

COOPERACIÓN HUMANA, RECIPROCIDAD Y CASTIGO. UN ENFOQUE EVOLUTIVO^{1,2,3}

HUMAN COOPERATION, RECIPROCITY, AND PUNISHMENT. AN EVOLUTIONARY APPROACH

Gustavo A. Silva C.^{4,5}

RESUMEN

En este artículo realizo una aproximación evolutiva al problema de la cooperación, con énfasis en la cooperación humana. Presento las diferentes propuestas que desde la teoría de la evolución se han esbozado para explicar el surgimiento del comportamiento altruista y termino esta parte con la postulación de la selección de grupos como única respuesta a tal comportamiento. Seguidamente, muestro cómo estas propuestas tienen sus propias limitaciones para explicar el comportamiento cooperativo humano. Al final del texto, como intención principal, propongo una solución alternativa a la selección de grupos como única explicación para la estabilidad de la cooperación humana mediante los castigos y las recompensas, pero bajo las consideraciones básicas de la selección individual y la racionalidad acotada. Para esto, quiero abordar dicha solución alternativa mediante una interpretación de la propuesta filosófica de Thomas Hobbes acerca del origen del Estado civil. Allí menciono los presupuestos teóricos de Hobbes relacionados con el papel del soberano en la constitución de un Estado. Gracias a este acercamiento, explico el valor filosófico del castigo y las recompensas como estabilizadores no altruistas de la cooperación humana, lo cual supera la solución que apela al mecanismo de la selección de grupos.

Palabras clave: cooperación humana, altruismo, castigos y recompensas, selección de grupos.

ABSTRACT

This article attempts, in general, an Evolutionary Approach to the problem of Cooperation, with special emphasis on Human Cooperation. Throughout the text I present the different proposals from the Theory of Evolution have been outlined to explain the emergence of altruistic behavior, ending with the application of Group Selection as the only response to such behavior. Then I show how these proposals has its own limitations in explaining Human Cooperative behavior. At the end of the text, as my main intention in this article, I propose an alternative to Group Selection as the only explanation for

1 Recibido: 26 de enero de 2015. Aceptado: 4 de abril de 2015.

2 Este artículo se debe citar: Silva, Gustavo A. "Cooperación humana, reciprocidad y castigo. Un enfoque evolutivo". *Rev. Colomb. Filos. Cienc.* 15.30 (2015): 81-121.

3 Agradezco los comentarios e importantes observaciones del evaluador anónimo que dictaminó el presente artículo.

4 Universidad Nacional de Colombia. Correo electrónico: gasilvac@unal.edu.co

5 Bogotá, Colombia.

the stability of Human Cooperation solution. For this, I want to address this alternative through a Philosophical Interpretation of Thomas Hobbes proposed about the origin of the Civil State solution. In this respect, I mention the theoretical assumptions of Hobbes related to the role of the Sovereign in the constitution of a State. Thanks to this approach, I show the philosophical value of Punishment and Rewards as non-altruistic stabilizers Human Cooperation, thus overcoming the solution that appeals to mechanism of Group Selection.

Key words: human cooperation, altruism, punishment and rewards, group selection.

1. INTRODUCCIÓN

Por décadas, la ortodoxia neoclásica de la economía ha supuesto que el comportamiento humano puede explicarse a partir de la presuposición de que los agentes poseen racionalidad perfecta y acceso a información completa, por lo cual, en toda interacción con los demás, pretenden maximizar sus beneficios. Este supuesto, denominado por los economistas como el modelo del *Homo economicus*, sostiene que el único determinante en la acción de las personas es el pago que esperan recibir por dicho comportamiento. Economistas como David Ricardo, Francis Edgeworth y Leon Walras reconocieron en el egoísmo y la racionalidad perfecta la herramienta más fuerte para explicar el comportamiento humano en el intercambio con los demás. Legado recogido de la postura de Adam Smith en su célebre trabajo *The Wealth of Nations*:

It is not from the benevolence of the butcher the brewer, or the baker that we expect our dinner, but from their regard to their own interest. We address ourselves, not to their humanity, but to their self-love, and never talk to them of our own necessities, but of their advantages (1776 19).

Pero si el comportamiento humano ha sido explicado de esta forma por siglos, existen ejemplos recurrentes que nos obligan a preguntarnos por la validez de dichos presupuestos antropológicos. Bajo este marco, ¿cómo se explica el comportamiento cooperativo en los seres humanos?, ampliamente extendido en las innumerables culturas y sociedades desde sus inicios, que además es un rasgo a gran escala que nos define como especie social. La cooperación se constituye en una fuerte objeción y evidente anomalía de la lógica imperante en las ciencias económicas que por siglos han interpretado el proceder humano desde una perspectiva eminentemente individualista, en una competencia del más fuerte.

Sorprendentemente, el comportamiento cooperativo en humanos no solo es paradójico para la explicación ortodoxa de los economistas, también es una

dificultad explicativa dentro de la interpretación clásica de la biología evolutiva. Darwin (1859) fundamenta su teoría de la evolución en el mecanismo de la selección natural que incentiva la competencia entre organismos para la adquisición constante de mejoras adaptativas. Estas mejoras individuales provocan que el organismo logre mayores posibilidades para su supervivencia y reproducción. La guerra del más fuerte la ganará el que más progenie establezca sobre la faz de la Tierra.

This is the doctrine of Malthus, applied to the whole animal and vegetable kingdoms. As many more individuals of each species are born than can possibly survive: and as, consequently, there is a frequently recurring struggle for existence, it follows that any being, if it vary however slightly in any manner profitable to itself, under the complex and sometime varying conditions of life, will have a better chance of surviving, and thus be *naturally selected* (68).

Tanto los presupuestos del *Homo economicus* como los de la selección natural emergen de la lógica individualista de la competencia entre personas u organismos. Sin embargo, esta lógica se revela como limitada cuando se requiere explicar el comportamiento cooperativo de la especie humana que a gran escala y de manera omnipresente se enfrenta a estas dos visiones de profunda influencia en la interpretación de la realidad.

En medio de esta tensión, múltiples experimentos realizados desde la década de 1980 han encontrado que los resultados contradicen el supuesto del *Homo economicus*. Al involucrar a los individuos en juegos como los del Dictador o del Ultimátum se ha encontrado que, en su comportamiento, tales personas no siempre buscan el beneficio último o la maximización de su propio interés (Rosas 2005 90). Como lo sostienen Ferh y Fischbacher (2005) en experimentos enmarcados en el dilema del prisionero (DP) de una sola ronda, con frecuencia la tasa de cooperación se encuentra entre el 40 y 60%. Más aún, sin necesidad de revisar los experimentos desarrollados por los economistas, con confianza podemos asegurar que en la vida diaria, ya no en el laboratorio, encontramos una gran variedad de situaciones en las que las personas parecen comportarse en contra de su propio interés, lo cual genera grandes dudas con respecto a la cobertura explicativa del supuesto canónico de la economía. En ocasiones las personas colaboran con los demás porque cree injusta una situación, porque considera que hacen parte de una comunidad o grupo social que debe colaborar mutuamente o porque en definitiva los seres humanos adscribimos a los demás puntos de vista y opiniones igualmente válidas a las propias que merecen toda consideración, entre otras razones.

El supuesto canónico parece implicar una complicación a menudo innecesaria de un intrincado y difícil cálculo racional para explicar el comportamiento

de las personas en sociedad. En lugar de esto, podemos encontrar situaciones en las que los individuos guían su comportamiento solo por preferencias sociales tales como la reciprocidad, la empatía o la aversión a la desigualdad. También, en ocasiones, como sostienen Elster (1991a) y Ostrom (1998), las interacciones de los individuos pueden deberse a procesos heurísticos que no alcanzan a involucrarse con la pesadez de un cálculo autointeresado ni a presuponer información completa en la situación determinada. Las anteriores razones pueden esgrimirse para explicar situaciones tales como la votación de los ciudadanos, el pago de los impuestos, el respeto a los acuerdos o a contratos incompletos, etc.

Esto no quiere decir que el supuesto canónico de la economía deba desecharse, pues como sostiene Cárdenas (2009) “el *Homo economicus* ha encontrado un muy limitado soporte empírico, excepto en casos de bienes privados puros y en los cuales el mercado funciona bajo los supuestos de la competencia perfecta” (20). En otras palabras, el modelo del hombre eminentemente calculador puede llegar a ser una buena herramienta explicativa en aquellas situaciones de mercado fuertemente competitivo (Ostrom 1998). Sin embargo, deberemos admitir nuevos modelos de explicación que nos permitan alcanzar una comprensión realista del comportamiento humano en sociedad, sin olvidar que aunque pueda que los individuos tengan una tendencia natural a cooperar, es evidente que no siempre siguen dicha tendencia, con frecuencia colaboran de manera condicional con los demás siguiendo normas de reciprocidad que se enmarcan dentro del conjunto de las preferencias sociales. En términos generales, podemos decir con Herbert Gintis (2000) que el éxito de las relaciones sociales depende, en buena medida, de una rica mezcla de virtudes morales y de intereses materiales, algo que el propio Adam Smith develó en sus obras fundamentales.

Lo anterior sugiere una nueva perceptiva de explicación del comportamiento cooperativo humano a partir de la consideración de una racionalidad acotada (Axelrod 1997; Ostrom 1998; Elster 1991b; Cárdenas 2009). Desde este punto de vista, podemos ver que es más factible llegar a explicaciones exitosas si asumimos la asimetría en la información, la introducción de preferencias sociales (Cárdenas 2009) y emociones morales (Frank), que facilitan el empleo de comportamientos heurísticos (Elster 1991b; Ostrom 1998) o de tipo ensayo y error (imitación) (Axelrod 1997) para asimilar normas cambiantes en la realidad que les permite a las personas interactuar de manera satisfactoria y benéfica. Esta perspectiva puede con facilidad enmarcarse, para su formalización, en una aproximación evolutiva que además de dar cuenta de dichos comportamientos puede proponer explicaciones de sus orígenes.

En este artículo realizo una aproximación evolutiva al problema de la cooperación, con énfasis en la cooperación humana. Presento las diferentes propuestas que desde la teoría de la evolución se han esbozado para explicar el surgimiento del comportamiento altruista y termino esta parte con la postulación de la selección de grupos como única respuesta a tal comportamiento. Seguidamente, muestro cómo estas propuestas tienen sus propias limitaciones para explicar el comportamiento cooperativo humano. Al final del texto, como intención principal, propongo una solución alternativa a la selección de grupos como única explicación para la estabilidad de la cooperación humana mediante los castigos y las recompensas, pero bajo las consideraciones básicas de la selección individual y la racionalidad acotada. Para esto, quiero abordar dicha solución alternativa mediante una interpretación de la propuesta filosófica de Thomas Hobbes acerca del origen del Estado civil. Allí menciono los presupuestos teóricos de Hobbes relacionados con el papel del soberano en la constitución de un Estado. Gracias a este acercamiento, explico el valor filosófico del castigo y las recompensas como estabilizadores no altruistas de la cooperación humana, lo cual supera la solución que apela al mecanismo de la selección de grupos.

2. EL CAMINO TEÓRICO DE LA EVOLUCIÓN DEL ALTRUISMO

2.1. Altruismo parental y selección individual

Un gran número de teóricos evolutivos han preferido explicar los comportamientos altruistas o cooperativos de algunas especies guiándose por la selección individual, dejando a un lado la complicada y poco probable selección de grupos. Por ejemplo, William Hamilton sostuvo que en ciertos casos, si no en la mayoría, los comportamientos altruistas o cooperativos en los animales pueden explicarse desde la perspectiva darwiniana de la lógica de la selección individual. En su renombrado trabajo de 1964, *The Genetical Evolution of Social Behaviour*, propone una iluminadora teoría que amplía el concepto darwiniano de éxito adaptativo y muestra que los individuos pueden mejorar su aptitud (definida por su tasa reproductiva), lo cual ayuda a incrementar la aptitud de parientes cercanos que definitivamente comparten sus mismos genes. Denominó esta estrategia *adaptación inclusiva*. Este tipo de comportamiento puede ser calificado como altruista, en cuanto existe un individuo que beneficia a otro, reportando para él mismo un costo directo e inmediato. Sin embargo, indirectamente el éxito reproductivo del actor altruista será efectivo, puesto que la donación de beneficios es efectuada en familiares cercanos

que llevan sus mismos genes. De esta forma el comportamiento altruista de tipo parental evoluciona a través de la adaptación inclusiva que favorece la propagación del rasgo comportamental⁶ en grupos familiares. Las castas estériles en los himenópteros pueden explicarse mediante la teoría de la adaptación inclusiva de Hamilton, pues el particular sistema de apareamiento que reportan hace que las hembras sean genéticamente más parecidas o cercanas a sus hermanas que a sus madres, lo cual facilita que las hembras estériles se dediquen de por vida al cuidado de sus hermanas fértiles.

