# Stochastic Approximation, Momentum, and Nash Play

H. Berglann[¤] and S. D. Flåm[y]

April 3. 2002

**Abstract.** Main objects here are normal-form games, featuring uncertainty and noncooperative players who entertain local visions, form local approximations, and hesitate in making large, swift adjustments. For the purpose of reaching Nash equilibrium, or learning such play, we advocate and illustrate an algorithm that combines stochastic gradient projection with the heavyball method. What emerges is a coupled, constrained, second-order stochastic process. Some friction feeds into and stabilizes myopic approximations. Convergence to Nash play obtains under seemingly weak and natural conditions, an important one being that accumulated marginal payo¤s remains bounded above.

Key words: Noncooperative games, Nash equilibrium, stochastic programming and approximation, the heavy ball method.

## 1. Introduction

Game theory has become an enormously important ...eld of study [18], [22]. It now o¤ers unity or coherence to many lines of inquiry. To wit, various social sciences, while steadily growing more game-theoretic, increasingly reckon Nash equilibrium as a focal point and key concept. Additional bonus and impetus comes with making those sciences more experimental [16].

The Nash solution concept nicely formalizes stable interaction among noncooperative players. In doing so it tends, however, to make larger demands on the players' knowledge and rationality than can easily be justi...ed. Usually, to become a clever strategist one needs substantial learning or much experience. Therefore, Nash equilibrium begs some justi...cation in dynamic terms [15], [19], [21].

Four features of human behavior then seem important. First, individuals always try to improve own welfare (or payo¤). Second, they typically do not quite know all strategic possibilities, intentions or consequences. Third, they are likely to form local perspectives and approximations. Fourth, they hesitate in making quick and large adjustments.

1

To embody all these features in dynamics, and to accommodate exogenous uncertainty as well, we advocate and illustrate here use of stochastic approximation [6]. That vehicle, which leans heavily on di¤erential equations, subscribes to a tradition that goes back to classical mechanics - and to Newton's claim that the initial state of a mechanical system determines its development.

The purpose of this paper is to put that claim to use for the study of stochastic games. Alternatively, one may read this paper as dealing with stochastic programming, parallel computation, or global optimization. Subsequent arguments are organized around a noncooperative stage game repeated time and again. Technicalities and proofs are found in the references.

## 2.   The Game

There is a …xed, …nite set $I$ of players. Agent $i \; 2 \; I$ is constrained to choose his strategy $x_i$ from a nonempty compact convex subset $X_i$ of a Euclidean space. He always seeks to improve own expected payo¤ $\frac{1}{4}_i(x_i; x_{\text{-}i}) := E \; ¦_i(x_i; x_{\text{-}i}; !)$. Here $x_{\text{-}i} =: (x_j)_{j \neq i}$ denotes the part of the overall strategy pro…le $x = (x_i)$ that is controlled by $i^0$s rivals. The elementary event $!$ belongs to a complete probability space $(\text{--}; \frac{3}{4}; ^1)$; with respect to which one takes the mathematical expectation $E$: Each bivariate function $(x; !) \; \mathcal{V} \; ¦_i(x; !) \; 2 \; R$ is concave, di¤erentiable in $x_i$; and integrable in $!$:

Of prime interest are points $x \; 2 \; X := ¦_{i2I} X_i$ where each marginal payo¤ $m_i(x) := \frac{@}{@x_i} \frac{1}{4}_i(x)$ is normal to $X_i$ - or quite simply, nil. That is, letting $P_i$ denote the orthogonal projection onto $X_i$ we seek a …xed point $x = (x_i)$ of the system

$$x_i \; \tilde{A} \; P_i[x_i + sm_i(x)] \text{ for all } i \text{ and arbitrary } s > 0: \tag{1}$$

Any such …xed point is a Nash equilibrium. (1) amounts to a decentralized projected gradient procedure. It portrays fairly myopic parties, each trying to improve his linear approximation of own expected payo¤. Modern, stochastic versions of such methods mirror two common aspects of human behavior: …rst, mean values (i.e. mathematical expectations $E$) are costly - and sometimes impossible - to compute; second, information concerning levels and gradients is readily available only at the current point. So, in our optic, letting $M_i(x; !) := \frac{@}{@x_i} ¦_i(x; !)$ denote $i$'s realized marginal payo¤, one might hope to have almost sure convergence of the following stochastic process: For each $i$ recursively posit

