# ALTRUISM, THE PRISONER'S DILEMMA, AND THE COMPONENTS OF SELECTION

**JEFFREY A. FLETCHER and MARTIN ZWICK**

Systems Science Ph.D. Program, Portland State University
Portland, OR, USA 97207-0751 (jeff@pdx.edu)

## Abstract

The n-player prisoner's dilemma (PD) is a useful model of multilevel selection for altruistic traits. It highlights the non zero-sum interactions necessary for the evolution of altruism as well as the tension between individual and group-level selection. The parameters of the n-player PD can be directly related to the Price equation as well as to a useful alternative selection decomposition. Finally, the n-player PD emphasizes the expected equilibrium condition of mutual defection in the absence of higher levels of organization and selection.

## Keywords

Altruism, free-rider problem, multilevel selection, n-player prisoner's dilemma, Price equation, tragedy of the commons.

## 1 Introduction

The mechanisms by which altruistic[1] behavior may evolve in biological systems has been vigorously debated over the last several decades. Alternative explanations include reciprocal altruism where the self-interest of individuals is served by the exchange of cooperation with others [2, 14], inclusive fitness where the self-interest of genes is served by benefiting copies of themselves in other organisms (usually relatives) [8], and multilevel selection (often called group selection) where the self-interest of groups may favor those with more altruistic members [13, 15]. Although these explanations have aspects that are mathematically equivalent [6, 13, 15], they clearly differ in their view of the level at which self-interest can select for self-sacrifice.

The purpose of our research is to demonstrate the usefulness of game theory, especially the n-player prisoner's dilemma (PD), as a framework for understanding the evolution of altruistic behaviors. Although computer simulations of the PD including n-player versions have been used to study reciprocal altruism [e.g. 2, see 3 for summary], surprisingly, as far as we know, the n-player PD has not been used explicitly to model multilevel selection. Previously we have shown that an n-player PD with minimal group structure can provide a very simple model of multilevel selection favoring altruism [4, 5]. This model highlights the non zero-sum outcomes needed for selection of altruistic traits as well as the tension between selection within groups (which favors selfish individuals) and selection between groups (which can favor altruistic individuals). Here we briefly review this earlier work and then extend it through use of the Price covariance equation [6, 10] which decomposes selection into within- and between-group components. We relate the parameters of our n-player PD to these components, and a similar but alternative selection decomposition.

## 2 N-Player Prisoner's Dilemma

The n-player PD offers a straightforward way of thinking about the tension between individual and group levels of selection. In real-world biological and social systems the effects of cooperation or defection are often distributed diffusely to other members of a group, i.e., they do not necessarily arise via pair-wise interactions. The n-player PD applies to both problems of conserving a common resource and to problems of equitable contributions towards a common good [7]. An n-player PD involving exploitation of common resources is also known as a "tragedy of the commons" [9], whereas a PD involving contributions is commonly known as "the free-rider problem."

When there is a common and finite resource, each individual benefits by using more than its share of that resource, but when all players apply this individual rationality it leads to a deficient (non-Pareto optimal) and hence collectively irrational outcome. For example, each country that fishes international waters can increase its utility by taking more of this common resource, but as more and more countries overfish, the common stock is depleted beyond the point where it can quickly replenish and

---

[1] Here we use the term *altruistic* to describe any behavior that gives benefit to others at a relative cost to the provider of the benefit. Psychological or moral aspects of altruism are not investigated nor implied by our use of this term. In our model cooperation is equivalent to altruism because the cooperate strategy always involves self-sacrifice.

so in subsequent years all have less. This leads to decreased utility both for countries that overfish (defectors) and those that don't (cooperators).

Tax evasion is a social example of the "free-rider problem." Biological examples include alarm calling [11] and female-biased sex ratios [1]. In alarm calling the contribution is made towards or withheld from a common group security. Alarm calling is altruistic assuming that this behavior lowers the caller's individual fitness as compared to living among alarm callers, but not exhibiting calling behavior oneself. Animals with female-biased sex ratios allow high group growth rates, but individuals that have more even (or male-biased) ratios are free-riders in that they benefit by the high growth rate provided by other group members while increasing their individual fitness by having more progeny of the rarer sex.

