more likely to have significant gaps between $\tilde{n}, \tilde{\tilde{n}},$, and hence more likely to get an interior solution for $x_c$.

### 4.1.2 Implications of Variations in $\rho$ for $\tilde{p}$

The first row of the results for Partial Learning show, for different values of risk aversion $\rho$, the critical value of $\tilde{p}$ such that for $p \geq \tilde{p}$ there exists a second Partial Learning equilibrium. The results of simulations 1–7 show that as risk aversion increases $\tilde{p}$ tends to falls, so there is a wider range of parameter values for which there exists a second stable PL equilibrium. Not surprisingly, columns 8 and 9 confirm that $\tilde{p}$ does not depend on $p$.

However, in looking at simulations 10–13 there is no clear pattern of results. In fact, a broader range of simulation results we have carried out shows that $\tilde{p}$ is quite sensitive to variations in these parameter values. While the decreasing relationship between $\tilde{p}$ and $\rho$ shown in Table 1 holds fairly generally, it is not universally true for all parameter values, and these broader simulations show that there is no systematic relationship between $\tilde{p}$ and the parameters $\gamma_l$, $\gamma_h$. The simulations we report in Sect. 4.2 give a better perspective on how risk aversion affects $\tilde{p}$ and hence the likelihood of there being a second Partial Learning equilibrium.

### 4.1.3 Implications for Ranking Expected Size of Stable IEA, Payoffs and Expected Welfare Across Different Models of Learning

It is straightforward to see from Table 1 that for all the parameter values we have used it is always the case that:

$$E\left(n^{PL_1}\right) < E\left(n^{NL}\right) < E\left(n^{FL}\right) < E\left(n^{PL_2}\right) \tag{11}$$

consistent with (8a). This reflects the limited range of parameter values we used in Table 1.

In terms of the rankings of the payoffs for both members and fringe countries in the high and low damage cost states of the world, it is readily checked that the payoffs in Table 1 are consistent with the rather limited set of comparisons we were able to make in (9a) and (9b). Importantly the results also confirm that these are the only general results that hold.

**Table 2** Expected size of NL IEA for different degrees of risk aversion: $B = 1.1$

| $\rho$ | | 00.05 | 00.50 | 00.99 | 02.50 | 05.00 | 10.00 | 20.00 |
|---|---|---|---|---|---|---|---|---|
| *% cases relating $\tilde{n}, \tilde{\tilde{n}}, n^{NL}$:* | | | | | | | | |
| 1 | $\tilde{n} \leq n^{NL} < \tilde{\tilde{n}}$ $(0 < x^c < 1)$ | 13.54 | 43.47 | 50.45 | 49.90 | 38.30 | 23.47 | 13.14 |
| 2 | $\tilde{n} < n^{NL} = \tilde{\tilde{n}}$ $(x^c = 0)$ | 0.45 | 0.93 | 0.87 | 0.63 | 0.30 | 0.11 | 0.03 |
| 3 | $\tilde{n} = n^{NL} = \tilde{\tilde{n}}$ $(x^c = 0)$ | 86.01 | 55.59 | 48.68 | 49.47 | 61.40 | 76.42 | 86.83 |
| *% of cases relating $n^{NL}, n_h, \bar{n}$:* | | | | | | | | |
| 4 | $n^h = n^{NL} < \bar{n}$ | 0.37 | 4.18 | 9.46 | 30.18 | 57.12 | 77.83 | 86.72 |
| 5 | $n_h < n^{NL} = \bar{n}$ | 69.88 | 27.47 | 16.56 | 6.75 | 3.08 | 1.20 | 0.37 |
| 6 | $n_h < n^{NL} < \bar{n}$ | 20.82 | 59.46 | 65.02 | 54.08 | 30.45 | 11.41 | 3.33 |
| 7 | $n_h = n^{NL} = \bar{n}$ | 8.93 | 8.89 | 8.94 | 8.99 | 9.35 | 9.56 | 9.58 |

In particular, for the different parameter values used in Table 1, all possible rankings of the payoffs for NL relative to FL and PL1 are possible.

Finally, the last row of Table 1 shows the rankings of expected welfare across the different models of learning for the different parameter values, where for Partial Learning it is only for the parameters in column 9 that we are able to select PL2 as the appropriate equilibrium. Columns 1–7 confirm that as $\rho$ increases, welfare for NL moves from being the highest relative to FL and PL1 to the lowest, while columns 8–13 show that keeping $\rho = 2.50$, introducing the other variations in parameter values, other than tightening the spread of $\gamma_l$, $\gamma_h$ around the median (column 12), also imply that the welfare from NL moves from being the highest relative to FL and PL to the lowest.

## 4.2 Results for Full Simulations

### 4.2.1 Implications for $x^c$ and $n^{NL}$

Recall that from Kolstad and Ulph (2008) the expected size of the stable IEA with No Learning when countries are risk neutral is $n^{NL} = \bar{n}$. The first 3 rows of Table 2 shows how increasing $\rho$ affects how $n^{NL}$ relates to $\tilde{n}, \tilde{\tilde{n}}$ and whether $0 < x^c < 1$ or $x^c = 0$; while rows 4–7 show how $n^{NL}$ relates to $n_h, \bar{n}$.

Row 1 shows that as $\rho$ increases the proportion of cases for which $x^c$ is an interior solution increases and then falls. The reason is the same as we argued in Sect. 4.1.1. As shown in rows 3, 5 and 7, for $\rho$ close to 0, the large majority of cases have $n^{NL} = \bar{n}$, with $\tilde{n}, \tilde{\tilde{n}}$ bunched close to $\bar{n}$ and $\tilde{n} = n^{NL} = \tilde{\tilde{n}}$. As $\rho$ increases $\tilde{n}, \tilde{\tilde{n}}$ steadily decline, the gap between them opens up, allowing more opportunity for $x^c$ to take an interior value, and the proportion of cases where $n^{NL} = \bar{n}$ declines. But as $\rho$ continues to increase above 1 $\tilde{n}, \tilde{\tilde{n}}$ tend towards $n_h$, the gap between them narrows, allowing less opportunity for $x^c$ to take an interior value.

### 4.2.2 Implications of Variations in $\rho$ for $\tilde{p}$

The first row of Table 3 below shows for each value of $\rho$ the average value of $\tilde{p}$, while the second row shows the proportion of cases for which $p \geq \tilde{p}$ and hence for which we can

select the second equilibrium for Partial Learning, the one with the highest expected aggregate utility across the different models of learning.[17] Consistent with the results in Table 1, as $\rho$ increases the average value of $\tilde{p}$ decreases and the proportion of cases for which $p \geq \tilde{p}$ increases. However for any given value of $\rho$ there is significant variation in the value of $\tilde{p}$ caused by the variation of the other parameters in the simulations; we do not report the details here, but as an example with $\rho = 0.99$, $\tilde{p}$ varies between 0.7864 and 0.9999.

So the general message is that as risk aversion increases the proportion of cases for which the second Partial Learning equilibrium exists also increases.

### 4.2.3 Implications for Ranking Expected Size of Stable IEA, Payoffs and Expected Welfare Across Different Models of Learning

Rows 3–5 of Table 3 show for each model of learning the maximum value (across all 500,000 simulations) of damage costs as a percentage of the value of output.[18] Not surprisingly this occurs in the high damage cost state of the world for which the equilibrium is the same for Full Learning and the first Partial Learning equilibrium, which is why this percentage is the same for Full Learning and Partial Learning.

For No Learning the maximum value of this percentage will tend to occur when $n^{NL}$ is close to $n_h$, and hence the number of signatories is low, and so would be similar to Full Learning and Partial Learning. But there is an additional consideration. As we have seen from Table 2, there are parameter values for which signatory countries will also pollute, which will raise damage costs, and this effect is particularly pronounced for intermediate values of $\rho$.

Rows 6–8 show the % of cases for which the expected IEA membership size is highest in the three models of learning, while rows 9–11 show the similar figures for expected aggregate utility. Consistent with the theoretical results in Sect. 3.4.1, rows 6 and 9 show that on both measures the highest values occur for Partial Learning only when the second Partial Learning equilibrium is feasible, and so these percentages are exactly the same as the percentage of cases for which $p \geq \tilde{p}$, as shown in row 2 of Table 3.