$$x_i^{k+1} := P_i \left[ x_i^k + s_k M_i(x^k; !^k) \right]: \tag{2}$$

Here, at stage $k = 0; 1; :::;$ arrives a new event $!^k \; 2 \; \text{--}$, independently sampled according to the prescribed measure $^1$: Input at that stage $k$ is also a positive stepsize $s_k$; selected a priori subject to

$$\sum s_k = +1 \text{ and } \sum s_k^2 < +1: \tag{3}$$

The hope that process (2) converges is well founded provided $x \; \mathcal{V} \; m(x) := [m_i(x)]_{i2I}$ be globally monotone; see [9], [10], [11], [12], [13]. Otherwise, there are good reasons

to be worried about convergence. Re‡ecting such worries our object here is to expand on gradient methods while preserving their many appealing properties.

Like (2) the procedure considered below does not presume much of foresight, experience, competence or optimization. In essence, it re‡ects iterated, noncoordinated pursuit of better payo¤s. It does, however, modify the …rst-order gradient dynamics (2) by adding a second-order heavy-ball momentum, just like the harmonic oscillator of classical mechanics. Essentially, instead of assuming that player i pursues the gradient method $0 = m_i(x) ¡ \underline{x}_i$ we posit that he rather drives the second-order process $\ddot{x}_i = m_i(x) ¡ \underline{x}_i$: The latter must be suitably modi…ed, of course, to account for discrete time, uncertainty, and constraints. This is done next.

## 3. Repeated Play

Let $\text{!}^k$ be a sequence of independent realizations of $\text{!}$; each having distribution $^1$: As model of repeated play we advocate that iteratively at stages $k = 0; 1; ::;$ each individual i updates his current strategy $x_i^k$ and velocity $v_i^k$ by the rule

$$
\begin{aligned}
x_i^{k+1} &:= P_i \left[ x_i^k + s_k v_i^k \right] \\
v_i^{k+1} &:= v_i^k + P_i \left[ x_i^k + s_k M_i(x^k; \text{!}^k) \right] ¡ P_i \left[ x_i^k + s_k v_i^k \right]
\end{aligned}
\tag{4}
$$

As earlier, $P_i$ denotes orthogonal projection onto $X_i$. Also like above, the parameter $s_k > 0$ is the stepsize used at stage k; selected a priori subject to (3). The initial points $(x_i^0; v_i^0); i \in I;$ are determined by accident or historical factors better discussed in each particular setting.

To appreciate process (4) it helps to endow it with a clock that shows accumulated "time" $t_k := s_0 + ¢¢¢ + s_{k¡1}; (¿_0 := 0)$ at the on-set of stage k: Then, upon writing $x_i(t_k) := x_i^k$ and $v_i(t_k) := v_i^k$ we see that (4) assumes the form

$$
\begin{aligned}
\{x_i(t_{k+1}) ¡ x_i(t_k)\} = s_k &:= \{P_i [x_i(t_k) + s_k v_i(t_k)] ¡ x_i(t_k)\} = s_k \\
\{v_i(t_{k+1}) ¡ v_i(t_k)\} = s_k &:= P_i [x_i(t_k) + s_k M_i(x(t_k); \text{!}^k) ] ¡ P_i [x_i(t_k) + s_k v_i(t_k)] = s_k
\end{aligned}
$$

Thus, since $s_k = t_{k+1} ¡ t_k \to 0^+;$ it turns out that behind (4) lurks - in expectation and the limit - a di¤erential system

$$
\begin{aligned}
\underline{x}_i &= P_{T_i x_i} [v_i] \\
\underline{v}_i &= P_{T_i x_i} [m_i(x)] ¡ P_{T_i x_i} [v_i]
\end{aligned}
\tag{5}
$$

Orthogonal projection $P_{T_i x_i}$ is here done onto the tangent cone $T_i x_i := clR_+(X_i ¡ x_i)$ of $X_i$ at $x_i$: By a solution to (5) we understand an absolutely continuous pro…le $0 \cdot t \mapsto [x(t); v(t)] = [x_i(t); v_i(t)]_{i \in I}$ which satis…es (5) almost everywhere. We suppose that the potential energy

$$
0 \cdot t \mapsto \int_0^t \sum_{i \in I} P_{T_i x_i(¿)} [m_i(x(¿))] ¢ \underline{x}_i(¿) d¿
\tag{6}
$$

remains bounded above along solution trajectories of (5). Following the arguments in [14] one may prove the following

**Theorem 1.** (Convergence of repeated play) Suppose system (5) has unique solution trajectories. Then, under the hypotheses above and the assumption that Nash equilibria are isolated, any discrete-time trajectory $(x^k; v^k)$ generated by (4) must be such that $x^k$ converges to a Nash equilibrium. 2

Examples and illustrations of process (4), when deterministic, are found in [14]. We remark that at any stage k player i might, quite reasonably, ...rst update his velocity $v_i^{k+1}$ as prescribed and thereafter set $x_i^{k+1} := P_i\,x_i^k + s_k v_i^{k+1}$ : In subsequent simulations we observe that this practice speeds up convergence.