Hamilton (1963) esbozó su teoría en una regla simple que denominó *coeficiente de parentesco*, que se expresa de la siguiente manera:

$$r \times b > c$$

En donde (*r*) es la probabilidad que tiene un organismo de poseer un gen que codifique el comportamiento altruista; (*b*) representa los beneficios brindados a los parientes del individuo que se comporta como altruista y (*c*) es el costo que se genera para un individuo que se comporta como altruista. Así el costo del altruismo debe compensarse o anularse; por eso el costo (*c*) aparece como algo menor a los beneficios entregados por el altruista a sus parientes. Dichos beneficios se incrementan a mayor grado de cercanía genética o parentesco. En otras palabras, *la regla de Hamilton* expresa que un gen altruista se propaga en una población si se replica en parientes consanguíneos del portador.

En términos de Trivers (1971), la teoría de Hamilton contempla una fuerte relación genética así:

Hamilton (1964) has demonstrated that degree of relationship is an important parameter in predicting how selection will operate, and behavior which appears altruistic may, on knowledge of the genetic relationships of the organisms involved, be explicable in terms of natural selection: those genes being selected for that contribute to their own perpetuation, regardless of which individual the genes appear in (35).

La teoría de selección por parentesco (*kin selection*) que Hamilton elaboró en su famoso ensayo es explotada hasta la saciedad por Richard Dawkins, quien ha sido uno de los más enconados críticos de la selección de grupos y defensor, además, del adaptacionismo a nivel inferior. Dawkins sostiene que dado que los comportamientos altruistas están regidos por el programa que

6 Me atrevo a denominar el altruismo como un rasgo comportamental, apoyado en la equiparación realizada por el propio Hamilton (1963) desde el punto de vista evolutivo alrededor del comportamiento y los rasgos morfológicos: “se admite por lo general que las características comportamentales de una especie son producto de la evolución en la misma medida que su morfología” (354).

los genes le imparten al individuo, es decir, que un gen o un grupo de genes son los responsables de codificar el comportamiento altruista en los individuos, entonces, en este sentido, la selección natural trabaja teniendo en cuenta el nivel genético.

En su extenso trabajo *The Selfish Gene* (1976), Dawkins afirma que un grupo de genes programa el organismo para que su comportamiento se dirija exclusivamente a beneficiar a sus parientes cercanos, en tanto que existe una alta probabilidad de que los parientes cercanos posean el gen o el grupo de genes altruistas también. Pero, a la vez, el gen puede programar al organismo para que beneficie a otros individuos que de alguna manera puedan ser reconocidos como altruistas. Esto obedece a que dichos organismos presentan una muy alta probabilidad de tener genes altruistas y, de cualquier modo, se está beneficiando la reproducción de los genes altruistas.

Se muestra entonces que la selección natural al nivel de genes condiciona el comportamiento altruista para que en últimas el beneficiado de dicho comportamiento sea el gen que lo codifica. Por tal motivo, desde este punto de vista, no puede sostenerse la existencia de un altruismo puro, incondicionado o extremo que de manera indiscriminada otorgue beneficios a cualquier organismo. Esto es algo que el mismo Dawkins sostiene cuando afirma que el gen es eminentemente egoísta, pues busca su beneficio hasta el punto de codificar comportamientos altruistas para su supervivencia.

2.2. Altruismo recíproco

Dawkins (1976) sostiene que los genes pueden codificar comportamiento que reporten beneficios a individuos que están más allá del círculo familiar del actor, algo que Hamilton no había considerado, pero que gracias a su influencia Trivers saca a la luz en 1971. En *The Evolution of Reciprocal Altruism* (1971), Trivers elabora una teoría que pretende extender los alcances del altruismo parental hacia un altruismo recíproco. Este modelo pretende salirse de los límites del parentesco para explicar comportamientos altruistas que a pesar de la teoría de Hamilton se mantuvieron en la oscuridad. Es el caso, por ejemplo, de la llamada de alarma de ciertas aves, el comportamiento involucrado en simbiosis de limpieza entre distintas especies y, por supuesto, el extendido comportamiento cooperativo en humanos.

The model presented here is designed to show how certain classes of behavior conveniently denoted as “altruistic” (or “reciprocally altruistic”) can be selected for even when the recipient is so distantly related to the organism performing the altruistic act that kin selection can be ruled out (35).

Es bueno anotar que el altruismo al que se refiere Trivers es evidentemente limitado o menos puro, pues la cooperación de los individuos está supeditada a la amplia posibilidad de que exista una retribución en un futuro no muy lejano.

Trivers sostiene que hay tres posibilidades o tipos de comportamiento altruista: *a)* cuando el actor dispersa los beneficios aleatoriamente por toda la población (altruismo puro o incondicional), *b)* cuando lo hace no aleatoriamente, considerando el grado de relación genética con los posibles receptores (altruismo parental) y *c)* cuando el altruista los dispersa no aleatoriamente, considerando las tendencias altruistas de los posibles beneficiarios o receptores (altruismo recíproco) (36). Para que se dé el último caso deben cumplirse tres condiciones, definidas por él así:

- i)* deben haber numerosas situaciones de tipo altruista en la expectativa de vida del altruista;
- ii)* un determinado altruista debe actuar repetidamente con el mismo conjunto de individuos, y
- iii)* pares de altruistas deben estar expuestos “simétricamente” a situaciones altruistas para que ambos presten de forma aproximada beneficios equivalentes uno al otro y sufran costos también equivalentes.

Si estas tres condiciones se presentan, sostiene Trivers, el comportamiento altruista del tipo recíproco puede ser seleccionado. En la naturaleza, dichas condiciones se cumplen así:

- i)* la primera condición puede verse representada en la expectativa de vida larga de ciertos organismos. Dada tal expectativa, se maximiza la posibilidad de que dos individuos puedan interactuar más de una vez, y que esta interacción sea una situación altruista que puede repetirse, lo cual genera la posibilidad de la reciprocidad en el futuro;
- ii)* la segunda condición puede cumplirse si observamos la baja tasa de dispersión de ciertas especies que posibilita que un individuo interactúe repetidamente con el mismo conjunto pequeño de individuos vecinos, lo que provoca que el comportamiento recíproco sea factible. Dawkins (1989) denomina a este fenómeno “viscosidad” y explica que es la tendencia que tienen los individuos de continuar viviendo en el lugar de su nacimiento;
- iii)* la tercera condición es un poco más compleja y, según Trivers, podemos encontrarla en casos típicos de interacción no jerárquica en donde no hay relaciones de poder involucradas o, por ejemplo, en casos en donde la jerarquía sí se presenta, pero los que ostentan el poder

pueden requerir la ayuda de los demás, y los que no lo ostentan pueden verse beneficiados de la protección del poderoso. La ayuda en combate resulta ser un buen ejemplo de simetría en relaciones jerárquicas (38).

El modelo de Trivers predice que el beneficio de la reciprocidad en un futuro será mucho mayor que el costo de la colaboración inicial y, evidentemente, el costo de la no colaboración será mucho mayor que el beneficio de mantener un comportamiento egoísta. Según este autor, gracias a la ventaja proporcional del comportamiento altruista recíproco con respecto a la del comportamiento egoísta, y a la posibilidad de que dicho comportamiento recíproco pueda seguir dándose entre dos o más actores, la presión selectiva de la naturaleza favorecerá a aquellos individuos que sean capaces de cooperar discriminadamente con individuos que a su vez también cooperan.

La clave del éxito del altruismo recíproco se fundamenta en la posibilidad de ejercer una interacción discriminada entre individuos, que permita reconocer con certeza a los egoístas y privilegiar las interacciones con los altruistas. Precisamente las tres condiciones del modelo de Trivers pretenden delimitar en un primer momento la interacción discriminada.

Para el caso del comportamiento cooperativo en la especie humana, tanto Trivers y Alexander (1987) como Sober y Wilson (1998) están de acuerdo con que este tipo de interacción involucra mucho más que las tres condiciones expuestas en el modelo del altruismo recíproco. Involucran “cierto nivel de sofisticación cognitiva para percibir las tendencias altruistas de los demás” (Sober & Wilson 115), y las tendencias explotadoras o tramposas, también.

Según afirma Trivers, la selección favorecerá comportamientos que puedan explotar los beneficios del altruismo siempre y cuando dicha actitud pueda quedar impune. Es decir, que si un organismo es capaz de interactuar con un altruista, recibiendo los beneficios que le otorga esta interacción, pero sin dar nada a cambio o, por lo menos, nada equitativo, dicho comportamiento será más apto que el del altruista mismo. Empero, a su vez, si el altruista genera un mecanismo que le permita identificar con precisión a este tipo de explotadores o tramposos, la selección favorecerá este rasgo puesto que dicho mecanismo incrementa el *fitness* del altruista y disminuye el del tramposo. En conclusión, así como la selección individual favorece mecanismos más elaborados para explotar la interacción (hacer trampa), también lo hace con mecanismos, igualmente elaborados, que contrarrestan esta explotación y que permiten identificar al tramposo, así como generar de forma confiable una interacción discriminada. Mecanismos cognitivos como la memoria o la capacidad de distinguir rostros y prevenirse ante la apariencia de un individuo son formas

que los seres humanos, y tal vez algunas especies de simios mayores, emplean para acceder a una interacción discriminada y a una cooperación más estable.

Pero, además de este tipo de mecanismos cognitivos, Trivers sostiene que gracias a la evolución de la cooperación humana, aparecen sentimientos morales como la amistad, la indignación, la culpa o la agresión moral que refuerzan el comportamiento cooperativo. Por ejemplo, en el caso de que un individuo sea identificado o reconocido como tramposo en una situación dada, la culpa y el arrepentimiento que este siente al no cooperar, o mejor dicho, al haber sido identificado como tramposo, le permitirá tener posibilidades nuevas de interacción con individuos altruistas que, gracias a sus sentimientos de amistad, aprecio y perdón al reconocer el arrepentimiento del otro, retomarán la interacción en vista de los beneficios futuros que esta reporta a ambas partes.

2.3. El dilema del prisionero iterado (DIP)

En el artículo conjunto *The Evolution of Cooperation* (1981), Robert Axelrod y William Hamilton retoman la propuesta que presentó Trivers en 1971. Realizan un estudio más profundo desde la perspectiva de la teoría de juegos (que el propio Trivers ya había utilizado en su famoso artículo) y proponen que el altruismo recíproco puede ser ejecutado por organismos que no poseen capacidades cognitivas complejas o sofisticadas, como es el caso de virus y en general de poblaciones microbianas. Inician su artículo refiriéndose al error que se ha cometido al considerar a la cooperación como una adaptación a nivel de grupos, es decir, un rasgo favorecido por un proceso de selección a nivel superior. Al respecto sostienen que “Recent reviews of the evolutionary process, however, have shown no sound basis for a pervasive group benefit view of selection; at the level of a species or a population, the processes of selection are weak” (Axelrod y Hamilton 1390). Su objetivo primordial, empero, tiene que ver con alcanzar una comprensión más rica de la evolución de la reciprocidad, revisando a partir de la teoría de juegos, particularmente del dilema del prisionero (DP), la estabilidad final de una estrategia cooperativa, su robustez y su viabilidad. Razonan que siempre y cuando las interacciones que se dan en la naturaleza sean aleatorias, es decir, con una muy pequeña posibilidad de que las interacciones se repitan entre los mismos actores, es de suponer que en un DP la estrategia más estable ha de ser siempre desertar, pues no hay incentivos para la cooperación, en cuanto la posibilidad de la reciprocidad es mínima.

Apart from being the solution in game theory, defection is also the solution in biological evolution. It is the outcome of inevitable evolutionary trends through mutation and natural selection: if the payoffs are in terms of fitness,

and the interactions between pairs of individuals are random and not repeated, then any population with a mixture of heritable strategies evolves to a state where all individuals are defectors. [...] In these respects the strategy of defection is stable (Axelrod & Hamilton 1992).

Los autores resaltan que existen muchas situaciones biológicas donde las interacciones que llevan a cabo la mayoría de los organismos tienen la forma del DP, pero con una particularidad: en tales casos las interacciones se pueden dar más de una vez y repetidamente, algo que Trivers ya había puesto de manifiesto en las dos primeras condiciones para su modelo de altruismo recíproco. Axelrod y Hamilton utilizan una nueva versión del DP que se acomoda más a estas situaciones, denominada dilema del prisionero iterado (DPI). Este nuevo modelo tiene en cuenta lo que Axelrod designa como la “sombra del futuro”, es decir, la posibilidad de realizar un número indeterminado de repetidas interacciones (jugadas) entre los mismos actores. En principio, estas situaciones biológicas evidentemente son posibles, siempre y cuando los individuos que interactúan posean una capacidad cognitiva algo elaborada que les permita recordar interacciones anteriores y reconocer a individuos con quienes han interactuado con antelación, algo sobre lo que tanto Trivers como Sober y Wilson han llamado la atención. Si estas capacidades cognitivas no están presentes en los individuos que interactúan, por más que dichas interacciones se den varias veces, para cada uno de ellos surgirán como una nueva interacción sin historial ni futuro.