Our key results are in rows 7–8 and 10–11, which show that as risk aversion increases from 0.0 to 20.0 the percentage of cases in which NL gives the highest number of expected signatories falls from just under 9.5% to close to zero, and the percentage of cases in which NL yields the highest level of expected utility falls from just under 87% to just over 2%. These trends reflect both the decline in the IEA membership from $\bar{n}$ to $n_h$ and the fact that, for intermediate values of $\rho$, IEA signatories produce positive emissions. The implications of the results for NL and PL are that as $\rho$ increases from 0.0 to 20.0 the percentage of cases where FL has the highest expected utility rises from just over 11% to just under 92%. These increases are monotonic between $\rho = 0.0$ and 10.0, but fall slightly when $\rho = 20.0$, due to the increase in the percentage of cases where PL has the highest expected signatories and expected utility from just under 4% to just under 6%. It is also important to note that percentage of cases where FL ranks first in expected utility exceeds that percentage for NL as $\rho$ increases from 0.0 to just 0.99.

So the general message is that quite modest increases in the level of risk aversion overturn the presumption from the literature using risk neutrality that No Learning yields the highest

---

[17] Of course, given our assumptions about the distribution of $p$ in the simulations the second measure is just 1 minus the first measure expressed as a percentage.

[18] In these simulations we report the maximum value of damage costs as a percentage of GDP, rather than the average value, since we believe this is a more useful indicator of the plausibility of the parameter values we have used.

**Table 3** Expected membership and aggregate utility for three scenarios of learning

| | $\rho$ | 00.00 | 00.05 | 00.50 | 00.99 | 02.50 | 05.00 | 10.00 | 20.00 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Average $\tilde{p}$ | .9795 | .9795 | .9794 | .9793 | .9775 | .9724 | .9607 | .9411 |
| 2 | Propn. cases where $p \geq \tilde{p}$ | 2.05 | 2.05 | 2.06 | 2.07 | 2.25 | 2.76 | 3.93 | 5.89 |
| 3 | $(\%D)^{PL}$ | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 |
| 4 | $(\%D)^{NL}$ | 6.67 | 8.50 | 8.80 | 9.13 | 9.36 | 9.28 | 8.96 | 6.67 |
| 5 | $(\%D)^{FL}$ | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 |
| *Highest* $E(n)$—% | | | | | | | | | |
| 6 | PL | 2.05 | 2.05 | 2.06 | 2.07 | 2.25 | 2.76 | 3.93 | 5.89 |
| 7 | NL | 9.46 | 9.09 | 6.68 | 5.14 | 2.73 | 1.31 | 0.49 | 0.14 |
| 8 | FL | 88.49 | 88.86 | 91.25 | 92.79 | 95.02 | 95.93 | 95.58 | 93.97 |
| *Highest* $E(U)$—% | | | | | | | | | |
| 9 | PL | 2.05 | 2.05 | 2.06 | 2.07 | 2.25 | 2.76 | 3.93 | 5.89 |
| 10 | NL | 86.84 | 81.80 | 51.31 | 36.30 | 16.63 | 7.70 | 4.02 | 2.39 |
| 11 | FL | 11.11 | 16.15 | 46.62 | 61.63 | 81.12 | 89.54 | 92.05 | 91.72 |

level of expected utility. This is also true for the expected utility of an individual country from an ex ante perspective, which, for symmetric countries, is simply aggregate expected utility divided by the number of countries. Thus, if countries were able to choose which model of learning they should adopt, then, for quite modest levels of risk aversion, they would favour leaving the decision to form an IEA and set their emissions until they had full information about the likely damage cost of climate change.

### 4.3 Robustness Checks

In this sub-section we briefly discuss the sensitivity of the results we derived in Sects. 4.1 and 4.2 to changes in underlying assumptions in three respects.

First, in the simulations in Sects. 4.1 and 4.2 we set the value of the parameter $B$ in $\pi_0$ to $B = 1.1$ to ensure that payoffs were sensitive to countries' abatement decisions. We have also run the simulation results reported in Tables 1, 2 and 3 above with higher values of $B$, whilst also ensuring that damage costs from emissions remain significant. We present the results in "Appendix B".[19] These show that for higher values of $B$, the key findings in Tables 1, 2 and 3 are slightly less sensitive to changes in risk aversion, but the overall results are robust; in particular, in relation to the results in Table 3, the percentage of cases where the payoff from Full Learning is the highest across all models of learning now increases monotonically in $\rho$ (from 11% to 92%), and is greater than the percentage of cases where the payoff from No Learning is highest for values of $\rho \geq 0.99$.

Second, we checked the limiting cases we discussed in Lemmas 2 and 3 in Sect. 3.4 above through numerical simulations. To do this we used essentially the same underlying parameter values as in Table 1, $N = 20$, $B = 1.1$, $\gamma_l = 0.0615$, $\gamma_h = 0.1289$, and $p = 0.5$, but then successively took limiting values for individual parameters of $p \approx 0.0$, $p \approx 1.0$, $\gamma_l \approx 1/N$, $\gamma_h \approx 1.0$, holding all other parameters at their original levels. To save space in Table 4 we report the results for 3 values of $\rho = 0.05$, 2.50, 10.0, which illustrate the limiting results derived in Lemmas 2 and 3. Moreover, the simulation results also indicate that for $p \approx 0$ not only $E(U^{FL}) > E(U^{PL})$ but also $E(U^{FL}) > E(U^{NL})$, reinforcing the message that learning leads to a better outcome in terms of expected utility. For $p \approx 1$, besides $E(U^{PL}) = E(U^{NL})$ the simulation results reveal a full ranking of expected utility, $E(U^{FL}) \leq E(U^{PL}) = E(U^{NL})$, which indicate a negative effect of learning. For the third limit parameter value, $\gamma_l = 1/N$, the results indicate that not only $E(U^{FL}) > E(U^{PL})$ but also the level of risk aversion affects the ranking of expected utility between FL and NL, with $E(U^{NL}) > E(U^{FL})$ for $\rho = 0.05$ and $E(U^{FL}) > E(U^{NL})$ for $\rho = 10$. Finally, for $\gamma_h = 1$ a full ranking of the learning scenarios was obtained: $E(U^{PL}) > E(U^{FL}) \geq E(U^{NL})$.

Finally, all the simulation results we have reported so far have been based on the constant relative risk aversion (CRRA) utility function:

$$u(\pi_i) = \frac{1}{(1-\rho)}\pi_i^{1-\rho}$$

where $\rho \geq 0$ measures the degree of relative risk aversion. We now test the robustness of our key results in Table 3 to the use of a constant absolute risk aversion (CARA) utility function:

$$u(\pi_i) = 1 - e^{-\lambda \pi_i} \qquad (12)$$

---

[19] To save space we report these results in "Appendix B" available on request from the authors.

brief

**Table 4** Results for limit cases: $\varepsilon \rightarrow 0$

| Limit case | $p = \varepsilon$ | | | $p = 1 - \varepsilon$ | | | $\gamma_l = (1/N) + \varepsilon$ | | | $\gamma_h = 1 - \varepsilon$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\rho$ | 0.05 | 2.50 | 10.0 | 0.05 | 2.50 | 10.0 | 0.05 | 2.50 | 10.0 | 0.05 | 2.50 | 10.0 |
| Results $E(n)$ | | | | | | | | | | | | |
| $E(n^{PL})$ | 8 | 8 | 8 | 17 | 17 | 17 | 8 | 8 | 8 | 17 | 17 | 17 |
| $E(n^{NL})$ | 8 | 8 | 8 | 17 | 17 | 17 | 12 | 9 | 8 | 2 | 2 | 2 |
| $E(n^{FL})$ | $8+\varepsilon$ | $8+\varepsilon$ | $8+\varepsilon$ | $17-\varepsilon$ | $17-\varepsilon$ | $17-\varepsilon$ | 14 | 14 | 14 | 9.5 | 9.5 | 9.5 |
| Rank | $E(n^{FL}) > E(n^{PL}) = E(n^{NL})$ | | | $E(n^{PL}) = E(n^{NL}) > E(n^{FL})$ | | | $E(n^{FL}) > E(n^{NL}) >$ $E(n^{PL}) = n_h$ | | | $E(n^{PL}) = n_l > E(n^{FL}) >$ $E(n^{NL})$ | | |
| Results $E(U)$ | | | | | | | | | | | | |
| $E(U^{PL})$ | 25.7 | 50.0 | 59.7 | 64.6 | 87.2 | 99.9 | 42.7 | 65.7 | 79.8 | 88.4 | 99.0 | 99.9 |
| $E(U^{NL})$ | 25.8 | 50.0 | 59.8 | 64.6 | 87.2 | 99.9 | 48.9 | 67.9 | 79.8 | 49.9 | 66.7 | 91.9 |
| $E(U^{FL})$ | 25.9 | 50.1 | 59.9 | 64.5 | 87.1 | 99.9 | 42.8 | 65.8 | 79.9 | 50.4 | 66.8 | 91.9 |
| Rank | $E(U^{FL}) > E(U^{NL}) \geq E(U^{PL})$ | | | $E(U^{FL}) \leq E(U^{PL}) = E(U^{NL})$ | | | $E(U^{FL}) >$ $E(U^{PL})$; $E(U^{NL}) \geq E(U^{PL})$ | | | $E(U^{PL}) > E(U^{FL}) \geq$ $E(U^{NL})$ | | |