## 4. Time-Homogeneous Play

When telling our tale about repeated play, we ...nd it di₵cult sometimes to argue in favor of a time-inhomogeneous system. So, what happens if $s_k$ is constant? Clearly, ...xing this parameter is risky - and notably with sti¤ systems. To address that issue assume, for simplicity, that there are no constraints - or alternatively, that these have already been incorporated by means of suitable penalties. Setting $d_i = d_i(x; v_i; !) := M_i(x; !) \; ; \; v_i;$ iteration (4) then comes in autonomous, more tractable form

$$
\begin{aligned}
x_i^{k+1} &:= x_i^k + s v_i^k \\
v_i^{k+1} &:= v_i^k + s d_i^k
\end{aligned}
\tag{7}
$$

In simulations of (7), to hedge against sti¤ness and facilitate convergence, we replace $d_i^k$ with a weighted sum $a_i d_i^k + b_i d_i^{k\,i\,1}$; using thus

$$
\begin{aligned}
x_i^{k+1} &:= x_i^k + s v_i^k \\
v_i^{k+1} &:= v_i^k + s(a_i d_i^k + b_i d_i^{k\,i\,1})
\end{aligned}
\tag{8}
$$

The parameters $a_i; b_i$ could account for accumulated learning on how to adapt in a complex dynamic environment.[1] We remark that the second equation in (8) is strikingly similar to a control algorithm commonly used in process industry, namely: the so-called Proportional-Integral-Controller; see [5], [17], [23], [24]. We take advantage of this to ...nd appropriate $a_i; b_i$-parameters in the example below:

## 5. An Example

Let each $i \, 2 \, I$ here be a Cournot oligopolist [8] who supplies the quantity $x_i \, ,$ 0 of a homogeneous, perfectly divisible good to a common market. Thereby he receives sales revenues $p x_i$ and incurs (di¤erentiable) production costs $c_i(x_i)$. The price is determined by a smooth inverse demand curve which is subject to stochastic

---

[1]In fact, if agent i were new to the kind of dynamic process in question, his optimal behavior might entail experimentation to determine the said parameters. Most likely agents would learn from each other and, if possible, imitate those who do well. Henceforth assume that each i has previous experience, perhaps from similar processes, and has appropriately tuned his $a_i$ and $b_i$.

‡uctuations. Speci...cally, the realized price equals $p = !\, P\,(Q)$ with $E! = 1$ and $Q = \sum_{i2I} x_i$. Thus

$$E\,|_i\,(x_i; x_{i\,i}; !) = P\,(Q)\,x_i\,i\,c_i(x_i) \text{ and } EM_i\,(x_i; x_{i\,i}; !) = P\,(Q) + P^0\,(Q)\,x_i\,i\,c_i^0(x_i)$$

The ...gures below depict individual supply $x_i$ over stages $k$ with constant stepsize $s = 1$, as generated by (8). There are ten players ($jIj = 10$), $!$ has a lognormal distribution with standard deviation $0:3$ , $P\,(Q) = 10\,i\,Q$, $c_i\,(x_i) = x_i$, $x_i^0 = 1$, $v_i^0 = 0$ and ...nally, $d_i^0 = 0$ for all $i$. The resulting Nash Equilibrium $x$ is unique with all $x_i = 0:818$.

The resemblance with controller algorithms made us look for methods to determine e¢cient values of $a_i$ and $b_i$ in the literature on Control Engineering. The wide spread use of such algorithms not withstanding, there exists no generally accepted method for how to tune these parameters. The empirical method developed by Ziegler and Nichols (1942) [24] is still holds good ground and has the great advantage of requiring very little information [23]. The said method gave the values of $a_i$ and $b_i$ that label Figure 1 and are used by all agents. The scantly dotted line shows behavior in the deterministic case (when $! , 1$) while in the more solid curve $!$ is sampled anew for each $k$.