En contra de esta posición, Axelrod y Hamilton sostienen que no es necesario que los organismos tengan cerebro o una capacidad de previsión para que puedan utilizar estrategias de interacción diferencial en la realidad biológica. Las bacterias son un ejemplo de esto:

- (i) Bacteria are highly responsive to selected aspects of their environment;
- (ii) this implies that they can respond differentially to what other organisms around them are doing;
- (iii) these conditional strategies of behavior can certainly be inherited; and
- (iv) the behavior of a bacterium can affect the fitness of other organisms around it, just as the behavior of other organisms can affect the fitness of a bacterium (1992).

Desde este punto de vista, el DPI ha de contemplar, en una forma mucho menos restrictiva, interacciones entre bacterias, por ejemplo, o entre una colonia de bacterias y primates, en tanto que la condición de la previsión para una interacción discriminada puede ser reemplazada por estrategias diversas en organismos menos complejos.

Para que la cooperación surja han de presentarse tres rasgos importantes, y en cierto modo imprescindibles para su selección y permanencia en la naturaleza.

Axelrod y Hamilton afirman, como ya lo mencionamos, que la estrategia de la cooperación debe ser viable, robusta y evolutivamente estable. La viabilidad se refiere a la posibilidad de que una estrategia, en este caso la cooperativa, pueda encontrar oportunidad de concretarse en medio de un ambiente del todo adverso, por ejemplo, un ambiente de egoístas. La segunda característica tiene que ver con que la estrategia de cooperación debe prosperar en medio de un ambiente altamente diversificado de estrategias que compiten con ella. La última característica enunciada fue definida en 1982 por el biólogo inglés Maynard Smith quien, utilizando la teoría de juegos, postuló que una estrategia evolutivamente estable (EEE) es aquella que es adoptada por la mayoría de integrantes de una población y no puede ser invadida o mejorada por ninguna estrategia mutante o alternativa.

Axelrod organizó una competencia computacional entre diversas estrategias para determinar cuál era la más eficiente en el DPI. Resultó que, de las 15 estrategias propuestas por expertos en teoría de juegos, la denominada *Tit for Tat*, elaborada por el psicólogo ruso Anatol Rapoport, fue la ganadora. Esta estrategia se caracteriza por iniciar la cooperación en la primera interacción y copiar el comportamiento previo de su adversario permanentemente. Así cuando el adversario no coopera, la estrategia le indicará al jugador que no coopere en la siguiente jugada y si el adversario lo hace, esto es garantía de que el jugador actuará de la misma manera. Esta forma de “ojo por ojo” permite que la represalia de un jugador estabilice la interacción, en tanto que el adversario no podrá aprovecharse inequitativamente de su contraparte, y, más bien, lo empujará a cooperar para recibir los beneficios de una interacción recíproca cooperativa y prolongada.

Para Axelrod y Hamilton no solamente *Tit for Tat* es eficiente frente a las demás estrategias, sino que es la mejor por ser viable, robusta y evolutivamente estable. El hecho de que haya triunfado frente a las demás estrategias, que por cierto en un segundo torneo fueron 62, demuestra su robustez pues prosperó en un ambiente altamente diversificado de estrategias rivales. Su éxito, como lo sostiene Axelrod, radica en que es una estrategia eminentemente “amable” pues inicia la cooperación y la mantiene a menos que su contraparte no coopere, en cuyo caso toma una posición de retaliación inmediata, para luego cooperar si su oponente así lo hace. Este carácter de “amabilidad” hace que *Tit for Tat* evite el conflicto o, si este llega a darse, evita que se prolongue innecesariamente pues nunca lo provoca al no desertar, a menos que el otro lo haya hecho antes. Por esto mismo, es una estrategia transparente y permite que el otro jugador actúe basado en la confianza de la interacción. También tiene la capacidad de “olvidar” la no cooperación si la contraparte decide volver a cooperar.

Si además de lo anterior se añade la sombra del futuro, es decir la inmensa posibilidad de que existan muchas interacciones en el futuro entre los mismos individuos, se puede colegir que la estrategia *Tit for Tat* es evolutivamente estable. La sombra del futuro les permite a los actores mantener una actitud cooperativa, en cuanto existe la confianza de que las ventajas de la cooperación puedan percibirse por tiempo indefinido. Si dicha interacción tiene un límite de repetición conocido por los actores, estos emplearán la estrategia de siempre desertar (SIEMPRE D) o no cooperar dada la seguridad de que la reciprocidad tiene un término definido y, por tanto, se querrá sacar más provecho desertando antes de llegar a ese límite.

Sin embargo, como lo sostienen Axelrod y Hamilton, *Tit for Tat* no es la única EEE, pues SIEMPRE D también lo es, sin importar la posibilidad de que continúe o no la interacción. En el DP de una sola ronda está demostrado que no cooperar es la estrategia más racional, de la misma manera, en una versión iterada del dilema, SIEMPRE D es evolutivamente estable (Boyd *et al.* 2005; Fehr y Fischbacher 2005) ya que puede sacar gran provecho de la incauta cooperación de otros y, además, aunque tal cooperación no se presente, el jugador que utilice SIEMPRE D nunca perderá.

¿Cómo es posible que *Tit for Tat* pueda surgir en medio de un ambiente dominado por SIEMPRE D? La viabilidad de *Tit for Tat* se ve comprometida en tanto que SIEMPRE D es evolutivamente estable y ampliamente extendida en la naturaleza. Para Axelrod y Hamilton existe una salida al equilibrio generado por SIEMPRE D: la cooperación surgida en pequeños grupos de individuos emparentados. La teoría del parentesco genético podría ser una respuesta a la viabilidad de la cooperación en el mundo natural, dado que dicha estrategia pudo generarse y contraponerse a la fuerza de SIEMPRE D gracias a los lazos genéticos entre pequeños grupos familiares. Una vez aparecida en estos grupos la estrategia cooperativa, sus ventajas adaptativas pudieron extenderse a grupos cada vez más numerosos y con menos individuos relacionados entre sí que, sin embargo, tenían una probabilidad alta de reencontrarse e interactuar. En otras palabras, SIEMPRE D puede ser invadida por *Tit for Tat* gracias al parentesco genético de un pequeño grupo y, luego, gracias a la alta probabilidad de que las interacciones entre dos individuos puedan ser repetidas. Para esto es necesario que exista una especie de “apiñamiento” de la población o, como lo propone Dawkins, se debe presentar viscosidad poblacional para que los integrantes de dicho grupo tengan una probabilidad no trivial de volver a interactuar.

La reciprocidad o altruismo recíproco así alcanzado puede caracterizarse por cuatro elementos, a saber: *a*) presencia de la sombra del futuro o posibilidad de futuras interacciones, *b*) probabilidad de que los individuos con los que tal

vez se interactuará sean cooperadores condicionales, *c*) decisión a cooperar inicialmente con el otro y *d*) rechazo a cooperar si el otro en interacciones previas ha defraudado los actos cooperativos con la deserción.

Fehr, Fischbacher y Gächter (2002) describen el altruismo recíproco como aquel en donde un actor condiciona su comportamiento sobre el conocimiento del comportamiento previo de los otros. En este sentido, el actor altruista recíproco, a diferencia del altruista incondicional, solo ofrecerá su ayuda y colaboración a otros individuos, inclusive a aquellos no emparentados con él, si espera algún tipo de beneficio en el futuro.

Hay que anotar, por lo demás, que el modelo descrito por Axelrod (1984) tiene como principal característica el hecho de que se desarrolla en interacciones bipersonales, tal como está diseñado el clásico DP. Por otro lado, y como una extensión de lo hecho por Axelrod y Hamilton (1981), múltiples estudios han concluido que las condiciones en las que el altruismo recíproco o la reciprocidad puede evolucionar se relacionan, a lo sumo, con la presencia de grupos pequeños en donde la viscosidad poblacional es alta (Joshi 1987; Bendor & Mookherjee 1987; Hirshleifer & Rasmusen 1989; Boyd & Richerson 1989-1992; Yamahashi & Takajashi 1994; Ostrom 1998; Fehr & Fishbacher 2005).

3. PROBLEMAS PARA LA ESTABILIDAD DE LA RECIPROCIDAD

Hemos visto cómo se han generado diversas propuestas desde las ciencias biológicas, pasando por la economía hasta la teoría de juegos, de modelos que explican la evolución de la cooperación. Sin embargo, gracias a evidencias experimentales y de la vida real se ha encontrado que cada uno de estos modelos tiene sus propias limitaciones para explicar el comportamiento cooperativo en los seres humanos. Asimismo, como ya se mencionó, todos los modelos están basados en interacciones bipersonales que no reflejan de forma adecuada la vida social en grupos en donde las relaciones se entretajan a través de múltiples interacciones simultáneas. Aunque es correcto afirmar que en interacciones bilaterales repetidas la reciprocidad puede llegar a ser una estrategia evolutivamente estable, también es cierto que cuando se pasa al nivel de interacciones *n*-personales dicha estrategia decae frente a la aparición de oportunistas que desertan aprovechándose de los beneficios otorgados por los reciprocadores (Boyd *et al.* 2005; Fehr & Fischbacher 2005). Estos oportunistas han sido denominados *free riders* en la vasta literatura sobre el tema.

Una posible explicación de este fenómeno se relaciona con que en grupos grandes es más probable el anonimato de sus miembros, dando lugar a que los incentivos materiales para la deserción dominen el comportamiento de

algunos individuos al poder actuar de forma injusta o en contra del bien público, al explotar la colaboración de los demás sin que sean detectados o con una muy baja probabilidad para esto. Los incentivos para el oportunista se potencializan cuando se presenta una situación en donde es imposible o muy difícil limitar el acceso al bien común o público (Fehr & Fischbacher 2005) y, además, existe la configuración contextual que permite cierto grado de anonimato en las acciones de los individuos que intervienen en la provisión, mantenimiento o aprovechamiento del bien.

Dado que, como lo afirman Fehr, Fischbacher y Gächter (2002), en las sociedades humanas parece ser la regla y no la excepción que se presenten permanentemente incentivos materiales para desertar, una estrategia como la reciprocidad está enfrentada a constantes amenazas de los *free riders*, hasta el punto de que a lo largo del tiempo dicha estrategia decaerá y cederá su posición a la desertión o la explotación.

Lo anterior ha sido demostrado en abundantes experimentos que evidencian que a largo plazo la estrategia del altruismo recíproco en juegos de bienes públicos o comunales decae de manera inevitable (Boyd 1988; Fehr & Fischbacher 2005). Esto se debe principalmente a que durante las repetidas interacciones las personas van adquiriendo más información de sus contrapartes, dejando al descubierto posibles desertiones pasadas que, dado el comportamiento recíproco, van minando la estabilidad en la cooperación hasta el punto de decaer por completo bajo la presión de los explotadores u oportunistas.

4. CASTIGO COMO SOLUCIÓN A LA INESTABILIDAD DE LA COOPERACIÓN POR RECIPROCIDAD

Como lo sostienen Boyd y colaboradores (2005), se ha encontrado un mecanismo evolutivo que puede dar solución al rompecabezas configurado por la reciprocidad, el oportunismo, el tamaño de los grupos sociales, el anonimato que este incita y la débil sombra de futuro que en ocasiones se presenta en la interacción humana. Los investigadores han denominado dicho mecanismo *castigo altruista*, definido como la disposición a castigar a aquellos que explotan la cooperación de los demás, forzándolos a cooperar para evitar futuros castigos, incluso si dicho castigo representa un determinado costo para el que lo ejecuta, de tal manera que, frente a individuos que no castigan y que no son castigados, aquellos que sí castigan poseen una desventaja adaptativa —precisamente este es el carácter altruista del castigo— (Boyd 1988; Boyd & Richerson 1992; Fehr, Fischbacher & Gächter 2002; Fehr & Fischbacher 2005; Gintis *et al.* 2005; Hauert *et al.* 2007; Gintis 2008). Investigadores

como Hernich y colaboradores (2006) o Boehm (1993) han encontrado pruebas relacionadas con que el castigo es una práctica extendida en múltiples y diversas sociedades. Los castigos pueden ir desde señalamientos, multas, daños físicos, suspensión de beneficios sociales, hasta el ostracismo a personas que comenten delitos o faltas contra el grupo social (Boehm 1993; Bowles & Gintis 2003; Fehr & Fischbacher 2005). Un ejemplo de señalamiento y ostracismo es el referido por Fehr y Fischbacher (2005): “During World War I, British men who did not volunteer for the army faced strong public contempt and were called ‘whimps’ ” (171).

Se ha encontrado que comportamientos de este orden pueden incrementar la tasa de cooperación en grupos en donde inicialmente hay una frecuencia alta de no cooperadores (Boyd & Richerson 1992). Esto se debe principalmente a que las ganancias que obtienen los *free riders* obtenidas por explotar la cooperación de los demás se ve altamente disminuida al recibir un castigo por su comportamiento, lo cual los deja en desventaja adaptativa frente a aquellos que cooperan mutuamente. De esta forma, dado que los castigadores también incurren en costos para implementar el castigo a los *free riders*, la estrategia con mejores resultados será la de cooperar y abstenerse de castigar.