A simple payoff scheme for an n-player PD is illustrated by Fig. 1. The horizontal axis specifies the fraction of individuals cooperating for the common good. The vertical axis gives the average utility or fitness to each individual. For convenience, we assume a linear relationship between utility and fraction of cooperators. The upper line denotes the utility for a defector (D) while the lower line is the utility for a cooperator (C). The defectors' line dominates the cooperators' line, i.e., selfish individual behavior always has a higher utility than altruistic behavior no matter what the fraction of altruists.
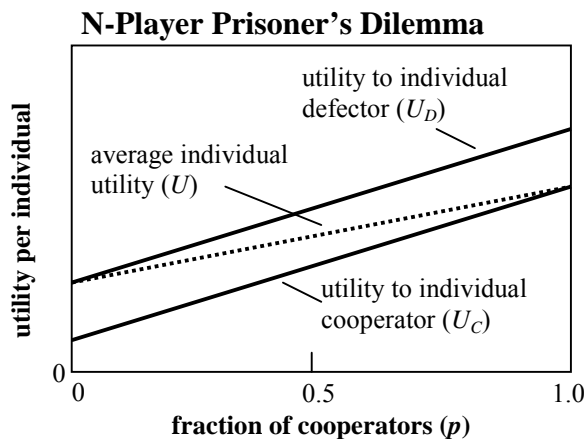
## N-Player Prisoner's Dilemma



Fig. 1. Utility lines for defectors and cooperators as a function of the fraction of cooperators (*p*) for a simple n-player PD. The dashed line indicates the average utility.

The deficient outcome of the PD here inheres in the fact that the utility to defectors when there is a minimum number of cooperators is lower than the utility to cooperators when there is a maximum number of cooperators, i.e. the average utility line has a positive slope. So even though for a given state of the system an *individual* benefits more by defection than cooperation, still cooperators in a *group* of cooperators get more benefit than defectors in a group of defectors.

The "tragedy" (and what makes this a PD) is that whatever the current state of the system, individual rationality or individual selection favors defection which tends to drive the system progressively towards a (boundary) equilibrium state less beneficial to all. This state is a non-Pareto optimal and irrational collective outcome. To summarize algebraically: $U_D(p) > U_C(p)$ (for all $p$) causes $p$ to decrease, but $U_C(1.0) > U_D(0.0)$. The co-parallel lines used here are the simplest of many cooperator and defector utility curves that can satisfy these PD conditions.

It may be tempting to think of cooperation as selfish rather than altruistic because a group of all cooperators gets more utility per individual than a group of all defectors, but this is incorrect and misses the crux of the PD. In the 2-player PD, the players would be better off if they both cooperate rather than both defect, but defecting is still the rational *individual* strategy for each player because the prisoners have no way to coordinate their actions and enforce any agreement to cooperate. Cooperating is disadvantageous no matter what the other player does. So in the absence of guarantees of cooperation by other players, cooperating is truly altruistic—it lowers one's individual utility (fitness) while raising the benefit to others. The same situation holds in the n-player game. Given the absence of binding agreements between players, each player is better off to defect, but benefits others by not doing so in that the system is kept at a state with a higher fraction of cooperators. Of course, this is the dynamic for a single set of players, or for a multi-group system viewed at the intra-group level. As we shall see, at the higher level of organization, i.e. that of the total population which includes two or more groups, cooperators can thrive, at least for a while, despite their inferior individual fitness.

## 3 The Model

In the simplest form of the model there are two groups with no migration between them. These groups initially are the same size and vary only in *p*, the fraction of cooperators and defectors. There are no other strategies besides always-cooperate and always-defect. We follow the fraction of cooperators in each group and across the whole population. In each group, the n-player PD (Fig. 1) is described by utility functions for cooperation and defection, $U_C$ and $U_D$, which are dependent on the fraction of

cooperators in each group. Each line is of the form $mp + b$ where slope $m$ is the same for both lines. Therefore there are two parameters to this n-player PD: the slope and $b_D - b_C$, the difference in the intercept for the defectors' and cooperators' utility lines. For simplicity we set $b_C = 0$ so the difference is $b_D$ which we simply call $b$ ($\geq 0$ in our simulations). The condition for a PD is thus $m > b$. In all runs reported in this paper this condition is satisfied.

We will show that the weighted average utility line, $U(p_i) = p_i(m - b) + b$, in Fig. 1 (which depends on both slope and intercept difference) determines the *between*-group selection force, while the increase in the number of defectors *within* each group is proportional to $b$. At the group level a higher fraction of cooperators confers an advantage. At the individual level defectors have an advantage over cooperators.