Figure 2 brings out responses when all agents half the values $a_i$ and $b_i$ employed in Figure 1. Absent uncertainty, the time needed to reach a steady level now becomes longer. Present uncertainty, (using the same series $!^k$ as in Figure 1) it causes less ‡uctuations than before.
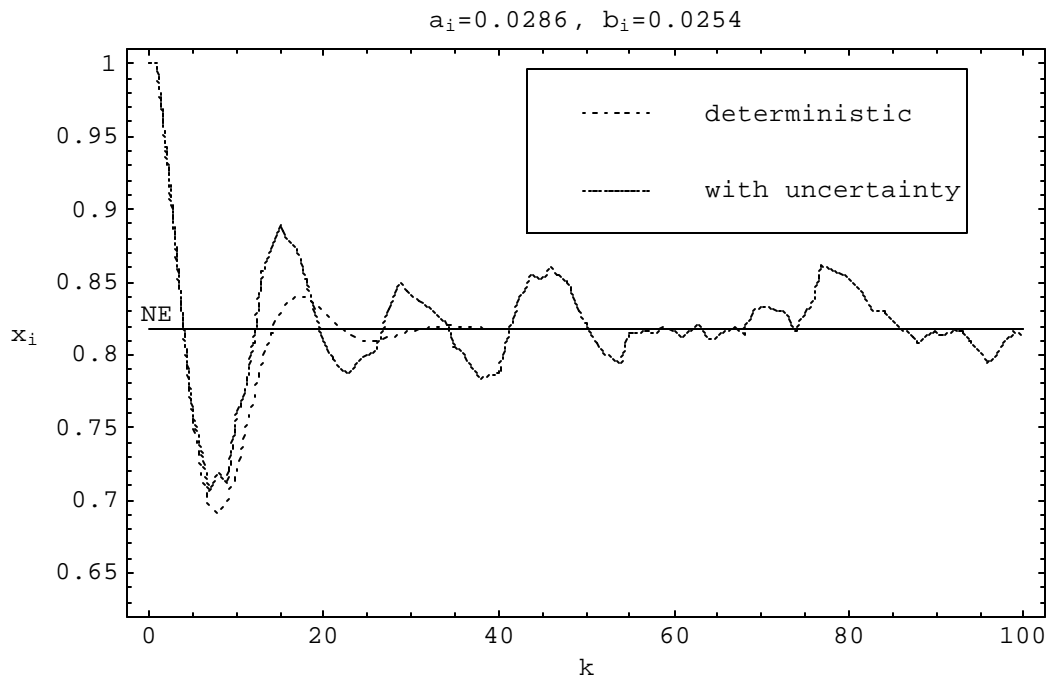


Figure 1: Supply $x_i$ over stages $k$ when all players employ parameters $a_i$ and $b_i$ determined by Ziegler-Nichols Method.
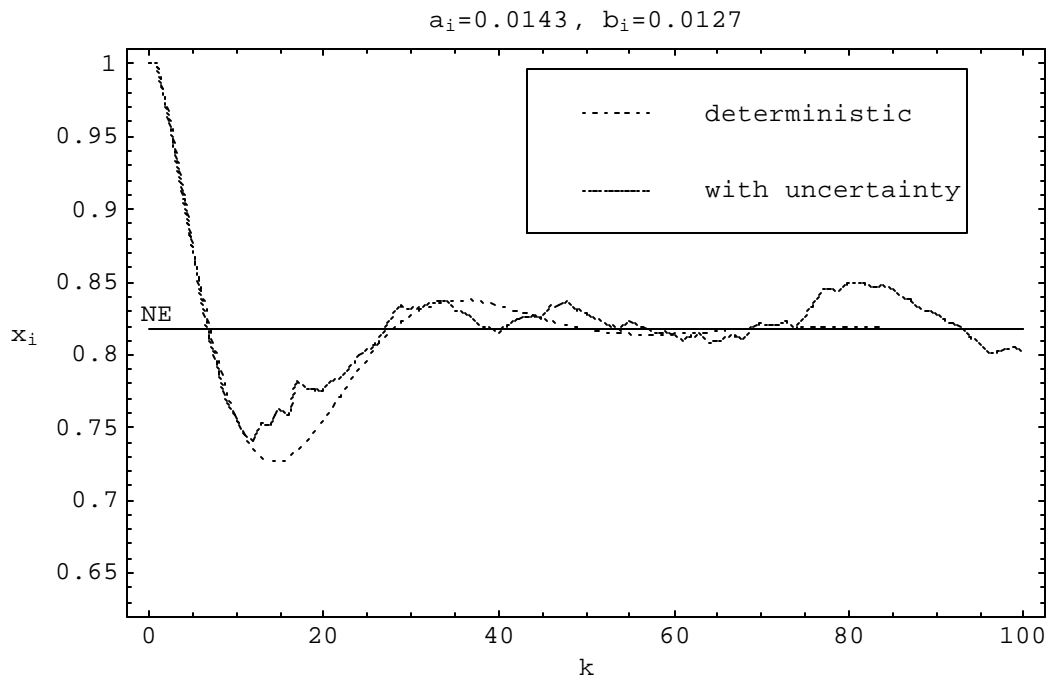
Figure 2: Supply $x_i$ over stages $k$ when all players employ parameters $a_i$ and $b_i$ with values half of the size determined by Ziegler-Nichols Method.