Supongamos que tenemos una población de individuos en donde se presentan las tres estrategias, a saber: cooperadores (C), *free riders* (F) y castigadores (P). Los (C) obtendrán un beneficio $b = 3$ por cada vez que interactúen con otros (C) o con otros (P) –no olvidemos que estos últimos también son cooperadores–, pero cuando interactúen con (F) estos, los *free riders*, obtendrán un beneficio de explotación $b_e = 6$ y los (C) no obtendrán ningún beneficio por haber sido explotado $e = 0$. De la misma forma, los (P) que interactúen con otros (P) o con otros (C) obtendrán un $b = 3$, pero cuando interactúen con (F) no obtendrán ningún beneficio; por el contrario, incurrirán en un costo por infringir castigo a los (F) de $c_p = -2$, mientras que a los (F) se les infringirá un castigo que representa un costo $c = -4$. Por último, si dos (F) interactúan, ambos se explotarán y, por tanto, no obtendrán ningún beneficio por haberse explotado mutuamente $e = 0$.

Si todos interactúan con todos, podemos hallar los siguientes resultados:

Tabla 1.

Resultados en interacción indiscriminada de cooperadores (C), <i>free riders</i> (F) y castigadores (P)
$R(C) \ b + e + b = 6$
$R(F) \ b_e + e + c = 2$
$R(P) \ b + c_p + b = 4$

Teniendo en cuenta esta situación hipotética que resume, en buena medida, los modelos elaborados por los investigadores para demostrar el poder del castigo altruista en el incremento de las tasas de cooperación (Boyd & Richerson 1992; Fehr, Fischbacher & Gächter 2002; Fehr & Gächter 2002; Boyd *et al.* 2005; Fehr & Fischbacher 2005; Sripada 2005; Gintis 2008; Axelrod 2007; Rosas 2008), podemos decir que la estrategia que obtiene mejores resultados de una interacción aleatoria es la cooperación (C), puesto que aunque puede llegar a ser explotada, por lo menos, no incurre en costos, dándole a sus resultados valores muy por encima de los otros. Al introducir el castigo altruista, una población con esta configuración puede derivar, con interacciones repetidas, en una composición alta de cooperadores y baja de desertores o *free riders*. La cooperación, entonces, puede ser fomentada por la aparición del castigo altruista (Boyd & Richerson 1992; Henrich & Boyd 2001; Boyd *et al.* 2005; Gintis *et al.* 2005; Rosas 2007). Desde otra perspectiva, si es castigado, el *free rider* tiene incentivos para dejar la explotación, de tal suerte que su mejor elección será evitar el castigo mediante un comportamiento cooperativo. En conclusión, el castigo puede ser un mecanismo que refuerza el seguimiento de la norma de reciprocidad e incrementa la tasa de cooperación en un grupo social.

Simulaciones de juegos de bienes públicos (DP de n-personas) realizadas por computador han encontrado que la cooperación declina a medida que transcurre el tiempo, pero que al introducir mecanismos de castigo a los desertores la tasa de cooperación se mantiene (Fehr & Gächter 2002). En la actualidad, para la gran mayoría de investigadores interesados en el tema, el comportamiento caracterizado por estar dispuesto a cooperar en interacciones sociales y a castigar la no cooperación, incluso si es necesario a costa de su propio beneficio, se conoce como *reciprocidad fuerte* (Sethi & Somanathan 1996; Fehr, Fischbacher & Gächter 2002; Bowles & Gintis 2003; Fehr & Henrich 2003; Gintis *et al.* 2005; Falk & Fischbacher 2006).

5. DILEMAS DE SEGUNDO ORDEN EN JUEGOS DE BIENES PÚBLICOS Y COMUNALES

Ahora bien, si revisamos con cuidado los resultados de la Tabla 1, podemos ver que los valores del comportamiento castigador (P) no son los mejores, dado que debe incurrir en costos para castigar a los *free riders* (F). Esta situación hace que, en términos evolutivos, la mejor estrategia sea la de cooperar (C). Desde este punto de vista, no solo los *free riders* castigados tendrán un incentivo negativo para comportarse de forma cooperadora, sino que también los castigadores (P) tendrán incentivos para dejar de castigar y simplemente dedicarse a cooperar, pues los resultados de los cooperadores no castigadores

son mejores que los demás. Así, la institución del castigo, que incrementa la cooperación en grupos con desertores, se convierte en un bien público de segundo orden (Boyd & Richerson 1992; Sober & Wilson 1998; Rosas 2007), donde el comportamiento que evita castigar por no incurrir en costos puede explotar el comportamiento castigador, generando de esta forma un dilema social de segundo orden (Oliver 1980; Yamagishi 1986; Axelrod 1986). Esto provocaría que se disminuyeran los castigadores, lo cual incrementaría correlativamente la tasa de desertores y, en consecuencia, haría decaer la cooperación dentro del grupo social y, como lo dice Axelrod (1986), destruiría toda moderación que se hubiera generado al inicio de la interacción.

Estudios experimentales (Axelrod 1986) han demostrado que el castigo altruista, aunque puede mantener la cooperación en largos periodos de tiempo, no es capaz de hacer que esta se convierta en un comportamiento evolutivamente estable, dado que a medida que los castigadores infringen sanciones a los no cooperadores, estos últimos van disminuyendo su frecuencia en la población, provocando, no solamente un incremento en la tasa de cooperación, sino que también una disminución de la población de los propios castigadores dada la escasez de *free riders*. Al disminuir la población de castigadores a niveles críticos, los incentivos para desertar y explotar la cooperación vuelven a ser importantes y de nuevo la tasa de free rider puede incrementarse.

La simulación realizada por Axelrod (1986) fue desarrollada de tal manera que grupos de jugadores interactuaban en dos periodos distintos bajo el marco de un juego de bienes públicos. En el primer periodo los individuos tenían la posibilidad palmaria de cooperar o no en la interacción con otros y en el segundo, aquellos que habían cooperado en el primer periodo tenían la posibilidad de castigar (incurriendo en un costo) a los individuos que no habían cooperado en la interacción previa. Lo que halló el autor en esta interesante simulación fue que el castigo puede incrementar las tasas de cooperación en donde existen desertores (“comportamiento audaz” como lo denominó Axelrod), lo cual ha sido demostrado en otros estudios (Boyd 1988; Boyd & Richerson 1992; Fehr, Fischbacher & Gächter 2002; Fehr & Gächter 2002; Fehr & Fischbacher 2004, 2005; Boyd *et al*, 2005; Sripada 2005; Gintis 2008; Axelrod 2007; Rosas 2008). También encontró que después de cien interacciones la estrategia de no castigar y solamente cooperar explota y disminuye la posibilidad de cooperar. Se puede colegir, con estos datos, que ser vengativo, castigar en nuestros términos, tampoco es una estrategia evolutivamente estable pues está expuesta, como cualquier bien público o comunal, a la explotación de los *free riders* en un nivel superior.

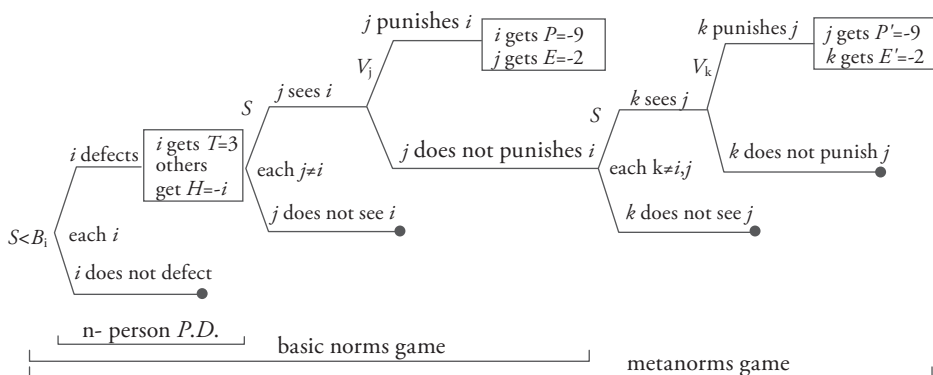
El propio Axelrod propone una solución a este dilema de segundo orden. En el mismo artículo de 1986, sostiene que una *metanorma* puede hacer que el

sistema de castigo alcance la autovigilancia, de tal suerte que el comportamiento castigador no decaiga en el transcurso del tiempo debido a la deserción de los castigadores hacia la mera cooperación. Según este autor, la metanorma no es otra cosa que un comportamiento que está dispuesto a castigar no solo a los no cooperadores, sino además a los que no castigan la deserción. En otras palabras, la *metanorma* es representada por un comportamiento que también castiga a los no castigadores. Boyd y Richerson (1992) la llamaron con posterioridad comportamiento de “retribución” o “estrategia moralista”. Este nuevo mecanismo se enmarca también en el denominado comportamiento de *reciprocidad fuerte* de los seres humanos.

Axelrod (1997) cita algunos ejemplos de castigos a no castigadores:

Las metanormas se han usado en los sistemas de denuncia en las sociedades comunistas. Cuando las autoridades acusan a alguien de hacer algo malo, los demás son impulsados a denunciar al acusado. No unirse a esta forma de castigo es tomado en sí mismo como una deserción contra el grupo (Axelrod 1997 72).

Axelrod (1986, 1997) cree que esta es una de las formas como se ejecuta una norma, a saber, no solamente castigando a aquellas personas que no la siguen, sino también a quienes no hacen respetar, mediante el castigo, el seguimiento de dicha norma. Este investigador propone un modelo para estudiar la estabilidad de la cooperación, el castigo altruista y la metanorma bajo la perspectiva de la teoría de juegos evolutivos. El diagrama de la simulación que llevó a cabo Axelrod (1986) es el siguiente:



Key: i, j, k individuals | S probability of a defection's being seen by any given individual
 B_i boldness of i | V_j vengefulness of j | T temptation to defect | H hurt suffered by others
 P cost of being punished | E punisher's enforcement cost
 P' cost of being punished for not punishing a defector
 E' cost of punishing someone for not punishing a defection

Figura 1. Diagrama de metanorma en teoría de juegos evolutivos (Axelrod 1986 1101).

En este diagrama (o árbol de decisión) se tiene en cuenta tanto la posibilidad o no de cooperar que posee un individuo, o la posibilidad que tiene de castigar la deserción, como la posibilidad que se presenta de castigar al que no ha castigado la deserción.

Axelrod (1986) llevó a cabo una simulación de cinco partidas, cada una con veinte jugadores y una duración por partida de cien generaciones. Según los resultados, al inicio de las interacciones el nivel de vengatividad (castigo) se incrementa y disminuye correlativamente el nivel de audacia (no cooperación) hasta umbrales muy bajos. Sin embargo, los jugadores mantienen niveles altos de vengatividad, aunque los *free riders* sean escasos. Se puede decir que si un cooperador detecta una deserción y no la castiga, este comportamiento provoca que otro cooperador tampoco castigue una deserción, y así hasta que el efecto cascada termine por desplazar la cooperación, instaurando de nuevo un ambiente altamente egoísta dentro del grupo. La *metanorma* puede bloquear el efecto cascada al controlar las deserciones de segundo orden (no castigar a los desertores) con la amenaza de un castigo por no castigar. Esto se debe a que los jugadores permanecen atentos a castigar la deserción en cualquier momento del juego, pues siempre quieren evitar el castigo por no sancionar una deserción detectada.

Para Axelrod (1986), este mecanismo provoca que el sistema sea autovigilante y hace de la norma de reciprocidad y, por tanto, de la cooperación, una estrategia evolutivamente estable. En la estabilidad de dicha estrategia también están de acuerdo Boyd y Richerson (1992).

Estos dos investigadores sostienen que hay dos formas cualitativamente distintas en las que la estrategia castigadora de segundo orden (el castigo a no cooperadores) puede reforzar la cooperación. En primer lugar, si los beneficios a largo plazo de la cooperación de un individuo que castiga son mayores que los costos en que incurre al aplicar la sanción, entonces las estrategias de cooperar y castigar a los no cooperadores, de cooperar solamente si lo han castigado y si no siempre desertar y, en ocasiones, de cooperar pero no castigar coexisten en un equilibrio oscilatorio permanente, equilibrio polimórfico. Es decir, cada una de esas tres estrategias llega a turnarse el liderazgo en el número de individuos que la ejecutan dentro del grupo en largos periodos de tiempo. Pero, como ellos mismos sostienen el equilibrio se caracteriza por una mayoría de desertores y una minoría de las demás estrategias; la tasa de cooperación se mantiene en niveles bajos (Boyd & Richerson 1992 181). En segundo lugar, si el costo de ser castigador es suficientemente alto, la estrategia “moralista” o la metanorma, como la denomina Axelrod (1986), puede ser evolutivamente estable dada la complementariedad de fuerzas entre dicha

estrategia, el castigo altruista de primer orden y la selección de grupos (Boyd & Richerson 1992 172-173).