In this model, at each timestep the number of cooperators within each group is increased by the number of individuals utilizing this strategy times its utility payoff per individual, and similarly for the number of defectors:

$$C_i' = C_i[1 + U_C(p_i)]$$

$$D_i' = D_i[1 + U_D(p_i)]$$

where $C_i$ and $D_i$ are the number of cooperators and defectors respectively in group $i$, and where primed terms represent values after selection. To aid in comparisons among runs, the population of each group is proportionally scaled back (preserving the ratio of cooperators and defectors) so that the total population size matches the original total. Scaling does not do anything substantive. For convenience we define $C = C_1 + C_2$, $D = D_1 + D_2$, and $N_i = C_i + D_i$.

Because the utility for defectors is always higher than that for cooperators, in the long run defectors will dominate both in each group and across the whole population. Yet while the fraction of cooperators *decreases* within each group, the overall fraction of cooperators in the whole population can *increase*. This seemingly anomalous possibility is known as Simpson's paradox [12] and is key to understanding the multilevel selection viewpoint of the evolution of altruism [13]. The effect is transient without mechanisms for reestablishing variation between groups. These mechanisms and their effects on altruism maintenance are explored thoroughly in [1]. An increase in the fraction of cooperators, $p$, also depends upon initial conditions. Specifically, given that in our model initially $N_1 = N_2 = C = D$, the condition for $p$ to increase overall despite decreasing in each group is:

$$m / b \; > \; N_i^2 \; / \; (C_i - D_i)^2$$

where $i = 1$ or 2 (see Appendix A in 4). This equation also implies that $m$ must be greater than $b$, and thus the PD is a necessary (but not sufficient) condition for the total fraction of cooperators to increase.

## 4 Experiments and Results

We have explored how combinations of our two parameters affect the magnitude and longevity of the Simpson's paradox effect [4, 5]. Fig. 2 shows a typical run with two groups where Simpson's paradox is evident. Here $m = 0.01$ and $b = 0.003$ and the initial conditions are 90 defectors and 10 cooperators in group 1 and 90 cooperators and 10 defectors in group 2. Note that the inequality above is satisfied, namely $0.01 / 0.003 > 100^2 / (90 - 10)^2$, and therefore we predict the resulting initial increase in overall cooperators.
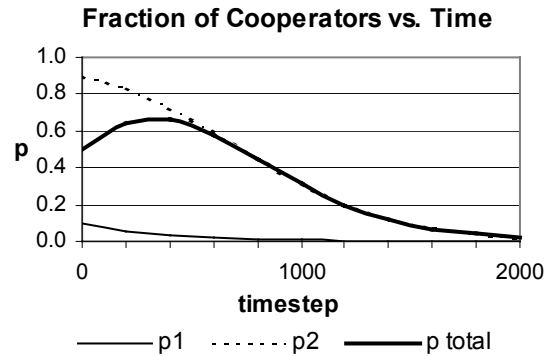
**Fraction of Cooperators vs. Time**



Fig. 2. Fraction of cooperators ($p$) in group 1, group 2, and total for $m = 0.01$ and $b = 0.003$.

Although the fraction of cooperators is decreasing in each group monotonically, the total fraction of cooperators is increasing until timestep 328. The overall increase in fraction of cooperators, despite the decrease within each group, is due to group 2 (cooperator dominated) expanding, while group 1 (defector dominated) is shrinking. After timestep 328 the continued decrease of cooperators in group 2 causes the overall fraction of cooperators finally also to decrease. By timestep 4,000 (not shown) the overall fraction of cooperators, $p$, is essentially zero ($< 0.01\%$). Relatively small modifications to our model can reestablish group variance in cooperator fraction and therefore maintain cooperators indefinitely. For instance, randomly reforming groups periodically with appropriate slope and intercept values allows cooperators to evolve to $p = 1.0$ fixation in our model. The size and number of groups as well as the frequency of group reformation affect this result.

## 5 Price Components of Selection

Price introduced a covariance equation [6, 10] which allows one to partition the change in overall cooperator fraction, $\Delta p = p' - p$, into two components.

$$\Delta p = \text{Cov}(s_i, p_i) / \text{E}(s_i) + \text{E}(s_i \, \Delta p_i) / \text{E}(s_i)$$

$$= \Delta p_b + \Delta p_w$$

where $s_i$ is a measure of group fitness, namely the growth rate of each group, $N_i'/N_i$. The Price equation components can be written as (Appendix A):

$$\Delta p_b = p^* - p$$

$$\Delta p_w = p' - p^*$$

where:

$$p^* = \Sigma(p_i \, N_i') / N'$$

This leads to the following interpretation. In the *between*-group term, $p^*$ plays the role of an idealized $p'$. This idealization predicts a new $p$ by applying the original $p_i$ values to the after-selection group sizes, ignoring $p_i$ changes *within* groups. In the corresponding *within*-group term, the real $p'$ is used, but instead of starting from $p$, we start with $p^*$. The within component thus corrects for the changes of $p_i$ within groups ignored by the idealization. Fig. 3 shows the same data used in Fig. 2 for overall $p$ expressed as a *change* from $t = 0$ along with the two components of selection given by the Price equation.