Figure 3 illustrates what happens when parameters $a_i; b_i$ di¤er across agents. The ...ve ...rst players use values listed in Figure 1; the others use those mentioned in Figure 2. Members of the ...rst group adapt fastest initially. Di¤erences across agents make for slower convergence in the deterministic case.
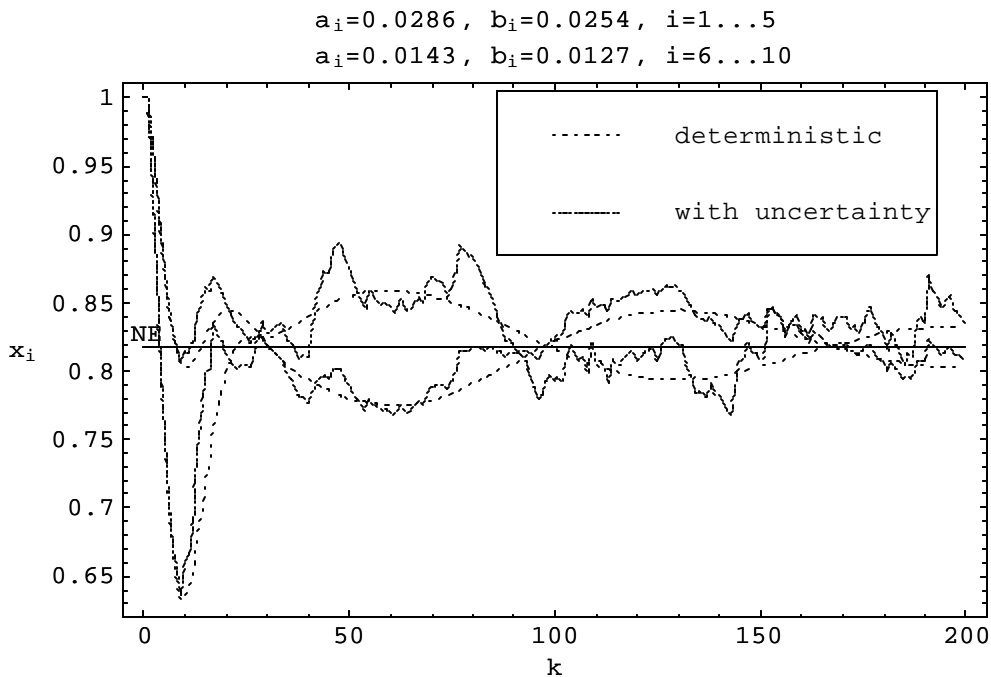
$$a_i=0.0286, \quad b_i=0.0254, \quad i=1\ldots5$$
$$a_i=0.0143, \quad b_i=0.0127, \quad i=6\ldots10$$



Figure 3: Supply $x_i$ over stages $k$ with di¤erent $a_i$; $b_i$. Group 1 uses parameters $a_i$ and $b_i$ determined by Ziegler-Nichols method - while group 2 half these values. Members of the …rst group adapt fastest initially.

## 6.  Concluding Remarks

We have presented (4) as a tale about repeated play of noncooperative, constrained games. A main motivation behind the heavy-ball philosophy was to reach beyond instances where the solution set is connected. A supplementary aim - referred to as equilibrium selection [21] - amounts to exploit uncertainty or blurred data so as to arrive at particularly stable solutions. In fact, randomness - if not already a key ingredient - could arti…cially be introduced to escape from unstable equilibria.