**Table 5** Expected membership and aggregate utility for three scenarios of learning—CARA utility function

| | $\rho$ | 00.05 | 00.50 | 00.99 | 02.50 | 05.00 | 10.00 | 20.00 |
|---|---|---|---|---|---|---|---|---|
| 1 | $\lambda$ | 0.017 | 0.167 | 0.333 | 0.833 | 1.667 | 3.333 | 6.667 |
| 2 | $p \geq \tilde{p}$ | 2.01 | 2.07 | 2.06 | 2.12 | 2.26 | 2.63 | 3.49 |
| 3 | $(\%D)^{PL}$ | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 |
| 4 | $(\%D)^{NL}$ | 8.55 | 8.69 | 8.84 | 8.94 | 9.29 | 9.24 | 8.64 |
| 5 | $(\%D)^{FL}$ | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 | 6.67 |
| *Highest $E(n)$—%* | | | | | | | | |
| 6 | PL | 2.01 | 2.07 | 2.06 | 2.12 | 2.26 | 2.63 | 3.49 |
| 7 | NL | 9.28 | 7.83 | 6.79 | 4.78 | 3.16 | 1.71 | 0.69 |
| 8 | FL | 88.71 | 90.10 | 91.15 | 93.10 | 94.58 | 95.66 | 95.82 |
| *Highest $E(U)$—%* | | | | | | | | |
| 9 | PL | 2.01 | 2.07 | 2.06 | 2.12 | 2.26 | 2.63 | 3.49 |
| 10 | NL | 82.57 | 55.90 | 42.29 | 22.87 | 10.46 | 4.34 | 2.16 |
| 11 | FL | 15.42 | 42.03 | 55.65 | 75.01 | 87.28 | 93.03 | 94.35 |

where $\lambda > 0$ measures the degree of absolute risk aversion. In checking how our results using CARA relate to those using CRRA we need to ensure that the values we choose for the parameters $\lambda$ are broadly consistent with those we chose for $\rho$. For any level of $\pi_i$ the coefficient of absolute risk aversion is related to the coefficient of relative risk aversion by the formula $\lambda = \rho/\pi_i$, and we discuss in "Appendix A2" how we use this relationship to choose values of $\lambda$ which are consistent with those we have used for $\rho$.

We have re-run the simulation results we presented in Table 3, which used the CRRA utility function, but using the CARA utility function. The choice of the ranges for all other parameter values other than $\rho$ are exactly the same as in Table 3. The results are shown in Table 5, where, as we have just discussed above, we have chosen the value of the parameter $\lambda$ in the CARA utility function to be consistent with the values of $\rho$ we used in Table 3.

It is clear to see that while the individual numbers in Table 5 are slightly different from those in Table 3 all the key findings we derived from Table 3 relating to the percentages of cases for which the expected size of IEA or expected utility are largest across the three models of learning carry over to those in Table 5.

So we conclude that the key findings from the simulation results reported in Sects. 4.1 and 4.2 are robust to a number of changes in the main aspects of the simulation modeling we have employed.

## 5 Summary and Conclusions

This paper bridges two strands of literature on the formation of IEAs under uncertainty by addressing the combined roles of learning and risk aversion. This approach allowed us to explore the impact of learning for any given level of risk aversion as well as the impact of changing risk aversion under various scenarios of learning.

We generalized the model of Kolstad and Ulph (2008) who showed that with risk neutrality the possibility of learning more information about environmental damage costs generally had rather pessimistic implications for the success of the formation of IEAs. The authors found

that learning reduces expected aggregate and individual payoffs for a wide range of parameter values. This suggests that countries are better off forming an IEA rather than waiting for better information.

In this paper, we have allowed countries to be risk averse using an expected utility approach which maps payoffs into utility. We first derived the theoretical results for each of our three scenarios of learning with risk aversion, confirming the main findings of Boucher and Bramoullé (2010) for the No Learning case. For No learning, risk leads to smaller stable coalitions and higher global emissions. In terms of equilibrium coalitions and global emissions, we showed that Full Learning remains unaffected by risk and changes for Partial Learning are small. However, even with special functional forms for the underlying utility functions there was limited scope for deriving analytical comparisons across our three scenarios of learning, primarily because welfare effects could differ for signatory and non-signatory countries. Our simulation results showed that contrary to the finding with risk neutrality, when countries become significantly risk averse, the set of parameter values for which countries are better off with No Learning compared to Full Learning shrinks significantly and those cases for which this is reversed increases accordingly. This may explain why it has taken so long for a proper climate agreement to be reached—countries are risk averse and waited till they had much better information about the risks of climate change.

In terms of future research, it would be desirable to use a model with asymmetric countries, though it is unlikely to be possible to derive analytical results; so it may be more useful to introduce different models of learning into Integrated Assessment Models of climate change. It would also be interesting to endogenise the process of learning by allowing countries to invest in research in order to obtain better information.

## Appendix A1: Proofs of Results

### Proposition 1: Full Learning

Since the true state of the world is revealed before countries decide whether to join an IEA, the results for stable IEAs, probabilities and payoffs in the four states of the world shown in Proposition 1 follow immediately from Kolstad and Ulph (2008). The difference from Kolstad and Ulph is that we now calculate aggregate expected utility using the expected utility approach which captures attitudes to risk, rather than expected payoff. This completes the proof of Proposition 1.

### Lemma 1: Emission Decisions with No Learning

### Fringe Country Output Decision

In stage 2, fringe country $i$ takes as given the output of all other countries, $X_{-i}$, and chooses $x_i$ to maximise $pu[\pi_0 + x_i - \gamma_l(x_i + X_{-i})] + (1-p)u[\pi_0 + x_i - \gamma_h(x_i + X_{-i})]$. Since $\gamma_l < \gamma_h < 1$, $x^f(n) = 1 \ \forall n, 2 \le n \le N$. This proves result (i).

### Member Country Emission Decision

From (4b) we obtain the following:

$$0 < \pi_{c,h}(x_c, n) < \pi_{c,l}(x_c, n) \ \forall \, n \geq 1, \ 0 \leq x_c \leq 1; \tag{A1a}$$

$$\frac{\partial \pi_{c,s}}{\partial x_c} \geq 0 \Leftrightarrow n \leq \frac{1}{\gamma_s}; \tag{A1b}$$

$$\pi_{c,s}(0, n) = \pi_0 - N\gamma_s + n\gamma_s; \qquad \pi_{c,s}(1, n) = \pi_0 - N\gamma_s + 1 \equiv \pi_{c,s}(1) \quad \forall n \tag{A1c}$$

and hence

$$\pi_{c,h}(1) < \pi_{c,h}(0, n) < \pi_{c,l}(0, n) < \pi_{c,l}(1) \text{ if } \frac{1}{\gamma_h} < n < \frac{1}{\gamma_l}. \tag{A1d}$$