Como lo sostienen tanto Boyd y Richerson (1992) como Axelrod (1986) lo importante de este mecanismo es que debe ser complementario con el del castigo altruista a los desertores, pues, si no se ligan ambos castigos la posibilidad de que un individuo deje de ser vengativo o castigador es alta, al ver que otros no han castigado una deserción (efecto cascada). Sin embargo, el modelo construido por Boyd y Richerson (1992) adolece, aunque con menor intensidad, del mismo problema presentado por los modelos de reciprocidad y castigo altruista de primer orden: el tamaño de los grupos incrementa el costo del castigo, lo cual fortalece los incentivos para los *free riders* en todos los niveles de la cooperación. Según estos dos autores, los modelos clásicos de castigo de primer orden tienen el problema de que al incrementar el tamaño del grupo también aumenta de forma exponencial la expectativa en el número de interacciones necesarias para alcanzar la cooperación, pues el costo del castigo también es exponencialmente mayor, mientras que en su modelo, en donde se tiene en cuenta el castigo de segundo orden, este incremento es solo lineal.

Con todo, la *reciprocidad fuerte*, que incluye el castigo de primer y segundo orden, requiere para su postulación de la existencia efectiva de una fuerza evolutiva a nivel de grupos que pueda explicar el hecho de por qué ciertos individuos asumen costos individuales que benefician a los demás en interacciones sociales (Boyd & Richerson 1992; Fehr, Fischbacher & Gächter 2002; Gintis *et al.* 2005). Así, en cualquier caso los comportamientos castigadores estudiados hasta el momento no dejan de ser, desde el punto de vista biológico, comportamientos altruistas.

El castigo altruista involucrado en la *reciprocidad fuerte* de los individuos se convierte en un nuevo rompecabezas, pues ¿cómo es posible que ciertos individuos se involucren en comportamientos de castigo a no cooperadores o a no castigadores, a pesar de que dichos comportamientos les reportan costos que los pueden poner en desventaja con otros? Algunos autores han intentado responder a este cuestionamiento defendiendo un mecanismo de la selección natural, que aún hoy en día es algo controvertido, la selección de grupos (Boyd & Richerson 1992; Sethi & Somanathan 1996; Fehr & Gächter 2002; Fehr & Fischbacher 2005; Gintis *et al.* 2005; Gintis 2008).

La selección de grupos se esgrime como razón por la cual un comportamiento evidentemente maladaptativo para un individuo, como el de incurrir en costos para castigar la deserción de otro y, a su vez, beneficiar a los integrantes de su

grupo, es viable en la naturaleza, específicamente en la interacción humana. Boyd y Richerson (1992) argumentan que los grupos con gran número de castigadores tienen altas tasas de comportamiento cooperativo y que esta situación representa ventajas frente a otros grupos que tienen bajos números de castigadores. Por tanto, la frecuencia del castigo y el comportamiento cooperativo es positivamente correlacionado a través de los grupos. En otras palabras, el castigo, como una “respuesta correlacionada” con la cooperación, favorece a los grupos, pues a mayor cooperación mayor *fitness* de grupo (Boyd & Richerson 1992 217). Los grupos con mayor cantidad de cooperadores entonces prevalecerán en la competencia natural frente a aquellos con una baja tasa de cooperación. Todo esto se debe, según los investigadores, a que la selección individual dentro de los grupos es débil frente a la fuerza selectiva de la competencia entre grupos, lo cual provoca que el comportamiento castigador pueda llegar a ser común dentro de los grupos.

En la última parte de este artículo intento ofrecer argumentos en favor de la estabilidad del castigo en la cooperación humana sin acudir a la explicación de la selección de grupos. Bajo los presupuestos metodológicos de la racionalidad acotada, presento una interpretación adecuada de las sanciones e incentivos como dispositivos que fomentan la cooperación en los seres humanos.

6. SELECCIÓN DE GRUPOS COMO EXPLICACIÓN AL ALTRUISMO

Para Sober y Wilson (1998) es muy posible que el comportamiento altruista, señal de cooperación en los organismos biológicos, definido en la biología evolucionista como el aumento de la aptitud de los demás y la disminución de la aptitud del actor, haya podido presentarse gracias a la fuerza de la selección natural que se ejerce entre los grupos y no dentro de ellos. En otras palabras, gracias al cumplimiento de ciertas condiciones naturales, la fuerza de la selección al nivel de los grupos puede superar la fuerza al nivel individual, provocando la permanencia de un rasgo o conjunto de rasgos que reporten un costo para el individuo pero que, a su vez, presenten un beneficio para el grupo. En este sentido, los autores defienden un pluralismo o mecanismo multinivel en los procesos de selección natural y, bajo esta perspectiva, un rasgo altruista puede prosperar aunque sea deletéreo para el individuo, pues puede ser beneficioso y adaptativo para el grupo.

Pareciera que un grupo que posee muchos individuos altruistas tiene baja aptitud frente a uno en donde predominan los individuos egoístas. Esto porque, según Sober y Wilson, a menudo se considera que el grado de aptitud

de un grupo se obtiene por el promedio de aptitud de los individuos que lo integran. Si esto fuera así, el promedio de aptitud de un grupo con muchos individuos altruistas, que por tener este rasgo, tiene baja aptitud individual, sería también bajo. Sin embargo, los autores sostienen que un grupo con estas características puede tener una aptitud global mayor a pesar de tener individuos con una aptitud relativa menor en su interior. La razón de esta conclusión es que este es el típico caso de la *paradoja de Simpson* o la *falacia del promedio* que indica que el promedio relativo de los individuos integrantes de un grupo no refleja acertadamente el promedio global del grupo, puesto que existe una contribución desigual a la media global del grupo; en tanto que los altruistas, aunque menos aptos, reportan a los egoístas y a los demás altruista beneficios y, por tanto, puntos adicionales de aptitud que no se ven reflejados en el promedio relativo, sino en la media global del grupo.

Para que el altruismo prospere en un grupo con egoístas que reciben beneficios de los altruistas y que, además, no incurren en ningún costo adicional por ayudar a otros, afirman Sober y Wilson, existen ciertas condiciones que deben cumplirse:

1. “Debe haber más de un grupo, una población de grupos” (11), pues la selección de grupos se da cuando hay competencia entre ellos.
2. “Es necesario que los grupos sean *distintos* en cuanto a proporción de altruistas” (11), de tal suerte que cada grupo aporte o contribuya de forma distinta a la media global.
3. “Debe existir una relación directa entre la proporción de altruistas en el grupo y el rendimiento del grupo” (11).
4. “Aunque por definición los grupos están aislados entre sí, en cierto sentido no lo están (la progenie de ambos grupos debe mezclarse, o bien competir en la formación de nuevos grupos)” (11).

En el modelo de Sober y Wilson esta es la única posibilidad de que las ventajas relativas de los altruistas superen las ventajas relativas de los egoístas. Es decir, solamente al conjugarse la aptitud relativa de cada grupo en un solo grupo más grande se puede identificar el aporte relativo de cada uno y, sobre todo, la ventaja que los altruistas le aportan a la aptitud global del grupo. De otra manera, al interior de cada grupo las ventajas relativas de los altruistas no podrán superar jamás las ventajas relativas de los egoístas. En conclusión, los dos autores argumentan que si se dan estas condiciones en un medio natural y si, además, se tiene cuidado de no caer en la *falacia del promedio*, se podrá identificar cómo el diferencial de aptitud de los grupos es más fuerte que el diferencial de aptitud de los individuos para la permanencia de rasgos como el comportamiento altruista.

7. CONTROVERSIA EN TORNO A LA SELECCIÓN DE GRUPOS

Desde su postulación inicial en 1962 por Vero C. Wynne-Edwards, la teoría de la selección de grupos ha sido objeto de constante debate. Algunos defensores como el filósofo australiano Kim Sterelny (2001, 2003) han extendido esta propuesta hacia la evolución de la cognición y la cooperación humana.

El impulso a la selección de grupos, y en general a la selección multinivel, se debe a la fundamental definición del llamado “esqueleto lógico de la selección natural” que propuso Richard Lewontin en su artículo de 1970, *The Units of Selection*. Allí, sostiene que para que exista evolución por selección natural deben cumplirse tres requisitos:

1. Different individuals in a population have different morphologies, physiologies, and behaviors (phenotypic variation).
2. Different phenotypes have different rates of survival and reproduction in different environments (differential fitness).
3. There is a correlation between parents and offspring in the contribution of each to future generations (fitness is heritable) (1).

Afirma entonces que la variación heredable es eficacia biológica. Pero, además de esto, resalta que cualquier tipo de entidad que cumpla con estos tres requisitos puede ser denominada “unidad de selección”. En otras palabras, los requisitos lógicos para que la selección natural se dé podrían ser cumplidos teóricamente por entidades de cualquier nivel, desde los genes hasta los superorganismos o grupos de organismos como lo proponen Sterelny y Griffiths (1999).

Debe anotarse, sin embargo, que la teoría de la selección de grupos aún hoy en día tiene dificultades para ser aceptada en el mundo científico. Uno de los argumentos en contra lo expone el filósofo Michael Ruse en su libro *Mystery of Mysteris. Is Evolution a Social Construction?*:

A group perspective has internal problems: although in the long run everyone might benefit, it is difficult to see why (other than in exceptional circumstances) a short-term devotion to self should not be preferred by selection. Adaptations directed toward immediate personal benefit would seem to be fitter than adaptations benefiting others, even if down the road all would gain with, and only with, the latter. Unfortunately, selection is necessarily a short-term process, without forethought for the future (129).

En este sentido, ¿cómo puede la selección natural emplear mecanismos que “piensen” o “consideren” ventajas y resultados adaptativos en etapas poste-

riores, si la materia prima de la que se vale no es otra que las mutaciones aleatorias “de corto plazo” que se presentan en todo organismo? El rechazo a la selección de grupos está encabezado por científicos tan prominentes como George Williams (1966, 1992), Lack (1966), Ghiselin (1974), Maynard Smith (1976) o Dawkins (1982). Según Dawkins (1980), por ejemplo “un error común consiste en creer que la cooperación dentro de un grupo en un nivel dado de organización, debe aparecer mediante la Selección de grupos [...] La teoría de la Estrategia Evolutivamente Estable ofrece una alternativa más austera” (360).

Además, muchos autores están de acuerdo en que si esta teoría llegase a existir, necesitaría de grupos pequeños que permanentemente se recombinen para permitir de esta forma que la fuerza de la selección de grupos pueda realizar una contraposición efectiva a la fuerza de la selección individual (Sober & Wilson 1998; Fehr & Fischbacher 2005; Gintis *et al.* 2005). Esta condición no ha sido lo suficientemente determinada en el marco de las relaciones sociales de los seres humanos, además de ser limitada frente al extendido comportamiento cooperativo en grupos humanos de gran tamaño. De otra parte, como lo sostiene Nesse (1994: 633), esta teoría aún requiere claridad conceptual y datos empíricos para su comprobación.

En cuanto a la claridad conceptual, filósofos y biólogos (Hull 1980; Lloyd 2001; Andrade & Fajardo 2008) han intentado desarrollar aparatajes conceptuales bien definidos –y por qué no, posturas ontológicas– que identifiquen con claridad elementos tan importantes para la selección como son la unidad de selección que difiere del nivel mismo de selección. Esto en palabras de Dawkins (1976) es la identificación clara del agente replicador y de la entidad que interactúa con el medio ambiente y a la que se le puede adscribir el cambio adaptativo. A pesar de todos estos esfuerzos, la cuestión sigue en debate.

8. CASTIGO NO ALTRUISTA. LA ESTABILIDAD DE LA COOPERACIÓN HUMANA SIN ACUDIR A LA SELECCIÓN DE GRUPOS

La *reciprocidad fuerte* se entiende como el comportamiento relacionado con dos aspectos: i) la cooperación condicional y ii) el castigo altruista, tanto para los no cooperadores como para los no castigadores (Fehr, Fischbacher & Gächter 2002; Bowles & Gintis 2003; Fehr & Henrich 2003; Fehr & Fischbacher 2003, 2005; Gintis *et al.* 2005; Sethi & Somanathan 2005). Como vimos en las secciones 5 y 6, la solución comúnmente aceptada en la actualidad como respuesta efectiva a los dilemas sociales de primer y segundo orden de la *reciprocidad fuerte* se enfrenta a la dificultad de suponer, en el segundo aspecto, la

necesidad de la selección de grupos (aún hoy en día altamente controvertida dentro de la ciencia biológica) (Boyd & Richerson 1992; Sethi & Somanathan 1996; Fehr & Gächter 2002; Gintis *et al.* 2005; Fehr & Fischbacher 2005; Gintis 2008).