These two components can also be written in a way that highlights their relationship to the n-player PD utility lines (see appendix B).

$$\Delta p_b = \Sigma C_i(1 + U(p_i)) / \Sigma N_i(1 + U(p_i)) - p$$

$$\Delta p_w = \Sigma C_i(U_C(p_i) - U(p_i)) / \Sigma N_i(1 + U(p_i))$$

where $U(p_i)$ is the average utility line and $U_C(p_i)$ is the cooperators' utility line in Fig. 1:

$$U(p_i) = p_i(m - b) + b$$

$$U_C(p_i) = p_i \, m$$

Note that the numerator in the $p^*$ term in $\Delta p_b$ which gives a predicted change in cooperator number thus depends on the average utility line from the n-player PD. If a control is run in our simulation where both cooperators and defectors are given the average

utility, the resulting change in $p$ during this run matches the Price $\Delta p_b$ decomposition exactly. In the numerator of the within term, which again gives a predicted change in cooperator number, $U_C(p_i) - U(p_i)$ simplifies to $pb - b$ and the $m$ term drops out. So the change in cooperator number predicted by the *within* term depends only on the intercept difference or the relative advantage defectors have over cooperators *within* each group. For this correction term the corresponding control is not possible because we are not making an idealization of the change from the real $p$ as we are in the between term.
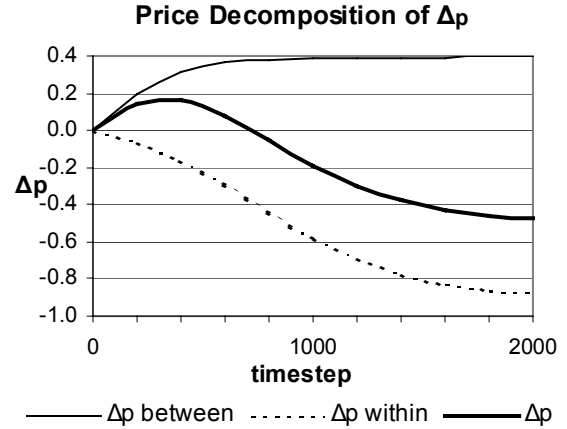
**Price Decomposition of Δp**



Fig. 3. Price decomposition of between- and within-group changes in cooperator fraction ($\Delta p$) for $m = 0.01$ and $b = 0.003$.

## 6 Alternative Decomposition

As we have noted, the Price decomposition makes an idealization about how $p$ changes in the *between* term and then corrects for the difference between $p'$ and this idealization in the *within* term. In this way the Price decomposition is *between-based*. An alternative decomposition is also possible that is *within-based*. That is, one can make a $p'$ idealization in the *within* term and then correct for the difference from the true $p'$ using the *between* term. In this case the idealization for the within term is to assume the fraction of cooperators in each group changes, but that the relative size (fitness) of groups does not. The degree to which group sizes do actually change is then handled as a correction in the between term. We will use $p^\#$ to denote this alternative idealization:

$$p^\# = \Sigma(p_i' \, N_i) / N$$

where the alternative components of selection are:

$$\Delta p_w = p^\# - p$$

$$\Delta p_b = p' - p^{\#}$$

Fig. 4 shows the within- and between-group components for this alternate decomposition. Notice that this gives quite different results than the Price decomposition shown in Fig. 3. In the Price decomposition the between-group selection force on $\Delta p$ rises to 0.4 which given that the starting point is $p$ = 0.5 matches the initial $p_2$ = 0.9, the fraction of cooperators in the cooperator dominated group. Since the total $\Delta p$ goes to $- 0.5$ ($p$ goes to zero) the within-group correction term goes to $- 0.9$. This result says that the equilibrium state is a balance between a strong between-group force (even after group 1 has disappeared) and a strong within-group force (even after all cooperators have disappeared). In contrast, this alternative decomposition of Fig. 4 more intuitively says that the between-group selection force rises as group 2 initially increases over group 1, but that this force goes to zero as the first group disappears. The alternative within-group component steadily decreases to $- 0.5$ to balance the initial $p$ of 0.5 as $p$ goes to zero.
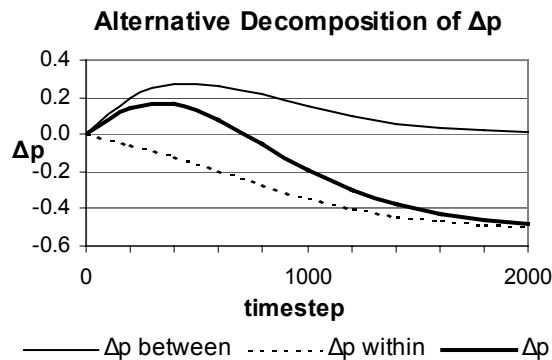
### Alternative Decomposition of Δp



Fig. 4. Alternative decomposition of between- and within-group components of change in cooperator fraction ($\Delta p$) for a run with $m$ = 0.01 and $b$ = 0.003.

The alternative decomposition is not inherently better or worse than the Price decomposition. Which is more appropriate will depend on the situation being studied. In the runs illustrated here where the within-group selection force eventually dominates, the alternative *within-based* decomposition may provide more insight. In situations where the between-group selection force dominates, the Price *between-based* decomposition may be more useful.

Finally note that the within term of the alternative decomposition has the idealized $p'$ and actual $p$ and this allows us to run a control in our model where cooperators are given a utility of $(U_C(p_i) - U(p_i)) / (1 + U(p_i))$ and defectors are given a utility of $(U_D - U(p_i)) / (1 + U(p_i))$. The resulting actual $\Delta p$ for this control run exactly matches the $\Delta p_w$ given by the alternative decomposition.

## 7 Conclusions

The n-player PD is a useful model for studying multilevel selection. It highlights the tension between individual selfishness and the common good that has long been recognized in both biological and social systems. The dominant model of multilevel selection, the Price covariance equation decomposition of selection forces, can be directly related to the parameters and utility lines of a simple n-player PD. In addition, a similar but alternative decomposition may add additional insight for multilevel selection dynamics.

Here we have purposefully only mentioned, but not discussed in detail how mechanisms for reestablishing group variance can preserve the group-level selection force which can counteract individual selection to increase and maintain altruistic behavior. However if between-group selection disappears mutual defection (selfish behavior) dominates. This is the essence of the "dilemma" or "tragedy" and is the expected result in the absence of a higher level of organization and selection pressure.

Therefore when cooperation is observed in nature it is not enough to recognize that cooperation ultimately must provide an individual advantage. This in itself does not imply that cooperation is a result of individual selection. If defectors can get more fitness than cooperators while living among cooperators, then even though both defectors and cooperators may be worse off when defection prevails, mutual defection will still become the equilibrium state. Thinking about the evolution of altruistic behavior in terms of the n-player PD in multiple groups may help reveal when higher levels of selection are operating.

## References

1. L. Avilés, Interdemic Selection and the Sex Ratio: A Social Spider Perspective, The American Naturalist, vol.142, no.2, pp.320-345, 1993.
2. R. Axelrod, The Evolution of Cooperation, New York:Basic Books, Inc, 1984.
3. L.A. Dugatkin, Cooperation Among Animals, An Evolutionary Perspective, New York:Oxford University Press, 1997.
4. J.A. Fletcher and M. Zwick, Simpson's Paradox Can Emerge from the N-Player Prisoner's Dilemma: Implications for the Evolution of Altruistic Behavior, Proceedings of The World Congress of the Systems Sciences and ISSS 2000, Toronto, Canada: International Society for the

Systems Sciences, 2000. (www.sysc.pdx.edu/download/papers/isss_fl_zw.pdf)

5. J.A. Fletcher and M. Zwick, N-Player Prisoner's Dilemma in Multiple Groups: A Model of Multilevel Selection, Artificial Life VII Workshop Proceedings, pp.86-89, 2000 (www.sysc.pdx.edu/download/papers/alife7.pdf).

6. S.A. Frank, Foundations of Social Evolution, Princeton:Princeton University Press, 1998.

7. H. Hamburger, N-Person Prisoner's Dilemma, Journal of Mathematical Sociology, vol.3, pp.27-48, 1973.

8. W.D. Hamilton, The Genetical Evolution of Social Behavior I and II, Journal of Theoretical Biology, vol.7, pp.1-52, 1964.