As stated in Sect. 3.2, for a given $n$, each coalition member chooses $x_c(n)$ to maximise $E(u_c(x_c(n), n))) \equiv pu(\pi_{c,l}(x_c(n), n)) + (1 - p)u(\pi_{c,h}(x_c(n), n))$ which leads to the following first and second order conditions:

$$\frac{\partial E(u_c)}{\partial x_c} = p(1 - \gamma_l n)u'[\pi_{c,l}(x_c, n)] + (1 - p)(1 - \gamma_h n)u'[\pi_{c,h}(x_c, n)], \tag{A2a}$$

$$\frac{\partial^2 E(u_c)}{\partial(x_c)^2} = p(1 - \gamma_l n)^2 u''[\pi_{c,l}(x_c(n), n)] + (1 - p)(1 - \gamma_h n)^2 u''[\pi_{c,h}(x_c(n), n)] < 0. \tag{A2b}$$

### Boundary Values for $x_c(n)$

From (A2a):

$$n \geq \frac{1}{\gamma_l} > \frac{1}{\gamma_h} \Rightarrow \frac{\partial E(u_c)}{\partial x_c} < 0 \ \forall \, x_c \Rightarrow x_c(n) = 0;$$

$$n \leq \frac{1}{\gamma_h} < \frac{1}{\gamma_l} \Rightarrow \frac{\partial E(u_c)}{\partial x_c} > 0 \ \forall \, x_c \Rightarrow x_c(n) = 1.$$

We now look for tighter bounds for $n$ that guarantees $x_c(n) = 0$ and $x_c(n) = 1$, respectively. So we now focus on the range $n_h - 1 < \frac{1}{\gamma_h} < n < \frac{1}{\gamma_l} \leq n_l$. From (A1b), in this range $\frac{\partial \pi_{c,l}}{\partial x_c} > 0; \quad \frac{\partial \pi_{c,h}}{\partial x_c} < 0$.

To make progress, we treat $n$ as if it was a real value, $z$. To save notation define:

$\tilde{\theta}_l \equiv u'[\pi_{c,l}(1)], \tilde{\theta}_h \equiv u'[\pi_{c,h}(1)], \tilde{\tilde{\theta}}_l(z) \equiv u'[\pi_{c,l}(0, z)]$ and $\tilde{\tilde{\theta}}_h(z) \equiv u'[\pi_{c,h}(0, z)]$ where, from (A1d):

$$\tilde{\theta}_l < \tilde{\tilde{\theta}}_l(z) < \tilde{\tilde{\theta}}_h(z) < \tilde{\theta}_h. \tag{A3}$$

We first define $\tilde{z}$ as the unique value of $z$ such that:

$$\frac{\partial E(u_c(1, \tilde{z}))}{\partial x_c} = p(1 - \gamma_l \tilde{z})\tilde{\theta}_l + (1 - p)(1 - \gamma_h \tilde{z})\tilde{\theta}_h = 0 \Rightarrow \tilde{z} = \frac{p\tilde{\theta}_l + (1 - p)\tilde{\theta}_h}{\gamma_l p\tilde{\theta}_l + \gamma_h(1 - p)\tilde{\theta}_h}. \tag{A4}$$

From (A4) we get: $\frac{1}{\gamma_h} < \tilde{z} < \frac{1}{\gamma_l}$ and $\frac{\partial E(u_c(1, \tilde{z}))}{\partial x_c} > 0 \quad \forall \, z < \tilde{z}$. Thus, $x_c(z) = 1 \ \forall \, z \leq \tilde{z}$.

We now define $\tilde{\tilde{z}}$ such that:

$$\frac{\partial E(u_c(0, \tilde{\tilde{z}}))}{\partial x_c} = p(1 - \gamma_l \tilde{\tilde{z}})\tilde{\tilde{\theta}}_l(\tilde{\tilde{z}}) + (1 - p)(1 - \gamma_h \tilde{\tilde{z}})\tilde{\tilde{\theta}}_h(\tilde{\tilde{z}}) = 0$$

$$\Rightarrow \tilde{\tilde{z}} = \frac{p\tilde{\tilde{\theta}}_l(\tilde{\tilde{z}}) + (1 - p)\tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})}{p\gamma_l \tilde{\tilde{\theta}}_l(\tilde{\tilde{z}}) + (1 - p)\gamma_h \tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})} \tag{A5}$$

and again $\frac{1}{\gamma_h} < \tilde{\tilde{z}} < \frac{1}{\gamma_l}$. Note, we have not been able to prove that there is a unique value of $\tilde{\tilde{z}}$ which solves (A5). We will discuss the implications shortly, but the next steps apply to any $\tilde{\tilde{z}}$ which solves (A5). As $\tilde{\theta}_l > 0$, $\tilde{\theta}_h > 0$, $\tilde{\tilde{\theta}}_l > 0$ and $\tilde{\tilde{\theta}}_h > 0$ then:

$$sign\left(\tilde{\tilde{z}} - \tilde{z}\right) = sign\left(\begin{array}{l} [p\tilde{\tilde{\theta}}_l(\tilde{\tilde{z}}) + (1 - p)\tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})][p\gamma_l \tilde{\theta}_l + (1 - p)\gamma_h \tilde{\theta}_h] \\ -[p\tilde{\theta}_l + (1 - p)\tilde{\theta}_h][p\gamma_l \tilde{\tilde{\theta}}_l(\tilde{\tilde{z}}) + (1 - p)\gamma_h \tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})] \end{array}\right)$$

$$= sign\left(p(1 - p)(\gamma_h - \gamma_l)\tilde{\theta}_l \tilde{\tilde{\theta}}_l(\tilde{\tilde{z}})[\frac{\tilde{\theta}_h}{\tilde{\theta}_l} - \frac{\tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})}{\tilde{\tilde{\theta}}_l(\tilde{\tilde{z}})}]\right) \tag{A6}$$

Using (A3), it can be shown that:

$sign\left(\tilde{\tilde{z}} - \tilde{z}\right) = sign\left\{p(1 - p)(\gamma_h - \gamma_l)\tilde{\theta}_l \tilde{\tilde{\theta}}_l(\tilde{\tilde{z}})[\frac{\tilde{\theta}_h}{\tilde{\theta}_l} - \frac{\tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})}{\tilde{\tilde{\theta}}_l(\tilde{\tilde{z}})}]\right\} \geq 0$. So any $\tilde{\tilde{z}}$ which solves (A5) must be at least as large as $\tilde{z}$.

Next we show that $\tilde{\tilde{z}} \leq \frac{1}{\bar{\gamma}}$. Define $\xi \equiv \tilde{\tilde{\theta}}_h(\tilde{\tilde{z}})/\tilde{\tilde{\theta}}_l(\tilde{\tilde{z}}) \geq 1$. Then:

$$\tilde{\tilde{z}} \leq 1/\bar{\gamma} \Leftrightarrow \frac{p + (1 - p)\xi}{p\gamma_l + (1 - p)\gamma_h \xi} \leq \frac{p + (1 - p)}{p\gamma_l + (1 - p)\gamma_h}$$

$$\Leftrightarrow p(1 - p)(\gamma_h + \xi\gamma_l) \leq p(1 - p)(\xi\gamma_h + \gamma_l) \Leftrightarrow (\xi - 1)(\gamma_h - \gamma_l) \geq 0. \tag{A7}$$

So any $\tilde{\tilde{z}}$ which solves (A5) must be no greater than $1/\bar{\gamma}$.

As $\tilde{n} = I(\tilde{z})$ and $\tilde{\tilde{n}} = I\left(\tilde{\tilde{z}}\right)$ then $n_h \leq \tilde{n} \leq \tilde{\tilde{n}} \leq \bar{n}$—result (ii).

From (A4) and (A5) it is straightforward to see that:

$p \to 0 \Rightarrow \tilde{n}, \tilde{\tilde{n}}, \bar{n} \to n_h$;   $p \to 1 \Rightarrow \tilde{n}, \tilde{\tilde{n}}, \bar{n} \to n_l$—result (iii).

Results (iv), $n < \tilde{n} \Rightarrow x_c(n) = 1$, and (vii), $n \geq \tilde{\tilde{n}} \Rightarrow x_c(n) = 0$, come directly from the definitions of $\tilde{n}$ and $\tilde{\tilde{n}}$, in (6a) and (6b), respectively.

We have not been able to prove analytically that there is a unique value of $\tilde{\tilde{z}}$. Thus, following definition (6b), we will treat $\tilde{\tilde{z}}$ and hence $\tilde{\tilde{n}}$ as the *largest* such value.