El castigo altruista involucrado en la *reciprocidad fuerte* se ve explotado con distinta intensidad, según el tamaño del grupo y el nivel del castigo (castigo a no cooperadores o a no castigadores) por parte de los individuos que prefieren no cooperar o castigar en este sentido, pero que de cualquier forma participan de los beneficios del sistema de castigos mantenido por otros (Boyd & Richerson 1992). En estas condiciones, la forma de explicar la estabilidad de la *reciprocidad fuerte* y, por tanto, de solucionar los dilemas a los que se enfrenta, es sostener que la selección de grupos es la razón por la cual un comportamiento no adaptativo a nivel individual, como es el castigo altruista, es viable en la naturaleza, en especial en el comportamiento humano (Boyd *et al.* 2005). Según esta interpretación, como ya lo vimos, la frecuencia del castigo y el comportamiento cooperativo es positivamente correlacionado a través de los grupos. De esta manera los grupos con una mayor cantidad de castigadores presentan ventajas competitivas frente a aquellos que tienen una tasa menor o que no tienen castigadores (Boyd *et al.* 2005 217).

Mi intención principal en este artículo es presentar una solución alternativa a la estabilidad de la cooperación humana mediante los castigos y las recompensas, pero bajo las consideraciones básicas de la selección individual y, para el caso de los seres humanos, la racionalidad acotada. Para esto, quiero abordar dicha solución alternativa a través de una interpretación, desde la teoría de juegos, de la propuesta filosófica de Thomas Hobbes acerca del origen del Estado civil. Para ello, mencionaré en esta última parte los presupuestos teóricos de Hobbes relacionados con el papel del soberano en la constitución de un Estado. Gracias a este acercamiento, mostraré el valor filosófico del castigo y las recompensas como estabilizadores no altruistas de la cooperación humana.

8.1. Sanciones e incentivos. El papel del soberano hobbesiano

El acuerdo entre hombres libres, en medio de un Estado de naturaleza peligrosa para la seguridad individual, puede lograr que un soberano se instaure como máxima autoridad llamada a pacificar las relaciones colectivas (XVII. 137/128). Esto genera que los individuos salgan de un DP, o más bien que se percaten de que no están en él, y puedan, o coordinar sus preferencias con mayor facilidad, o incursionar en un juego de DPI. Este papel es desempeñado por el gobernante, dado su máximo poder (entregado por acuerdo), a través de distintos mecanismos, tales como proferir leyes civiles, amenazar

a los potenciales infractores y no cooperadores, castigar la deserción frente al cumplimiento de las leyes o incentivar su seguimiento con la finalidad de hacer racional la cooperación. Así, el soberano, después de ser instaurado, tiene el poder de transformar lo que en apariencia se veía en el Estado de naturaleza como un DP en una dinámica colaborativa.

Junto con el poder de transformar estructuralmente los DP, el soberano puede hacer que se dé efectivamente la mutua realización de promesas, de tal manera que con ellas se expresen las preferencias de las partes participantes en un acuerdo y, así, se afirme la adquisición de un compromiso. Para Hampton (1998 209), la promesa es fundamental para transformar la estructura de un DP, pues desde una perspectiva humeana el temor al castigo nos obliga a señalar, mediante una promesa, que preferimos determinado curso de acción que nos comprometemos con otra persona.

Hobbes sostiene que los derechos y facultades que constituyen la soberanía en ejecución son doce (XVIII. 142/133); entre los cuales los más relevantes para nuestra actual finalidad son proferir leyes civiles que regulen la vida en comunidad, junto con el diseño de incentivos y sanciones que le permitan al soberano obligar el seguimiento de los preceptos civiles establecidos por el gobierno. En este punto es importante señalar que el seguimiento de las leyes naturales, según Hobbes, es connatural a la racionalidad de cada uno. Todo hombre, por un lado, debe buscar a toda costa su supervivencia, incluso si para esto se tienen que enfrentar al mismo soberano (XXI. 177/167) y, por otro, debe tratar a los demás como él quisiera ser tratado (XIV. 107/100). Claro, en este último caso, prima la ley de autopreservación. De otra parte, también es facultad fundamental del soberano impartir castigos y recompensas a quienes violen las leyes por él instauradas. Los castigos y recompensas como dispositivos “sociales” son para Hobbes, resalta Hardin (1990), de suma importancia para mantener cohesionado al Estado, ese hombre artificial (Leviatán) con el poder suficiente para hacer que los individuos se mantengan en paz.

[...] se asigna al soberano el poder de recompensar con riquezas u honores, y de castigar con penas corporales o pecuniarias, o con la ignominia, a cualquier súbdito, de acuerdo con la ley que él previamente estableció; o si no existe ley, de acuerdo con lo que el soberano considera más conducente para estimular los hombres a que sirvan al Estado, o para apartarlos de cualquier acto contrario al mismo (XVIII. 148/138).

Hardin (1990 79) afirma que desde el punto de vista de Hobbes, representante del contractualismo de los siglos XVII y XVIII, y teniendo en cuenta las propuestas de Locke, fundador de las teorías del derecho natural, los gobiernos fracasan si en ellos no se ejerce coerción en los comportamientos

individuales. En este sentido, el proyecto político de constituir un Estado y mantenerlo en el tiempo para consolidar la situación civil, alejada del desorden y el conflicto del hipotético Estado de naturaleza, no es viable si en él (en el Estado político o civil) no se presuponen elementos coercitivos que, ejecutados por el soberano, mantengan a los individuos temerosos de incumplir con el pacto de vivir en sociedad. Yo añadiría, como creo que lo asume Hobbes también, que además de los elementos coercitivos son necesarios elementos o dispositivos de recompensa o incentivos. Hobbes sostiene que el castigo y la recompensa “vienen a ser los nervios y tendones que mueven los miembros y articulaciones de un Estado” (XXVIII. 262/246).

Ahora bien, ¿cuál es la justificación de la coerción y la recompensa en la teoría hobbesiana del Estado? En principio, como el mismo Hardin (1990) lo resalta, es importante señalar que la instauración de un soberano mediante el contrato social suscrito por los individuos en situación de guerra (al estilo hobbesiano) se desarrolla de manera racional mediante un pleno consentimiento de estos. Así que este consentimiento racional debe extenderse también a las facultades y dispositivos de gobierno necesarios que tal soberano deberá ejecutar para cumplir con el objetivo para el cual ha sido elegido.

Así, en general, la mayor parte de nuestras interacciones se llevan a cabo frente a otra persona (relación diádica) y, por lo común, estas relaciones están atravesadas por intercambio de promesas, explícitas o no, que mantienen la colaboración en un punto estable. Teniendo en cuenta esto, se presentan situaciones difíciles en las que consentimos que nuestros gobernantes asuman sistemas de coerción para ejecutar castigos, de tal forma que podamos disfrutar de los beneficios de los contratos y las promesas.

Sabemos, y lo consentimos racionalmente, que los beneficios de la cooperación instaurada mediante contratos y promesas nos ofrecen mayores réditos que el no cumplir. Por esta razón aceptamos la posibilidad y la amenaza de castigo por el incumplimiento, si una de las dos partes decide desertar. Consentimos esta posibilidad a pesar de que el castigado puede ser uno mismo. Somos conscientes de los beneficios del sistema, incluso si yo resulto castigado por su ejecución.

Piénsese, por ejemplo, en situaciones comerciales en donde las tarifas y los precios han sido definidos de antemano por un ente regulador (gobierno o soberano). Para mí, como cliente está bien que las tarifas de algunos servicios no se definan al gusto del prestador (podría abusar de su monopolio). De la misma forma, para el prestador no es bueno competir con tarifas que otros, con mayores capacidades y recursos, pueden ofrecer, explotando su sistema de tarifas con precios más bajos. Ahora, ambas partes sabemos que si uno de

los dos no respeta estas tarifas preestablecidas, caerá sobre aquel una sanción estatal, así que será mejor que coopere y respete dichas tarifas.

Como lo afirma Hardin (1990 80), no podemos concebir a las sociedades modernas sin respaldar sus relaciones contractuales con sistemas de coerción. Hobbes diría que sin la espada, los pactos solo son meras palabras (XXVIII). En esto radica el consentimiento racional. Al asumir y aceptar los sistemas coercitivos, comprendiendo la superioridad de los beneficios que dichos sistemas pueden asegurar en las relaciones contractuales, actuamos de manera racional y buscamos nuestros propios intereses.

De la misma forma, en interacciones multiagente como la provisión de bienes públicos, la influencia de los sistemas de coerción es fundamental para poder disfrutar de los beneficios que ofrece la cooperación. Casos como los sistemas de control de tráfico en donde aparecen vinculados comparendos y sanciones para infractores son consentidos racionalmente, pues el simple hecho de que todos contemplan la posibilidad del castigo en estas situaciones nos permite beneficiarnos de la buena conducción de la mayoría de las personas que están al frente del volante. Aunque siempre conducimos con precaución, ya que anticipamos alguna posible violación de la norma, de alguna manera sentimos cierta confianza en que, en general, ese hecho constituye una probabilidad muy baja. De lo contrario, nunca conduciríamos en las calles. Si consideramos que las posibilidades de la infracción son altas, no pondríamos en riesgo nuestra vida, ni las de nuestros acompañantes.

Todos estos elementos se siguen enmarcando en la teoría hobbesiana del auto-interés como fundamento psicológico del comportamiento humano, además del principio de autopreservación que gobierna a los hombres. De acuerdo con esto, y como acertadamente lo sostiene Hardin (1990 80), el interés racional es suficiente para explicar los sistemas de coerción.

Desde el punto de vista de Hobbes, los súbditos admiten la coerción de un soberano porque mediante ella se protegen de las pasiones de otros y, por qué no decirlo, de ellos mismos, al extraerlos, con esto, de una dinámica del DP. Vistas así las cosas, podría decirse que buena parte del bienestar del Estado y de las personas que en él interactúan se debe al gobierno coercitivo que lo dirige. Recuérdese que para Hobbes la vida sin gobierno es severa e improductiva, incluso peligrosa para la seguridad personal. De este modo, la justificación de la coerción en el gobierno hobbesiano se encuentra en el hecho de que mediante los castigos (de la misma forma que con los incentivos) la estructura de un juego de DP se transforma fácilmente en una de coordinación de preferencias o también en un de DPI.

Así, si dos o más personas se encuentran aportando para un bien público, pongamos por caso impuestos para la construcción de caminos, sin un poder

coercitivo las preferencias de los individuos serán las de SIEMPRE D, siempre desertar, dado que encuentran más rentable explotar la cooperación de los otros o, en su defecto, prefieren no cooperar para que otros no los exploten. Esta interacción tiene la típica matriz de preferencias de un DP, o de un juego de bienes públicos. Sin embargo, si existe un poder coercitivo que amenace con infringir castigos a los individuos no cooperadores y que les haga ver que explotar la cooperación puede llegar a ser menos ventajosa de lo que creían, entonces sus preferencias se transforman en actitudes cooperativas que, por supuesto, son más rentables que la desertión asumiendo el castigo. De esta manera, un sistema de castigos transforma la estructura de un juego de DP o bienes públicos en uno de coordinación de preferencias o DPI. Hobbes sostiene algo muy cercano a esto, pues al menos en una ocasión del *Leviatán* (XXVII. 241/226) afirma que los castigos que no son suficientemente grandes para disuadir al infractor lo inducirán a violar la norma porque este solo calculará que el beneficio de desertar es mayor que el propio castigo:

En efecto, si de antemano se conoce el castigo, cuando éste no es bastante grande para disuadir de la acción, constituye un estímulo para ella, porque cuando los hombres comparan el benéfico de la injusticia por ellos cometida con el daño que representa su castigo, por razón de naturaleza eligen lo que resulta preferible para ellos (XXVII. 241/226).

Con todo, creo que Hobbes suscribiría un punto de vista como el sostenido anteriormente y que se refiere al poder que tienen los castigos, y de la misma forma los incentivos, en transformar un DP en un problema de coordinación de preferencias. Esto es fundamental puesto que así se ratifica filosóficamente el importante papel del castigo y los incentivos en la estabilización de la cooperación.

8.2. Castigos y recompensas institucionales no altruistas

Aquí me quiero detener un momento en la obvia solución que extrañamente pocos han considerado. Por lo menos, en el caso de la especie humana es evidente que un sistema de castigos pudo influir de modo sustantivo en la evolución de la cooperación⁷. Hemos encontrado un seguro que protege los sistemas de castigo

⁷ Por sugerencia del evaluador anónimo de este artículo quiero resaltar que existen otros mecanismos distintos al castigo que también influyen de alguna manera en la estabilización de la cooperación. Uno de estos mecanismos puede ser la reputación o búsqueda de buena imagen ante los demás. Este mecanismo también provoca que el individuo que lo ejecuta logre beneficios al ser preferido para interacciones cooperativas futuras, por ejemplo. Sin embargo, creo que el castigo es una fuerza más poderosa y ampliamente extendida para lograr la estabilidad de la cooperación; nadie quiere enfrentarse a un castigo, mientras que la reputación puede no ser prioridad en una estrategia para la interacción.

de la explotación de los desertores, evitando que el sistema se disuelva en una reducción ad infinitum de castigos de primer, segundo, tercer orden, etc., para luchar en contra de los *free riders* que aparecen en todos los niveles.