9. G. Hardin, The Tragedy of the Commons, Science, vol.162, pp.1243-48, 1968.

10. G.R. Price, Selection and Covariance, Nature, vol.277, pp.520-521, 1970.

11. P.W. Sherman, Nepotism and the Evolution of Alarm Calling, Science, vol.197, pp.1246-1253, 1977.

12. E.H. Simpson, The Interpretation of Interaction in Contingency Tables, Journal of the Royal Statistical Society B, vol.13, no.2, pp.238-241, 1951.

13. E. Sober and D.S. Wilson, Unto Others, The Evolution and Psychology of Unselfish Behavior, Cambridge, MA:Harvard University Press, 1998.

14. R.L. Trivers, The Evolution of reciprocal Altruism, Quarterly Review of Biology, vol.46, pp.35-57, 1971.

15. M.J. Wade, Kin Selection: Its Components, Science, vol.210, pp.665-667, 1980.

## Appendix A

Here we show that the Price equation between- and within-group selection components can be expressed in the form $\Delta p_b = p^* - p$ and $\Delta p_w = p' - p^*$ where $p^* = \Sigma(p_i, N_i') / N'$.

The between-group component is:
$$\Delta p_b = \text{Cov}(s_i, p_i) / E(s_i)$$

Using the definition of covariance, we get:
$$E(s_i p_i) / E(s_i) - (E(s_i) E(p_i)) / E(s_i)$$

In the first term denominator we substitute the definition of $s$ and replace the numerator expectation by its summation definition. In the second term the $E(s_i)$ terms cancel.
$$(1/N) \Sigma(N_i s_i p_i) / (N'/N) - p$$

The $N$ terms cancel and $s_i$ is replaced by its definition.

$$\Sigma(N_i (N_i'/N_i) p_i) / N' - p$$

The $N_i$ terms cancel to give:
$$\Delta p_b = \Sigma(N_i' p_i) / N' - p = p^* - p$$

The within-group component is:
$$\Delta p_w = E(s_i \Delta p_i) / E(s_i)$$

Expanding the $\Delta p_i$ term, writing the expectations as summations, and using the definition of $E(s_i)$ gives:
$$1/N \ \Sigma(N_i s_i p_i') / (N'/N) - 1/N \ \Sigma(N_i s_i p_i) / (N'/N)$$

Again canceling $N$ terms and substituting for $s_i$ gives:
$$\Sigma(N_i' p_i') / N' - \Sigma(N_i' p_i) / N'$$

$N_i' p_i' = C_i'$ and $C'/N' = p'$ which gives:
$$\Delta p_w = p' - \Sigma(N_i' p_i) / N' = p' - p^*$$

## Appendix B

Here we show the Price equation between- and within-group selection components expressed in terms of the n-player PD utility lines. First it is useful to show that $N_i' = N_i (1 + U(p_i))$.
$$\begin{aligned} N_i' &= C_i' + D_i' \\ &= C_i (1 + U_C(p_i)) + D_i (1 + U_D(p_i)) \\ &= C_i + D_i + (C_i U_C(p_i) + D_i U_D(p_i)) \\ &= (C_i + D_i) + (C_i + D_i) U(p_i)) \\ &= N_i (1 + U(p_i)) \end{aligned}$$

Now from Appendix A, the between-group component of the Price equation can be written as:
$$\Delta p_b = \Sigma(N_i s_i p_i) / N' - p$$

$N_i p_i = C_i$, and using the definition of $s_i$ and of $N_i'$ gives:
$$\Delta p_b = \Sigma C_i (1 + U(p_i)) / \Sigma N_i(1 + U(p_i)) - p$$

For the within term we start with the following from appendix A:
$$\Delta p_w = \Sigma(N_i s_i p_i') / N' - \Sigma(N_i s_i p_i) / N'$$

The second term and denominators can be rewritten as above and given that $N_i s_i = N_i'$ and $N_i' p_i' = C_i$ we get:
$$\Sigma[C_i' - C_i (1 + U(p_i))] / \Sigma N_i(1 + U(p_i))$$

From the definition of $C_i'$ we get:
$$\Sigma[C_i (1 + U_C(p_i)) - C_i (1 + U(p_i))] / \Sigma N_i(1 + U(p_i))$$

which reduces to:
$$\Delta p_w = \Sigma C_i(U_C(p_i) - U(p_i)) / \Sigma N_i(1 + U(p_i))$$