## Interior Values for $x^c(n)$

If $\tilde{z}$ is an integer then $\tilde{n} = \tilde{z}$ and $x^c(\tilde{n}) = 1$ as $\frac{\partial E(u_c(1, \tilde{n}))}{\partial x_c} = 0$. If $\tilde{z}$ is not integer then $0 \leq x^c(\tilde{n}) < 1$ as $\frac{\partial E(u_c(1, \tilde{n}))}{\partial x_c} < 0$. Thus, $n = \tilde{n} \Rightarrow 0 \leq x^c(n) \leq 1$—result (v).

From the definition of $\tilde{n}$ in (6a), $n > \tilde{n} \Rightarrow 0 \leq x^c(n) < 1$ as $\frac{\partial E(u_c(1, n))}{\partial x_c} < 0$. From the definition of $\tilde{\tilde{n}}$ in (6b), $n < \tilde{\tilde{n}} \Rightarrow 0 \leq x^c(n) \leq 1$. Thus, $\tilde{n} < n < \tilde{\tilde{n}} \Rightarrow 0 \leq x^c(n) < 1$—result (vi). Moreover, if there is a unique value of $\tilde{\tilde{z}}$ then $n < \tilde{\tilde{n}} \Rightarrow 0 < x^c(n) \leq 1$ and hence $\tilde{n} \leq n < \tilde{\tilde{n}} \Rightarrow 0 < x_c(n) < 1$. This completes the proof of Lemma 1.

## Proposition 2: No Learning

From Lemma 1 we know the values of $x_f(n)$ and $x_c(n)$ for all $n$. We now solve Stage 1. Define:
$$E(u_f(n)) = pu(\pi_0 + 1 - \gamma_l[N - n + nx_c(n)]) + (1 - p)u(\pi_0 + 1 - \gamma_h[N - n + nx_c(n)]),$$
$$E(u_c(n)) = pu(\pi_0 + x_c(n) - \gamma_l[N - n + nx^c(n)]) + (1 - p)u(\pi_0 + x_c(n) - \gamma_h[N - n + nx_c(n)])$$
and $\Delta(n) = E(u_c(n)) - E(u_f(n - 1))$, noting that $n^{NL} = E(n^{NL})$ is a stable IEA iff $\Delta(n^{NL}) \geq 0$ and $\Delta(n^{NL} + 1) < 0$.

To save notation, define $\pi_s = \pi_0 - \gamma_s N$  $s = l, h$.

(i)  $n \geq \tilde{\tilde{n}} + 1$

$$\Delta(n) = pu(\pi_l + \gamma_l n) + (1 - p)u(\pi_h + \gamma_h n) - pu(\pi_l + \gamma_l n + (1 - \gamma_l))$$
$$- (1 - p)u(\pi_h + \gamma_h n + (1 - \gamma_h)) < 0$$

So no such $n$ could be stable.

(ii)  $n < \tilde{n}$

All countries set $x = 1$, so $\Delta(n) = 0$. We exclude these trivial cases.

(iii)  $\tilde{n} + 1 \leq n < \tilde{\tilde{n}}$

$$\Delta(n) = p[\pi_l + n\gamma_l + x_c(n)(1 - \gamma_l) - \gamma_l(n - 1)x_c(n)]$$
$$+ (1 - p)[\pi_h + n\gamma_h + x_c(n)(1 - \gamma_h) - \gamma_h(n - 1)x_c(n)]$$
$$- p[\pi_l + n\gamma_l + (1 - \gamma_l) - \gamma_l(n - 1)x_c(n - 1)]$$
$$- (1 - p)[\pi_h + n\gamma_h + (1 - \gamma_h) - \gamma_h(n - 1)x_c(n - 1)]$$

So, if $x_c(n) \geq x_c(n - 1)$, then $\Delta(n) < 0$ and so $n$ could not be a stable IEA. If $x_c(n) < x_c(n - 1)$, then the sign of $\Delta(n)$ depends on the precise values of $x_c(n)$ and $x_c(n - 1)$.

(iv)  $n = \tilde{n}$

From (ii) we know that $x_c(\tilde{n} - 1) = x_f(\tilde{n}) = 1 \Rightarrow E(u_c(\tilde{n} - 1)) = E(u_f(\tilde{n}))$. So when $n = \tilde{n}$, IEA members could have set $x_c(\tilde{n}) = 1$, but, by definition of $\tilde{n}$, they did not, so $E(u_c(\tilde{n})) > E(u_c(\tilde{n} - 1)) = E(u_f(\tilde{n} - 1)) \Rightarrow \Delta(\tilde{n}) > 0$.

To complete the argument, successively increase $n$ from $\tilde{n}$ until $\Delta(n) < 0$, in which case $n^{NL} = n - 1$ is a stable IEA. There must exist such a stable IEA since we know from (iv) that $\Delta(\tilde{n}) > 0$ and from (i) that $\Delta(\tilde{\tilde{n}} + 1) < 0$. We cannot rule out the possibility that there is more than one stable IEA. This completes the proof of Proposition 2.

## Proposition 3: Partial Learning

Because emission decisions in stages 2 are taken under certainty, the results are the same as in Kolstad and Ulph (2008), namely:

*Stage 2*

(a)  $x_{f,s}(n) = 1$  $s = l, h$  $\forall n$
(b)  $n \geq n_l \Rightarrow x_{c,s}(n) = 0$  $s = l, h$
(c)  $n_h \leq n < n_l \Rightarrow x_{c,h} = 0, \ x_{c,l} = 1$
(d)  $n < n_h \Rightarrow x_{c,s}(n) = 1$  $s = l, h$

So using the notation $\pi_s = \pi_0 - \gamma_s N, s = l, h$, introduced in the proof of Proposition 2, the payoffs are:

(i)  $n \geq n_l$

$$E(u_f(n)) = pu[\pi_l + n\gamma_l + 1] + (1 - p)u[\pi_h + n\gamma_h + 1]$$
$$E(u_c(n)) = pu[\pi_l + n\gamma_l] + (1 - p)u[\pi_h + n\gamma_h]$$

(ii)  $n_h \leq n < n_l$

$$E(u_f(n)) = pu[\pi_l + 1] + (1 - p)u[\pi_h + n\gamma_h + 1]$$
$$E(u_c(n)) = pu[\pi_l + 1] + (1 - p)u[\pi_h + n\gamma_h]$$

(iii)  $n < n_h$

$$E(u_f(n)) = E(u_c(n)) = pu[\pi_l + 1] + (1 - p)u[\pi_h + 1]$$

*Stage 1*

(i)  $n \geq n_l + 1$

$$\Delta(n) = pu(\pi_l + n\gamma_l) + (1 - p)u(\pi_h + n\gamma_h) - pu(\pi_l + n\gamma_l + (1 - \gamma_l))$$
$$- (1 - p)u(\pi_h + n\gamma_h + (1 - \gamma_h)) < 0$$

So no $n$ in this range can be a stable IEA.

(ii)  $n = n_l$

$$\Delta(n_l) = pu(\pi_l + n_l\gamma_l) + (1 - p)u(\pi_h + n_l\gamma_h) - pu(\pi_l + 1)$$
$$- (1 - p)u(\pi_h + n_l\gamma_h + (1 - \gamma_h))$$

As $n_l\gamma_l \geq 1$ so:

$$\Delta(n_l) \geq 0 \Leftrightarrow \frac{p}{(1 - p)} \geq \frac{u(\pi_h + n_l\gamma_h + (1 - \gamma_h)) - u(\pi_h + n_l\gamma_h)}{u(\pi_l + n_l\gamma_l) - u(\pi_l + 1)} \equiv \chi > 0$$

i.e. $\Delta(n_l) \geq 0 \Leftrightarrow p \geq \frac{\chi}{1+\chi} \equiv \tilde{p} > 0$.

Since, from (i), $\Delta(n_l + 1) < 0$, $n_l$ is stable iff $p \geq \tilde{p}$.

(iii)  $n_h + 1 \leq n < n_l$

$$\Delta(n) = pu(\pi_l + 1) + (1 - p)u(\pi_h + n_h\gamma_h) - pu(\pi_l + 1) - (1 - p)u(\pi_h + n\gamma_h + (1 - \gamma_h))$$
$$< 0$$

So no $n$ in this range can be stable.