Así, de acuerdo con Hobbes, el gobernante goza de un gran honor y dignidad por desempeñar sus funciones en el mantenimiento del orden social (XVIII. 149/140). El poder que se le entrega junto con la dignidad y el honor es un mayúsculo pago por sus servicios como dirigente del Estado y, por supuesto, como supremo ejecutor del sistema de coerción. De la misma forma, menciona Hobbes, que a aquellos quienes por su labor entregan servicios al Estado, entre los que podemos contar a los cuadros de control o cuadros de coerción institucional, se les deben entregar recompensas (salarios o sueldos) en reconocimiento a las labores realizadas por contrato (XXVIII. 260/245). Es más, Hobbes sostiene que a algunos individuos se les entregarán incentivos en liberalidad (sin contrato de por medio) por el mero hecho de que con ello se incentiva el mejor servicio al Estado (XXVIII. 260/245), pongamos por caso, personas no vinculadas a los servicios del Estado que denuncian las deserciones.

En el primer caso el individuo queda obligado por contrato a desempeñar las funciones del cargo, sean las de legislador, juez o integrante de un cuerpo de coerción. Este es el punto que quiero resaltar. Es evidente que, manteniendo los presupuestos de la psicología egoísta de Hobbes y aquellos enmarcados en la teoría de la elección racional, los cuadros de control o coercitivos deben ser recompensados con incentivos tipo salario y demás, lo suficientemente bien como para que los costos de desempeñar funciones de castigadores sean menores frente al beneficio obtenido, que en muchos casos también incluye el honor y el poder. De esta manera, son los incentivos los que transforman la estructura de pagos de los DP de segundo orden, para generar juegos de coordinación, juegos en donde los intereses de los participantes se coordinan para cooperar entre sí. Tanto el soberano hobbesiano, como los individuos que llevan a cabo el “trabajo sucio” de castigar, son recompensados por el Estado, superando los costos de sus funciones, con recursos generados por la comunidad para tal fin. Aparece aquí el racional consentimiento de las personas para aceptar y aportar a un sistema de coerción que se autorrefuerza mediante los aportes que el sistema mismo vigila para que se entreguen.

En otras palabras, los *castigos institucionales* derivados de sistemas coercitivos estatales rompen con los DP de órdenes superiores, puesto que mediante las recompensas instituidas como sueldos, más aquellas vinculadas a su misión como el honor y el poder, se superan los costos que deben asumir los individuos castigadores institucionales. De tal suerte que el castigo ya no es un fenómeno extraño de tipo altruista, ni biológico ni psicológico, que aporta a

la estabilidad de la cooperación, y no puede explicarse desde un punto de vista racional. Más específicamente, dado que los castigos institucionales reportan más beneficios que costos a quienes los ejecutan, dicho comportamiento no puede enmarcarse, desde un punto de vista evolutivo, bajo el altruismo biológico de ceder beneficios a los demás en contra de su propio éxito adaptativo (invertir en el fortalecimiento de un sistema de castigos para que todos se beneficien). Además, este tipo de castigos, así considerados, tampoco pueden enmarcarse bajo la denominación de altruismo psicológico, pues el ejecutor del castigo institucional está motivado más por su propio interés de lograr recompensas y reconocimientos que por el mero hecho de castigar a alguien por aversión al delito o por empatía y solidaridad hacia los explotados. Por tanto, bajo las dinámicas de la elección racional, el castigo institucional en la especie humana es la solución más consistente a la estabilidad de la cooperación. Como lo hemos visto, ya Hobbes lo había contemplado.

Los castigos institucionales, enmarcados en un sistema de control y coerción que se autorrefuerza con la vigilancia del aporte a este sistema y a las demás empresas comunales, transforman efectivamente la estructura de pagos de un DP, o de bienes públicos en interacciones multiagente, para convertirlos en problemas de coordinación que tienen solución en puntos de equilibrio colaborativos. Gracias a que los juegos de coordinación contemplan la cooperación como una solución posible, la estrategia que desemboca en este tipo de equilibrio es autorreforzada, en tanto que los individuos no presentan preferencias distintas a seguir cooperando.

Puede interpretarse que Hobbes, en algún sentido, consideraba los castigos y, por supuesto, los incentivos como dispositivos que coordinan preferencias. Para él, los castigos establecen ejemplos para que los posibles infractores prefieran cooperar, “[...] la finalidad del castigo no es la venganza y la descarga de la ira, sino el propósito de corregir tanto al ofensor *como a los demás, estableciendo un ejemplo*” (XXX. 286/269; cursiva fuera de texto).

En otras palabras, los castigos institucionales pueden hacer que las personas sean más conscientes de sus preferencias y las coordinen con otros individuos. Las sanciones permiten identificar una *prominencia* o *“saliencia”* de un equilibrio colaborativo, por el que en principio la mayoría optará. Piénsese en el caso del control del tráfico: si comprendo que las multas o castigos por conducir bajo los efectos del alcohol son sustancialmente altas, no tendremos muchos problemas en tomar un taxi y realizar una transacción con el conductor. Confiamos (esta es la coordinación de preferencias) en que al conductor no le interesa estar ebrio, tampoco violar peligrosamente las normas de tránsito, de tal forma que nos llevará a salvo a nuestro destino.

En otras palabras, es necesario tener en cuenta que en muchos casos, como lo sostiene Gilbert (2001), las sanciones no solamente ajustan el comportamiento de las personas al del grupo, por el simple hecho de que los individuos anticipan o prevén un posible castigo a la desertión; también, y tal vez más importante que lo anterior, los castigos hacen que los individuos recuerden que hacen parte de una comunidad con ciertas reglas de juego y que están comprometidos con su funcionamiento. La identidad de grupo es una de las soluciones a los dilemas sociales, una de las más efectivas maneras de fomentar la cooperación.

De otra parte, es bueno señalar que no es necesario explicar el caso de los *castigos institucionales*, que ya existían en las sociedades humanas más tempranas (Lyons 2003), como castigos biológicamente altruistas que dependen de un supuesto y controvertido proceso de selección de grupos. Más bien este tipo de mecanismos se explican mejor bajo la clásica selección individual, al encontrar que el beneficio por castigar es mucho mayor que el costo. Además, ese beneficio es asegurado por la misma labor que realizan los agentes de control frente a los aportadores al bien público.

Quiero, por último, referirme a la importancia de los incentivos como dispositivos estabilizadores de la cooperación. De la misma forma que los castigos, si los incentivos son dispensados por un ente de control, como lo menciona Hobbes, los dilemas de órdenes superiores se desvirtúan de inmediato, puesto que también ellos se alimentan de aportes controlados y vigilados por los entes estatales a los que se les paga por dicha tarea. En principio, vimos como los incentivos o recompensas son fundamentales para romper con el círculo vicioso de los sistemas de control que explotaban los *free riders*. Gracias a las recompensas adecuadas, los participantes de cuadros de castigo tienen el incentivo suficiente para realizar su labor sin hacer de sus acciones estrategias irracionales de altruismo. Se comprende entonces cómo los incentivos también transforman la estructura de DP en coordinaciones de intereses que fomentan la cooperación.

Quizás el vínculo entre sanciones e incentivos es más fuerte de lo que normalmente percibimos. No solamente ambos tienen un papel similar en la estabilización de la cooperación, sino que también pueden explicarse mediante los presupuestos de la teoría de la elección racional. Además, a nivel neuronal, castigar produce en los ejecutores de la sanción efectos similares a cuando ellos mismos reciben recompensas (Nowak 286). Según Nowak, los castigadores experimentan “agradables oleadas” de irrigación sanguínea en los centros cerebrales de la recompensa. Esto podría reflejar el correlato cerebral de experiencias relacionadas con el placer a castigar que casi todos sentimos

cuando le “damos su merecido” a alguien, un tipo básico de *Shadenfreude*. Así pues, la ciencia neurobiológica puede ofrecernos pruebas de la evolución del castigo y los incentivos, como caminando hombro con hombro en medio de nuestra filogenia.

En cualquier caso, la recompensa, si la comparamos con el castigo, es un dispositivo que lleva a la cooperación de manera más creativa. Esto es así porque, en el caso de la recompensa, las personas que transforman sus DP gracias a ella, en general no solo cumplen con los acuerdos o normas, en muchas ocasiones superan estas expectativas y generan nuevos vínculos de cooperación. La recompensa incentiva, no solamente presiona, la cooperación. En otro sentido, los castigos estabilizan la cooperación y generan una coordinación de preferencias en puntos de equilibrio determinados, pero, también habitualmente los individuos en estos casos solo se limitan a cumplir con los tratos y tareas, sin ir más allá de lo prescrito o prometido (Nowak 301). Por esto que Nowak sostiene que “Reward, not necessity, is the true mother of invention” (301).

9. CONCLUSIONES

Como lo mencioné en la sección 7, aunque hay muchos científicos que proponen una salida más inmediata a la paradoja de la evolución de la cooperación postulando un mecanismo de selección a nivel de grupos, considero que no es necesario adoptar tal solución por dos razones principalmente. Si bien la solución a partir de la selección de grupos explica el fenómeno, el mecanismo en sí tiene muchos “cabos sueltos”. Algunos de estos son por ejemplo: *i*) que los científicos aún no se ponen de acuerdo ni en el estatus epistemológico, ni en el estatus ontológico de este tipo de selección; *ii*) que a pesar de algunos desarrollos teóricos para demostrar la selección de grupos, los ejemplos de dichos procesos de selección son altamente controvertidos y, en su mayoría, se pueden explicar sin acudir a tal propuesta; *iii*) que la selección de grupos aún hoy carece de pruebas empíricas que la corroboren. La otra razón por la que no es necesario acudir a la selección de grupos como única vía de solución de la cooperación tiene que ver con que, como lo he intentado mostrar, la cooperación en seres humanos puede explicarse bajo los presupuestos concomitantes de la selección individual y la teoría de la elección racional.

Así, como se ha mostrado, juegos evolutivos como el dilema del prisionero iterado (DPI) presentan buenas razones para aceptar que la selección individual puede soportar la cooperación, aunque de manera inestable. Además, si se incluye la evolución de las sanciones y las recompensas, un sistema de coope-

ración puede llegar a alcanzar una cierta estabilidad (estabilidad polimórfica en la mayoría de los casos). Como lo expliqué (sección 5), esta estabilidad polimórfica, lograda de mejor manera por la inclusión de un sistema de vigilancia, sigue siendo insuficiente para entender la cooperación extendida en seres humanos, puesto que, teniendo en cuenta las condiciones impuestas por la racionalidad individual de los agentes, también el sistema de vigilancia puede llegar a ser explotado, convirtiéndolo en un dilema de segundo orden. Sin embargo, valiéndonos de una racionalidad acotada (racionalidad imperfecta, que emplea información incompleta y que, además, está influida por las emociones, especialmente aquellas que se derivan de los castigos y las recompensas), podemos solventar los problemas de explotación de órdenes superiores del sistema de cooperación y sus sistemas de vigilancia, lo cual genera mejores posibilidades para explicar la cooperación desde dos puntos de vista, uno evolutivo y otro de racionalidad individual.

Creo que solucionar el problema paradójico de la evolución de la cooperación humana es solucionar, en últimas, el problema de la explotación de los sistemas de vigilancia que estabilizan la cooperación. En este caso mi propuesta es simple y directa. Postulo que si encontramos un tipo de sanciones y recompensas que entreguen a sus ejecutores beneficios que superen la inversión o el costo de ejecutar dichas sanciones o recompensas, podemos alcanzar un sistema de cooperación autovigilante y estable. ¿Cómo se logra esto? Concentrándonos en los castigos institucionales que ya Hobbes había propuesto como solución al dilema del contrato social y el seguimiento de leyes en un Estado de civilidad. Afirmo que los incentivos institucionales no deben ser considerados mecanismos altruistas, ni a nivel biológico, ni a nivel psicológico. En el primer caso, porque los ejecutores de incentivos institucionales (negativos o positivos) son recompensados por el statu quo para que cumplan su función a cabalidad. Así, en lugar de ser un esfuerzo costoso, castigar o recompensar es más bien “al final del día” una labor llena de beneficios, recompensas u honores. En el segundo caso, puede pensarse de una manera más coherente con la racionalidad individual que los castigos y recompensas institucionales no son ejecutados por motivaciones relacionadas con la empatía o la solidaridad, pues dado que el agente ejecutor recibe beneficios claros y superiores por realizar su función, este castiga o recompensa siguiendo la lógica que le señala su propio interés.

Con un sistema de vigilancia estable, lejos del altruismo y, por tanto, de los oportunistas o *free riders*, la cooperación entre humanos a escalas mayores se extiende a través de múltiples factores que, en últimas, lo único que provocan es que los individuos puedan coordinar eficientemente sus diversas preferencias para sacar el mejor provecho de la cooperación mutua. Las recompensas y

los castigos generan, presionan o acondicionan el ambiente para que factores de solución de dilemas sociales aparezcan. Así uno de los más importantes factores de solución de dilemas sociales que es generado por los incentivos, tanto negativos como positivos, tiene que ver con la modificación de la estructura del juego de DP en juegos del tipo DPI o juegos de confianza (AG). Eso abre la posibilidad de varios equilibrios de Nash que pueden coincidir en la cooperación mutua.