Clearly, when I is a singleton, this paper …ts the frames of single-agent optimization under uncertainty. In that regard (4) has something to o¤er in three respects. First, the heavy-ball method of Polyak [20] has, to our knowledge, not come fully into stochastic programming. Second, process (7) is amenable to parallel computing [7]. Third, following Attouch et al.[3], the same process is applicable for global optimization - or the selection of "good" stationary points; see also [1], [2], [4]. Approach (8) may require, for its e¢ cient operation, some auto-tuning of the parameters $a_i$; $b_i$: Appropriate routines to that e¤ect are found in the engineering literature on control; see for instance [5], [23].

## References

[1] F. Alvarez, On the minimizing property of a second order dissipative system in Hilbert spaces, SIAM J. Control Opt. 38, 4, 1102-1119 (2000).

[2] F. Alvarez and J. M. Pérez, A dynamical system associated with Newton's method for parametric approximations of convex minimization problem, Appl. Math. Optm. 38, 193-217 (1998).

[3] H. Attouch, X. Goudou, and P. Redont, The heavy ball with friction method I, The continuous dynamical system: Global exploration of the local minima of a real-valued function by asymptotic analysis of a dissipative dynamical system, Communications in Contemporary Mathematics 2, 1, 1-34 (2000).

[4] H. Attouch and P. Redont, The second-order in time continuous Newton method, in M. Lassonde, Approximation, Optimization and Mathematical Economic, Physica-Verlag, Heidelberg 4-36 (2001).

[5] R. Bandyopadhyay and D. Patranabis, A fuzzy logic based PI auto-tuner, ISA Transactions 37, 227-235 (1998).

[6] M. Benaim, A dynamical system approach to stochastic approximation, SIAM J. of Control and Optimization 34, 437-472 (1996).

[7] D. P. Bertsekas and J. N. Tsitsiklis, Parallel and Distributed Computation: Numerical Methods, Prentice-Hall, New York (1989).

[8] A. Cournot, Recherches sur les principes mathématiques de la théorie des richesses, Riviere & Cie, Paris (1838).

[9] S. D. Flåm, Approaches to economic equilibrium, Journal of Economic Dynamics and Control 20, 1505-1522 (1996).

[10] S. D. Flåm, Restricted attention, myopic play, and the learning of equilibrium, Annals of Operations Research 82, 473-482 (1998) .

[11] S. D. Flåm, Learning equilibrium play: a myopic approach, Computational Optimization and Appl.14, 87-102 (1999) .

[12] S. D. Flåm, Repeated play and Newton's method, Int. Game Theory Review 2, 141-154 (2001).

[13] S. D. Flåm, Approaching equilibrium in parallel, in D. Butnariu, Y. Censor and S. Reich (eds.) Inherently Parallel Algorithms in Feasibility and Optimization and their Applications, North-Holland 267-278 (2001).

[14] S. D. Flåm and J. Morgan, Newtonian mechanics and Nash Play, manuscript (2001).

[15] D. Fudenberg and D. K. Levine, The Theory of Learning in Games, MIT Press, Cambridge Mass. (1998).

[16] J. K. Goeree and C. A. Holt, Ten little treasures of game theory and ten intuitive contradictions, The American Economic Review 91, 5, 1402-1422 (2001).

[17] C. C. Hang, K. J. Åstrøm, Q. G. Wang, Relay feedback auto-tuning of process controllers – a tutorial review, Journal of Process Control 12, 143-162 (2002).

[18] J. Hofbauer and K. Sigmund, Evolutionary Games and Population Dynamics, Cambridge University Press (1998).

[19] J. Hofbauer, From Nash and Brown to Maynard Smith: Equilibria, dynamics and ESS, Selection 1, 81-88 (2000).

[20] B. T. Polyak, Some methods of speeding up the convergence of iteration methods, Z. VyCisl. Math. i Mat. Fiz. 4, 1-17 (1964).

[21] L. Samuelson, Evolutionary Games and Equilibrium Selection, The MIT Press, Cambridge, Massachusetts (1997).

[22] J. Watson, Strategy, Norton, New York (2002).

[23] K. J. Åström, H. Panagopoulos and T. Hägglund, Design of PI Controllers based on Non-Convex Optimization, Automatica 34, 5, 585-601 (1998).

[24] J. G. Ziegler and N. B. Nichols, Optimum settings for automatic controllers, Trans. ASME 64, 759-768 (1942).