(iv)  $n = n_h$

$$\Delta(n_h) = pu(\pi_l + 1) + (1 - p)u(\pi_h + n_h\gamma_h) - pu(\pi_l + 1) - (1 - p)u(\pi_h + 1)$$

$\Delta(n_h) \geq 0$ as $n_h\gamma_h \geq 1$. From (iii) $\Delta(n_h + 1) < 0$ and hence $n_h$ is always a stable IEA.

(v)  $n < n_h$

$\Delta(n) = 0$; as in Proposition 2, we ignore these cases.

So there always exists a stable IEA with $n^{PL} = n_h$ and if $p \geq \tilde{p}$ there is a second stable IEA with $n^{PL} = n_l$. This completes the proof of Proposition 3.

## Lemma 2: Expected Membership

(i) Let $p = \varepsilon \approx 0$. From Proposition 1, Lemma 1 and Proposition 3 we get:

$$E(n^{FL}) = \varepsilon n_l + (1 - \varepsilon)n_h, \ E(n^{NL}) = n_h, \ E(n^{PL}) = n_h.$$

This proves $E(n^{FL}) > E(n^{PL}) = E(n^{NL})$.

(ii) Let $p = 1 - \varepsilon \approx 1$. From Proposition 1, Lemma 1 and Proposition 3 we get:

$$E\left(n^{FL}\right) = (1 - \varepsilon)n_l + \varepsilon n_h, \; E\left(n^{NL}\right) = n_l, \; E\left(n^{PL}\right) = n_l.$$

This proves $E(n^{PL}) = E(n^{NL}) > E(n^{FL})$.

(iii) and (iv). From Proposition 1 we get:

$$E(n^{FL}) = pn_l + (1 - p)n_h.$$

From Proposition 3, under PL there always exists a stable IEA with $n^{PL_1} = n_h$ members. If $\tilde{p} \leq p \leq 1$, then there is a second stable IEA with $n^{PL_2} = n_l$ members, where $\tilde{p} = \frac{\chi}{1+\chi}$ and $\chi = \frac{u[\pi_0 - (N-n_l)\gamma_h + (1-\gamma_h)] - u[\pi_0 - (N-n_l)\gamma_h]}{u[\pi_0 - (N-n_l)\gamma_l] - u[\pi_0 - N\gamma_l + 1]} = \frac{\Omega}{\Psi}$.

Thus, $\tilde{p} = \frac{\chi}{1+\chi}$ can be rewritten as:

$$\tilde{p} = \frac{\Omega}{\Omega + \Psi}$$

$$= \frac{u[\pi_0 - (N-n_l)\gamma_h + (1-\gamma_h)] - u[\pi_0 - (N-n_l)\gamma_h]}{u[\pi_0 - (N-n_l)\gamma_h + (1-\gamma_h)] - u[\pi_0 - (N-n_l)\gamma_h] + u[\pi_0 - (N-n_l)\gamma_l] - u[\pi_0 - N\gamma_l + 1]}$$

If $\gamma_l \approx 1/N$ then $n_l = N$. Thus: $\tilde{p} = \frac{u[\pi_0 + (1-\gamma_h)] - u[\pi_0]}{u[\pi_0 + (1-\gamma_h)] - u[\pi_0 - N\gamma_l + 1]}$.

As $\gamma_l \approx 1/N$ we get $\tilde{p} \approx 1$ and consequently the second stable IEA with $n^{PL_2} = n_l$ members will not occur. Therefore, $n^{PL} = n_h$ and $E\left(n^{FL}\right) = pn_l + (1-p)n_h > n^{PL}$, which proves (iii).

If $\gamma_h = 1$ then $\tilde{p} \approx \frac{0}{u[\pi_0 - (N-n_l)\gamma_l] - u[\pi_0 - N\gamma_l + 1]} = 0$. Hence, the second stable IEA with $n^{PL_2} = n_l$ members will always occur, which Pareto-dominates $n^{PL_1} = n_h$. Therefore, $n^{PL} = n_l$ and $n^{PL} > E\left(n^{FL}\right) = pn_l + (1-p)n_h$, which proves (iv).

## Lemma 3: Expected Utility

(i) If $p \approx 0$ then from Propositions 1 and 3, the difference of expected utilities between PL and FL is given by:

$$\Delta_1 = E\left(U^{PL}\right) - E\left(U^{FL}\right)$$

$$= pn_h u\left(\pi_{c,l}^{PL_1}\right) + p(N - n_h)u\left(\pi_{f,l}^{PL_1}\right) + (1-p)n_h u\left(\pi_{c,h}^{PL_1}\right) + (1-p)(N - n_h)u\left(\pi_{f,h}^{PL_1}\right)$$

$$- \left[pn_l u\left(\pi_{c,l}^{FL}\right) + p(N - n_l)u\left(\pi_{f,l}^{FL}\right) + (1-p)n_h u\left(\pi_{c,h}^{FL}\right) + (1-p)(N - n_h)u\left(\pi_{f,h}^{FL}\right)\right]$$

where:

$$\pi_{c,l}^{PL_1} = \pi_0 + 1 - \gamma_l N < \pi_{c,l}^{FL_1} = \pi_0 - \gamma_1(N - n_1), \pi_{f,l}^{PL_1} = \pi_0 + 1 - \gamma_l N < \pi_{f,l}^{FL_1} = \pi_0 + 1 - \gamma_1(N - n_1),$$

$$\pi_{c,h}^{PL_1} = \pi_{c,h}^{FL} = \pi_0 - \gamma_h(N - n_h), \pi_{f,l}^{PL_1} = \pi_{f,l}^{PL} = \pi_0 + 1 - \gamma_1(N - n_h).$$

As $\pi_{c,h}^{PL_1} = \pi_{c,h}^{FL}$ and $\pi_{f,h}^{PL_1} = \pi_{f,h}^{FL}$, $\Delta_1$ can be simplified to:

$$\Delta_1 = p\left[n_h u\left(\pi_{c,l}^{PL_1}\right) - n_l u\left(\pi_{c,l}^{FL}\right)\right] + p\left[(N - n_h)u\left(\pi_{f,l}^{PL_1}\right) - (N - n_l)u\left(\pi_{f,l}^{FL}\right)\right]$$

Thus,

$$\frac{d\Delta_1}{dp} = \left[ n_h u\left(\pi_{c,l}^{PL_1}\right) - n_l u\left(\pi_{c,l}^{FL}\right) \right] + \left[ (N - n_h) u\left(\pi_{f,l}^{PL_1}\right) - (N - n_l) u\left(\pi_{f,l}^{FL}\right) \right]$$

$$= n_h \left[ u\left(\pi_{c,l}^{PL_1}\right) - u\left(\pi_{f,l}^{PL_1}\right) \right] + n_l \left[ u\left(\pi_{f,l}^{FL}\right) - u\left(\pi_{c,l}^{FL}\right) \right] + N \left[ u\left(\pi_{f,l}^{PL_1}\right) - u\left(\pi_{f,l}^{FL}\right) \right]$$

The 1$^{\text{st}}$ and the 3$^{\text{rd}}$ square brackets are negative and the second is positive.
Note also that $\pi_{c,l}^{FL} = \pi_{c,l}^{PL} = \pi_0 - \gamma_l(N - n_l) - [\pi_0 + 1 - \gamma_l N] = \gamma_l n_l - 1 > 0$.
Hence:

$$\pi_{c,l}^{PL_1} = \pi_{f,l}^{PL_1} = \pi_{c,l}^{FL} = \pi_{f,l}^{FL} \text{ and } u\left(\pi_{f,l}^{FL}\right) - u\left(\pi_{c,l}^{FL}\right) < u\left(\pi_{f,l}^{FL}\right) - u\left(\pi_{f,l}^{PL_1}\right).$$

Therefore $\frac{d\Delta_1}{dp} < 0$.

As $\Delta_1(p = 0) = 0$, $\frac{d\Delta_1}{dp} < 0$ implies $E(U^{FL}) > E(U^{PL})$.