Otro factor importante que es provocado por los incentivos institucionales es la denominada identidad de grupos porque compartir un conjunto definido de reglas de juego puede hacer que los individuos se sientan parte de un grupo al que hay que aportar y hacer respetar. Otro factor estudiado y que es provocado por un sistema de vigilancia se relaciona con las denominadas *prominencias* o "*salinecias*" que la posibilidad de castigos o recompensas provoca en las preferencias de los jugadores autointeresados. Puede suceder que el simple hecho de saber que la mayoría de las personas conocen la posibilidad de un castigo por realizar una acción explotadora nos hace evidente que la mejor opción para actuar es cooperar, puesto que la mayoría de la gente no se arriesgaría a la sanción, como nosotros mismos no lo haríamos. Por último, es importante señalar que las emociones generadas por las sanciones y recompensas refuerzan cada uno de los factores antes mencionados y provocan que las motivaciones se conjuguen con las decisiones racionales para la cooperación. Esta es una de las principales características de la racionalidad acotada pues su desarrollo está íntimamente influido por factores de tipo emocional.

Con todo, quiero aclarar que mi explicación se aparta de la consideración clásica de la economía que presupone a un hombre eminentemente calculador, con racionalidad perfecta e información completa. Mi propuesta asume la naturaleza del hombre racional, pero imperfecto que actúa influido por sus emociones y, claro está, por su propio interés que acaso no es más que emoción racionalizada, pero al fin de al cabo razón en ejecución.

Por lo demás, es importante resaltar que en este trabajo quiero distinguir entre las explicaciones últimas y próximas de la cooperación. Así, sostengo que no es necesario apelar a la selección de grupos como explicación última de la cooperación, puesto que podemos seguir manteniendo una explicación al nivel orgánico o individual. Y esto es así porque podemos concebir sin problemas los castigos institucionales como una explicación próxima de la cooperación, ya que ellos la estabilizan y además se encuentran desprovistos de todo carácter altruista, tanto desde el punto de vista biológico como desde el psicológico.

Por último, mis conclusiones me llevan a pensar que, a diferencia de algunos autores como Nowak (2011) (quien sostiene que gracias a la cooperación

hemos llegado a generar pensamiento racional, complejo y discursivo), la evolución de la cooperación, por lo menos el tipo de cooperación que empleamos los seres humanos, debe presuponer una coevolución de otros mecanismos necesarios para su estabilización. Mecanismos como la propia razón y las emociones debieron evolucionar juntos, en paralelo, con la cooperación, para poder generar la posibilidad de una cooperación expandida en grandes grupos humanos.

TRABAJOS CITADOS

- Alexander, R. D. *The Biology of Moral Systems*. Transaction Publishers, 1987.
- Andrade, L. E. & D. Fajardo. “Niveles de selección y la cuestión ontológica acerca de las jerarquías biológicas”. *El hombre y la máquina* 30 (2008): 86-99.
- Axelrod, A. *The Complexity of Cooperation: Agent-Based Models of Competition and Collaboration*. Princeton, NJ: Princeton University Press, 1997.
- _____. “An Evolutionary Approach to Norms”. *The American Political Science Review* 80.4 (dic. 1986): 1095-111.
- _____. *The Evolution of Cooperation*. New York: Basic Books, 1984.
- Axelrod, R. & W. D. Hamilton. “The Evolution of Cooperation”. *Science* 211.4489 (1981): 1390-396.
- Bendor, J. & D. Mookherjee. “Institutional Structure and the Logic of Ongoing Collective Action”. *American Political Science Review* 81 (1987): 129-54.
- Boehm, C. “Egalitarian Behavior and Reverse Dominance Hierarchy”. *Current Anthropology* 34.3 (1993): 227-54.
- Bowles, S. & H. Gintis. “The Origins of Human Cooperation”. *The Genetic and Cultural Origins of Cooperation*. Ed. P. Hammerstein. Cambridge: MIT Press, 2003.
- Boyd, R. “The Evolution of Reciprocity in Sizable Groups”. *Journal Theory of Biology* 132 (1988): 337-56.
- Boyd, R., H. Gintis, S. Bowles & P. J. Richerson. “The Evolution of Altruistic Punishment”. Ed. H. Gintis. *Moral Sentiments and Material Interests: The Foundations of Cooperation in Economic Life*. MIT Press, 2005.
- Boyd, R. & P. J. Richerson. “Punishment Allows the Evolution of Cooperation (or anything else) in Sizable Groups”. *Ethology and Sociobiology* 13 (1992): 171-95.

- _____. "The Evolution of Reciprocity in Sizable Groups". *Journal Theory of Biology* 132 (1989): 337-56.
- Cárdenas, J. C. *Dilemas de lo colectivo: instituciones, pobreza y cooperación en el manejo local de los recursos de uso común*. Bogotá: Universidad de Los Andes, 2009.
- Darwin, Ch. *The Descent of Man* [1871]. New York: Penguin Books, 2004.
- _____. *El origen del hombre* [1871]. Bogotá: Panamericana, 1994.
- _____. *The Origin of Species by Means of Natural Selection: or, the Preservation of Favored Races in the Struggle for Life* [1859]. Ed. J. W. Burrow. London: Penguin Books, 1968.
- Dawkins, R. *The Extended Phenotype: The Long Reach of the Gene*. New York: Oxford University Press, 1989.
- _____. *The Selfish Gene*. New York: Oxford University Press, 1976.
- Elster, J. *Juicios salomónicos: las limitaciones de la racionalidad como principio de decisión*. Barcelona: Gedisa, 1991a.
- _____. *Tuercas y tornillos: una introducción a los conceptos básicos de las ciencias sociales*. Barcelona: Gedisa, 1991b.
- Falk, A. & U. Fischbacher. "A Theory of Reciprocity". *Games and Economic Behavior* 54(2) (2006): 293-315.
- Fehr, E. & U. Fischbacher. "The Economics of Strong Reciprocity". *Moral Sentiments and Material Interests. The Foundations for Cooperation in Economic Life* (2005): 151-93.
- _____. "Third Pary Punishment and Social Norms". *Evolution y Human Behavior* 25 (2004): 63-87.
- _____. "The Nature of Human Altruism". *Nature* 425 (2003): 785-91.
- Fehr, E., U. Fischbacher y S. Gächter. "Altruistic Punishment in Humans". *Nature* 415 (2002): 137-40.
- Fehr, E. & J. Henrich. "Is Strong Reciprocity a Maladaptation? On the Evolutionary Foundations of Human Altruism". *Genetic and Cultural Evolution of Cooperation*. Ed. P. Hammerstein. Cambridge: MIT Press, 2003.
- Frank, R. H. *Passions within Reason: The Strategic Role of the Emotions*. WW Norton y Co., 1988.
- Ghiselin, M. T. *The Economy of Nature and the Evolution of Sex*. Berkley: University of California Press, 1974.

- Gilbert, M. "Collective Preferences, Obligations, and Rational Choice". *Economics and Philosophy* 17 (2001): 109-19.
- Gintis, H. (ed). *Moral Sentiments and Material Interests: The Foundations of Cooperation in Economic Life*. Cambridge: MIT Press, 2005.
- _____. "Strong Reciprocity and Human Sociality". *Journal of Theoretical Biology* 206.2 (2000): 169-79.
- Gould, S. J. *The Structure of Evolutionary Theory*. Cambridge: The Bepnap Press, 2002.
- Hamilton, W. D. "The Genetical Evolution of Social Behavior. I". *Journal of Theoretical Biology* 7.1 (1964): 1-16.
- _____. "The Evolution of Altruistic Behavior". *American Naturalis* (1963): 354-56.
- Hampton, J. *Hobbes and the Social Contract Tradition*. Cambridge University Press, 1986.
- Hardin, R. "Racionallity Justifying Political Coercion". *Journal of Philosophical Research* 15 (1990): 79-91.
- Hauert, C. *et al.* "Via Freedom to Coercion: The Emergence of Costly Punishment". *Science*, 316.5833 (2007): 1905-907. doi: 10.1126/science.1141588.
- Henrich, J. & R. Boyd. "Why People Punish Defector". *Journal Theory of Biology* 208 (2001): 79-89.
- Hernich, J. *et al.* "Costly Punishment across Human Societies". *Science* 312 (2006): 1767-770.
- Hirshleifer, D. & E. Rasmusen. "Cooperation in a Repeated Prisoner's Dilemma with Ostracism". *Journal of Economic Behavior and Organization*, 12 (1989): 87-106.
- Hobbes, T. *Leviatán* [1651]. México: Fondo de Cultura Económica, 1980.
- _____. *Hobbes's Leviathan* [1651]. Oxford: Oxford at The Clarendon Press, 1909.
- Hull, D. "Individuality and Selection". *Annual Review of Ecology and Systematics* 11 (1980): 311-32.
- Joshi, N. "Evolution of Cooperation by Reciprocation within Structured Demes". *Journal of Genetics* 66 (1987): 69-84.
- Lack, D. *Population Studies of Birds*. Oxford: Oxford University Press, 1966.

- Lewontin, R. "The Units of Selection". *Annual Review of Ecology and Systematics* 1 (1970): 1-16.
- Lyons, L. *The History of Punishment*. London: Amber Books, 2003.
- Maynard Smith, J. "Group Selection". *Quarterly Review of Biology* 51.2 (1976): 277-83.
- Nesse, R. "Why is Group Selection such a Problem?" *Behavioral and Brain Sciences* 17.4 (1994): 633-34.
- Nowak, M. A. & R. Highfield. *SuperCooperators: Altruism, Evolution, and Why We Need each other to Succeed*. New York: Free Press, 2011.
- Oliver, P. "Rewards and Punishment as Selective Incentives for Collective Action: Theoretical Investigations". *American Journal of Sociology* 85 (1980): 1356-375.
- Ostrom, E. "A Behavioral Approach to the Rational Choice Theory of Collective Action". *American Political Science Review* 92.1 (1998): 1-22.
- Rosas, A. "The Return of Reciprocity: a Psychological Approach to the Evolution of Cooperation". *Biology y Amp; Philosophy* 23.4 (2008): 555-66.
- _____. "Las emociones morales como adaptaciones para la cooperación en dilemas sociales". *Ludus Vitalis* 15.28 (2007): 97-118.
- _____. "La moral y sus sombras: la racionalidad instrumental y la evolución de las normas de equidad". *Crítica: Revista Hispanoamericana de Filosofía* 37.110 (2005):79-104.
- Ruse, M. *Mystery of Mysteries: Is Evolution a Social Construction?* Cambridge, Mass: Harvard University Press, 1999.
- Sethi, R. & E. Somanathan. "The Evolution of the Social Norms in Common Property Resource Use". *The American Economics Review* 86.4 (1996): 766-88.
- Smith, A. *The Theory of Moral Sentiments [1759]*. New York: Penguin, 2010.
- _____. *An Inquiry into the Nature and Causes of the Wealth of Nations: a Selected Edition [1776]*. New York: Oxford University Press, 2008.
- _____. *Teoría de los sentimientos morales [1759]*. Madrid: Alianza Editores, 2004.
- _____. *Investigación sobre la naturaleza y causas de la riqueza de las naciones [1776]*. México: Fondo de Cultura Económica, 1958.
- Sober, E. & S. D Wilson. *Unto Others. The Evolution and Psychology of Unselfish Behavior*. Cambridge: Harvard University Press, 1998.

- Sripada, C. S. "Punishment and the Strategic Structure of Moral Systems". *Biology and Philosophy* 20 (2005): 767-89.
- Sterelny, K. *The Evolution of Agency and other Essays*. Cambridge: Cambridge University Press, 2001.
- _____. *Thought in a Hostile World*. New York: Blackwell, 2003.
- Sterelny, K. & Griffiths, P. E. *Sex and Death*. Chicago: The University of Chicago Press, 1999.
- Trivers, R. L. "The Evolution of Reciprocal Altruism". *Quarterly Review of Biology* (1971): 35-57.
- Williams, G. *Adaptation and Natural Selection*. New Jersey: Princeton University Press, 1966.
- _____. *Natural Selection: Domains, Levels, and Challenges*. Oxford: Oxford University Press, 1992.
- Wynne-Edwards, V. C. *Animal Dispersion in Relation to Social Behavior*. London: Oliver y Boyd, 1962.
- Yamagishi, T. "The Provision and Sanctioning System as a Public Good". *Journal of Personality and Social Psychology* 51 (1986): 110-16.
- Yamahashi, T. & N. Takajashi. "Evolution of Norms without Metanorms". *Social Dilemmas and Cooperation*. Ed. Schulz, U., W. Albers y U. Muller. Berlin: Springer-Verlag, 1994.