(ii) If $p \approx 1$ then from Propositions 1 and 3, the difference of expected utilities between NL and PL is given by:

$$\Delta_2 = E\left(U^{NL}\right) - E\left(U^{PL}\right)$$

$$= pn_l u\left(\pi_{c,l}^{NL}\right) + p(N - n_l) u\left(\pi_{f,l}^{NL}\right) + (1 - p)n_l u\left(\pi_{c,h}^{NL}\right) + (1 - p)(N - n_l) u\left(\pi_{f,h}^{NL}\right)$$

$$- \left[ pn_l u\left(\pi_{c,l}^{PL_2}\right) + p(N - n_l) u\left(\pi_{f,l}^{PL_2}\right) + (1 - p)n_l u\left(\pi_{c,h}^{PL_2}\right) + (1 - p)(N - n_l) u\left(\pi_{f,h}^{PL_2}\right) \right]$$

where

$$\pi_{c,l}^{NL} = \pi_0 - \gamma_l(N - n_l) = \pi_{c,l}^{PL_2}, \quad \pi_{f,l}^{NL} = \pi_0 + 1 - \gamma_l(N - n_l) = \pi_{f,l}^{PL_2}$$

$$\pi_{c,h}^{NL} = \pi_0 - \gamma_h(N - n_l) = \pi_{c,h}^{PL_2}, \quad \pi_{f,h}^{NL} = \pi_0 + 1 - \gamma_h(N - n_l) = \pi_{f,h}^{PL_2}.$$

As the payoffs of coalition members and non-members are the same under both learning scenarios, for both states of the world, we get $\Delta_2 = 0$ and hence $E(U^{PL}) = E(U^{NL})$.

(iii) If $\gamma_l \approx 1/N$ then $n_l = N$ and $n^{PL} = n_h$, as shown above in the proof of Lemma 2. From Propositions 1 and 3:

$$\Delta_3 = E(U^{PL}) - E(U^{FL})$$

$$= pn_h u(\pi_{c,l}^{PL_1}) + p(N - n_h) u(\pi_{f,l}^{PL_1}) + (1 - p)n_h u(\pi_{c,h}^{PL_1}) + (1 - p)(N - n_h) u(\pi_{f,h}^{PL_1})$$

$$- \left[ pn_l u(\pi_{c,l}^{FL}) + p(N - n_l) u(\pi_{f,l}^{FL}) + (1 - p)n_h u(\pi_{c,h}^{FL}) + (1 - p)(N - n_h) u(\pi_{f,h}^{FL}) \right]$$

where:

$$\pi_{c,l}^{PL_1} = \pi_0 + 1 - \gamma_l N < \pi_{c,l}^{FL} = \pi_0 - \gamma_l(N - n_l) = \pi_0, \quad \pi_{f,l}^{PL_1} = \pi_0 + 1 - \gamma_l N < \pi_{f,l}^{FL} = \pi_0 + 1$$

$$\pi_{c,h}^{PL_1} = \pi_{c,h}^{FL} = \pi_0 - \gamma_h(N - n_h), \quad \pi_{f,h}^{PL_1} = \pi_{f,h}^{FL} = \pi_0 + 1 - \gamma_h(N - n_h)$$

As $\pi_{c,h}^{PL_1} = \pi_{c,h}^{FL}$ and $\pi_{f,h}^{PL_1} = \pi_{f,h}^{FL}$, $\Delta_3$ can be simplified to:

$$\Delta_3 = p \left[ n_h u(\pi_{c,l}^{PL_1}) - n_l u(\pi_{c,l}^{FL}) \right] + p \left[ (N - n_h) u(\pi_{f,l}^{PL_1}) - (N - n_l) u(\pi_{f,l}^{FL}) \right].$$

Using $n_l = N$ we get:

$$\Delta_3 = p\left[n_h\left(u(\pi_{c,l}^{PL_1}) - u(\pi_{f,l}^{PL_1})\right) + N\left(u(\pi_{f,l}^{PL_1}) - u(\pi_{c,l}^{FL})\right)\right].$$

The sign of $\Delta_3$ is negative as $\pi_{c,l}^{PL_1} = \pi_{f,l}^{PL_1}$ and $\pi_{f,l}^{PL_1} - \pi_{c,l}^{FL} = 1 - \gamma_l N < 0$.

This proves that $E(U^{FL}) > E(U^{PL})$.

(iv) If $\gamma_h \approx 1$ then $n_h = 2$ and $n^{PL} = n_l$, as shown above in the proof of Lemma 2. From Propositions 1 and 3:

$$\Delta_4 = E\left(U^{PL}\right) - E\left(U^{FL}\right)$$

$$= pn_l u\left(\pi_{c,l}^{PL_2}\right) + p(N - n_l)u\left(\pi_{f,l}^{PL_2}\right) + (1 - p)n_l u\left(\pi_{c,h}^{PL_2}\right) + (1 - p)(N - n_l)u\left(\pi_{f,h}^{PL_2}\right)$$

$$- \left[pn_l u\left(\pi_{c,l}^{FL}\right) + p(N - n_l)u\left(\pi_{f,l}^{FL}\right) + (1 - p)n_h u\left(\pi_{c,h}^{FL}\right) + (1 - p)(N - n_h)u\left(\pi_{f,h}^{FL}\right)\right]$$

where:

$$\pi_{c,l}^{PL_2} = \pi_0 - \gamma_l(N - n_l) = \pi_{c,l}^{FL}, \ \pi_{c,l}^{PL_2} = \pi_0 + 1 - \gamma_l(N - n_l) = \pi_{f,l}^{FL},$$

$$\pi_{c,h}^{PL_2} = \pi_0 - \gamma_h(N - n_l) > \pi_{c,h}^{FL} = \pi_0 - \gamma_h(N - n_h),$$

$$\pi_{f,h}^{PL_2} = \pi_0 + 1 - \gamma_h(N - n_l) > \pi_{f,h}^{FL} = \pi_0 + 1 - \gamma_h(N - n_h).$$

As $\pi_{c,h}^{PL_2} = \pi_{c,h}^{FL}$ and $\pi_{f,h}^{PL_2} = \pi_{f,h}^{FL}$, $\Delta_4$ can be simplified to:

$$\Delta_4 = (1 - p)\left[n_l u\left(\pi_{c,h}^{PL_2}\right) - n_h u\left(\pi_{c,h}^{FL}\right)\right] + (1 - p)\left[(N - n_l)u\left(\pi_{f,h}^{PL_2}\right) - (N - n_h)u\left(\pi_{f,h}^{FL}\right)\right]$$

Thus, $sign(\Delta_4) = sign\left\{\frac{\Delta_4}{1-p}\right\}$.

Let us rewrite $\frac{\Delta_4}{1-p}$:

$$n_h[u\left(\pi_{f,h}^{FL}\right) - u\left(\pi_{c,h}^{FL}\right)] + n_l[u\left(\pi_{c,h}^{PL_2}\right) - u\left(\pi_{f,h}^{PL_2}\right)] + N[u\left(\pi_{f,h}^{PL_2}\right) - u\left(\pi_{f,h}^{FL}\right)]$$

The 1st and the 3rd square brackets are positive and the second is negative.

Comparing $\pi_{f,h}^{FL}$ and $\pi_{c,h}^{PL_2}$ we get:

$$\pi_{f,h}^{FL} - \pi_{c,h}^{PL_2} = 1 - \gamma_h(n_l - n_h).$$

As $\gamma_h \approx 1$, $\pi_{f,h}^{FL} - \pi_{c,h}^{PL_2} \approx 1 - (n_l - n_h)$.

The sign of this expression is negative as $n_l - n_h \geq 2$.

Thus: $\pi_{c,h}^{FL} < \pi_{f,h}^{FL} < \pi_{c,h}^{PL_2} < \pi_{f,h}^{PL_2}$.

Consequently:

$$|u\left(\pi_{c,h}^{PL_2}\right) - u\left(\pi_{f,h}^{PL_2}\right)| < |u\left(\pi_{f,h}^{PL_2}\right) - u\left(\pi_{f,h}^{FL}\right)|.$$

Hence, $\Delta_4 > 0$, that is, $E(U^{PL}) > E(U^{FL})$.

## Appendix A2: Choice of Parameter Values

### CRRA

1. We use CRRA utility function and payoff function:

$$u(\pi_i) = \frac{1}{1 - \rho}\pi_i^{1-\rho}; \pi_i = \frac{B\gamma_h(N - 1) + x_i - \gamma X}{\gamma_h(N - 1)} \text{ where } B > 1.$$

For the results in Sect. 4 we have used B = 1.1; however to ensure results are robust we have done simulations also for B = 1.5. These results are presented in "Appendix B" available to readers on request to the authors.

2. We choose coefficient of relative risk aversion $\rho = 0.0, 0.05, 0.5, 0.99, 2.50, 5.0, 10.0, 20.0$.
3. The other key parameters are $N, p, \gamma_l, \gamma_h$. We describe the choice of these for the full simulations then briefly say what we do for the simple simulation.
4. We choose 4 random variables: $z_1,\ldots,z_4$ lying strictly between 0 and 1.
5. We will work with the median $\bar{\bar{\gamma}}$ : $2\bar{\bar{\gamma}} = \gamma_l + \gamma_h$. This allows us to choose $p$ independently of $\gamma_l, \gamma_h$, so we set as $p = 0.0001 + z_1 \cdot 0.9998$.
6. We want to ensure that it is possible to have $2 < n_h < \bar{\bar{n}} - 1 < \bar{\bar{n}} + 1 < n_l < N$.
7. This first of all requires that $N$ is greater than 7; and we set an upper limit for $N = 100$. So we set $NN = 7.0 + z_2 \cdot 93; N = I(NN)$.
8. Next we set $\bar{\bar{\xi}} \equiv 1/\bar{\bar{\gamma}}$. To allow for $2 < n_h < \bar{\bar{n}} - 1 < \bar{\bar{n}} + 1 < n_l < N$ we set $\bar{\bar{\xi}}$ to lie in range $[4, N - 3]$; so $\bar{\bar{\xi}} = 4 + z_3(N - 7)$. Then we have $\bar{\bar{\gamma}} = 1/\bar{\bar{\xi}}, \bar{\bar{n}} = I(\bar{\bar{\xi}})$.
9. Finally we choose $\gamma_l, \gamma_h$, ensuring the possibility that $2 < n_h < \bar{\bar{n}} - 1 < \bar{\bar{n}} + 1 < n_l < N$ and that $\bar{\bar{\gamma}}$ : $2\bar{\bar{\gamma}} = \gamma_l + \gamma_h$. We choose the smaller of two intervals: $(1/N, \bar{\bar{\gamma}}), \bar{\bar{\gamma}}, 1)$.

(i) If $\bar{\bar{\gamma}} < 0.5(1 + 1/N)$ choose:

$$\xi_h = (\bar{\bar{\xi}} - 2.0) + z_4(\bar{\bar{\xi}} - 3); \gamma_h = 1/\xi_h; \quad \gamma_l = 2\bar{\bar{\gamma}} - \gamma_h; \quad n_s = [\gamma_s], \ s = l, h$$

(ii) If $\bar{\bar{\gamma}} \geq 0.5(1 + 1/N)$:

$$\xi_l = (\bar{\bar{\xi}} + 2) + z_4(N - \bar{\bar{\xi}} - 2);$$
$$\gamma_l = 1/\xi_l; \quad \gamma_h = 2\bar{\bar{\gamma}} - \gamma_l; \quad n_s = [\gamma_s], \ s = l, h.$$

10. For simple simulations we choose $N = 20$; $p = 0.1, 0.5, 0.9$; $z_3 = 0.1, 0.5, 0.9$; $z_4 = 0.1, 0.5, 0.9$; $z_3$ determines where the median of $\gamma_s$ lies in the relevant range, while $z_4$ determines the spread of values of $\gamma_s$ around that median.

## CARA

$\lambda$ is the coefficient of absolute risk aversion; for any given level of payoff $\pi$ it is related to the coefficient of relative risk aversion by the formula: $\lambda = \rho/\pi$. For $\pi$ we have chosen the average value of the payoff in the FL scenario: $\bar{\pi} = p\bar{\pi}_l + (1 - p)\bar{\pi}_h$ where $\bar{\pi}_s$ is average payoff (averaged over signatories and non-signatories) in state of the world $s$, given by: $\bar{\pi}_s = \pi_0 - (N - n_s)\gamma_s + (N - n_s)/N$, where $\pi_0 = B\gamma_h(N - 1)$. We use parameter values $p = 0.5$ (both states of the world are equally likely); and $1/\gamma_h = 0.5(N + 1) - 0.25(N - 1)$, $1/\gamma_l = 0.5(N + 1) + 0.25(N - 1)$, consistent with our use of the median value of $1/\gamma_s$ being the midpoint $(N + 1)/2$ and $1/\gamma_s, s = l, h$ lying in the mid-points of the subsequent sub-ranges. By taking values of $N$ in a wide range between 25 and 100 we have shown that $\bar{\pi}$ does not vary significantly with $N$ but does vary with $B$ according to the following approximation: $\bar{c} = 3.0 + 4 \cdot (B - 1.1)/1.1$. We have used this to calculate the values of $\lambda$ in Table 5 in the main text.

# References

Arrow K, Fisher A (1974) Environmental preservation, uncertainty and irreversibility. Quart J Econ 88:312–319

Barrett S (1994) Self-enforcing international environmental agreements. Oxf Econ Pap 46:878–894

Boucher V, Bramoullé Y (2010) Providing global public goods under uncertainty. J Public Econ 94:591–603

Bramoullé Y, Treich N (2009) Can uncertainty alleviate the commons problem? J Eur Econ Assoc 7(5):1042–1067

Carraro C, Siniscalco D (1993) Strategies for the international protection of the environment. J Public Econ 52:309–328

Dellink R, Finus M (2012) Uncertainty and climate treaties: does ignorance pay? Resour Energy Econ 34:565–584

Endres A, Ohl C (2003) International environmental cooperation with risk aversion. Int J Sustain Dev 6:378–392

Epstein L (1980) Decision-making and the temporal resolution of uncertainty. Int Econ Rev 21:269–284

Finus M, Caparrós A (2015) Game theory and international environmental cooperation: essential readings. Edward Elgar, Cheltenham

Finus M, Pintassilgo P (2013) The role of uncertainty and learning for the success of international climate agreements. J Public Econ 103:29–43

Gollier C, Jullien B, Treich N (2000) Scientific progress and irreversibility: an economic interpretation of the 'Precautionary Principle'. J Public Econ 75:229–253

Hong F, Karp L (2014) International environmental agreements with exogenous and endogenous risk. J Assoc Environ Resour Econ 1:365–394

Karp L (2012) The effect of learning on membership and welfare in an international environmental agreement. Clim Change 110:499–505

Kolstad C (1996a) Fundamental irreversibilities in stock externalities. J Public Econ 60:221–233

Kolstad C (1996b) Learning and stock effects in environmental regulations: the case of greenhouse gas emissions. J Environ Econ Manag 31:1–18

Kolstad C (2007) Systematic uncertainty in self-enforcing international environmental agreements. J Environ Econ Manag 53:68–79

Kolstad C, Ulph A (2008) Learning and international environmental agreements. Clim Change 89:125–141

Kolstad C, Ulph A (2011) Uncertainty, learning and heterogeneity in international environmental agreements. Environ Resour Econ 50:389–403

Meyer D, Meyer J (2006) Measuring risk aversion. Found Trends Microecon 2:107–203

Na S-L, Shin HS (1998) International environmental agreements under uncertainty. Oxf Econ Pap 50:173–185

Narain U, Fisher A, Hanemann M (2007) The irreversibility effect in environmental decision making. Environ Resour Econ 38:391–405

Rubio S, Casino B (2005) Self-enforcing international environmental agreements with a stock pollutant. SpanEconRev 7:89–109

Rubio S, Ulph A (2007) An infinite-horizon model of dynamic membership of international environmental agreements. J Environ Econ Manag 54:296–310

Ulph A (2004) Stable international environmental agreements with a stock pollutant, uncertainty and learning. J Risk Uncertain 29:53–73

Ulph A, Maddison D (1997) Uncertainty, learning and international environmental policy coordination. Environ Resour Econ 9:451–466

Ulph A, Ulph D (1996) ch. 3: Who gains from learning about global warming? In: van Ierland E, Gorka K (eds) The economics of atmospheric pollution. Springer, Heidelberg, pp 31–62

Ulph A, Ulph D (1997) Global warming, irreversibility and learning. Econ J 107:636–650

Weitzman M (2009) On modeling and interpreting the economics of catastrophic climate change. Rev Econ Stat 91:1–19