

# Dynamics of Learning Summary - Hanchun Wang

Tuesday, July 19, 2022 1:39 PM

<p>Equilibrium</p>	<table border="1"> <tr> <td>Action domain</td> <td><math>\Delta_n = \{x \in \mathbb{R}^n; 0 \leq x_i \leq 1, x_1 + \dots + x_n = 1\}</math> be the <math>(n - 1)</math> dimensional simplex</td> </tr> <tr> <td>Linear Payoff</td> <td><math>x \cdot Ay, e_i \cdot Ax = (Ax)_i</math></td> </tr> <tr> <td>Indifference set</td> <td><math>Z_{ij} = \{x; (Ax)_i = (Ax)_j\}</math></td> </tr> <tr> <td>Best Response</td> <td><math>BR(x) = \operatorname{argmax}_{y \in \Delta} (y \cdot Ax) \stackrel{\text{def}}{=} \{y' \in \Delta; y' \cdot Ax \leq y \cdot Ax, \forall y \in \Delta\}</math></td> </tr> <tr> <td>(strict) Nash Equilibrium</td> <td> <math>\hat{x} \in \Delta</math> is NE <math>\stackrel{\text{def}}{=} x \cdot A\hat{x} \leq \hat{x} \cdot Ax, \forall x \in \Delta</math>, strict inequality for strict NE  <math>\hat{x} \in \Delta</math> is NE <math>\leftrightarrow \hat{x} \in BR(\hat{x})</math>  <math>\hat{x} \in \Delta</math> is NE <math>\leftrightarrow \exists c \in \mathbb{R}, \hat{x}_i &gt; 0, (A\hat{x})_i = c \rightarrow \forall ij, \hat{x} \in Z_{ij}</math> </td> </tr> <tr> <td>Evolutionary stable equilibrium (ESS)</td> <td> <math>\hat{x} \in \Delta</math> is ESS <math>\stackrel{\text{def}}{=} \exists \epsilon &gt; 0, x \cdot A(\epsilon x + (1 - \epsilon)\hat{x}) &lt; \hat{x} \cdot A(\epsilon x + (1 - \epsilon)\hat{x})</math>  <math>\hat{x} \in \Delta</math> is ESS <math>\stackrel{\text{def}}{=} \exists \epsilon &gt; 0, \forall y \in B_\epsilon(x), y \cdot Ay &lt; \hat{x} \cdot Ay</math>  <math>\hat{x} \in \operatorname{int}(\Delta)</math> is an ESS <math>\rightarrow \hat{x}</math> is the unique NE                 </td> </tr> <tr> <td>Coarse Correlated Equilibrium (CCE, no-regret set)</td> <td><math>(p_{ij}) \in CCE \leftrightarrow \sum_{i,j} a_{i'j} p_{ij} \leq \sum_{i,j} a_{ij} p_{ij}</math> and <math>\sum_{i,j} b_{ij'} p_{ij} \leq \sum_{i,j} b_{ij} p_{ij}</math></td> </tr> <tr> <td>Correlated Equilibrium (CE)</td> <td><math>(p_{ij}) \in CE \leftrightarrow \sum_k a_{i'k} p_{ik} \leq \sum_k a_{ik} p_{ik}</math> and <math>\sum_l b_{lj'} p_{lj} \leq \sum_l b_{lj} p_{lj}</math></td> </tr> <tr> <td></td> <td> <math>\{\text{strict NE}\} \subset \{\text{ESS}\} \subset \{\text{NE}\}</math>  <math>\{\text{NE}\} \subset \{\text{CE}\} \subset \{\text{CCE}\}</math> </td> </tr> </table>	Action domain	$\Delta_n = \{x \in \mathbb{R}^n; 0 \leq x_i \leq 1, x_1 + \dots + x_n = 1\}$ be the $(n - 1)$ dimensional simplex	Linear Payoff	$x \cdot Ay, e_i \cdot Ax = (Ax)_i$	Indifference set	$Z_{ij} = \{x; (Ax)_i = (Ax)_j\}$	Best Response	$BR(x) = \operatorname{argmax}_{y \in \Delta} (y \cdot Ax) \stackrel{\text{def}}{=} \{y' \in \Delta; y' \cdot Ax \leq y \cdot Ax, \forall y \in \Delta\}$	(strict) Nash Equilibrium	$\hat{x} \in \Delta$ is NE $\stackrel{\text{def}}{=} x \cdot A\hat{x} \leq \hat{x} \cdot Ax, \forall x \in \Delta$ , strict inequality for strict NE $\hat{x} \in \Delta$ is NE $\leftrightarrow \hat{x} \in BR(\hat{x})$ $\hat{x} \in \Delta$ is NE $\leftrightarrow \exists c \in \mathbb{R}, \hat{x}_i > 0, (A\hat{x})_i = c \rightarrow \forall ij, \hat{x} \in Z_{ij}$	Evolutionary stable equilibrium (ESS)	$\hat{x} \in \Delta$ is ESS $\stackrel{\text{def}}{=} \exists \epsilon > 0, x \cdot A(\epsilon x + (1 - \epsilon)\hat{x}) < \hat{x} \cdot A(\epsilon x + (1 - \epsilon)\hat{x})$ $\hat{x} \in \Delta$ is ESS $\stackrel{\text{def}}{=} \exists \epsilon > 0, \forall y \in B_\epsilon(x), y \cdot Ay < \hat{x} \cdot Ay$ $\hat{x} \in \operatorname{int}(\Delta)$ is an ESS $\rightarrow \hat{x}$ is the unique NE	Coarse Correlated Equilibrium (CCE, no-regret set)	$(p_{ij}) \in CCE \leftrightarrow \sum_{i,j} a_{i'j} p_{ij} \leq \sum_{i,j} a_{ij} p_{ij}$ and $\sum_{i,j} b_{ij'} p_{ij} \leq \sum_{i,j} b_{ij} p_{ij}$	Correlated Equilibrium (CE)	$(p_{ij}) \in CE \leftrightarrow \sum_k a_{i'k} p_{ik} \leq \sum_k a_{ik} p_{ik}$ and $\sum_l b_{lj'} p_{lj} \leq \sum_l b_{lj} p_{lj}$		$\{\text{strict NE}\} \subset \{\text{ESS}\} \subset \{\text{NE}\}$ $\{\text{NE}\} \subset \{\text{CE}\} \subset \{\text{CCE}\}$
Action domain	$\Delta_n = \{x \in \mathbb{R}^n; 0 \leq x_i \leq 1, x_1 + \dots + x_n = 1\}$ be the $(n - 1)$ dimensional simplex																		
Linear Payoff	$x \cdot Ay, e_i \cdot Ax = (Ax)_i$																		
Indifference set	$Z_{ij} = \{x; (Ax)_i = (Ax)_j\}$																		
Best Response	$BR(x) = \operatorname{argmax}_{y \in \Delta} (y \cdot Ax) \stackrel{\text{def}}{=} \{y' \in \Delta; y' \cdot Ax \leq y \cdot Ax, \forall y \in \Delta\}$																		
(strict) Nash Equilibrium	$\hat{x} \in \Delta$ is NE $\stackrel{\text{def}}{=} x \cdot A\hat{x} \leq \hat{x} \cdot Ax, \forall x \in \Delta$ , strict inequality for strict NE $\hat{x} \in \Delta$ is NE $\leftrightarrow \hat{x} \in BR(\hat{x})$ $\hat{x} \in \Delta$ is NE $\leftrightarrow \exists c \in \mathbb{R}, \hat{x}_i > 0, (A\hat{x})_i = c \rightarrow \forall ij, \hat{x} \in Z_{ij}$																		
Evolutionary stable equilibrium (ESS)	$\hat{x} \in \Delta$ is ESS $\stackrel{\text{def}}{=} \exists \epsilon > 0, x \cdot A(\epsilon x + (1 - \epsilon)\hat{x}) < \hat{x} \cdot A(\epsilon x + (1 - \epsilon)\hat{x})$ $\hat{x} \in \Delta$ is ESS $\stackrel{\text{def}}{=} \exists \epsilon > 0, \forall y \in B_\epsilon(x), y \cdot Ay < \hat{x} \cdot Ay$ $\hat{x} \in \operatorname{int}(\Delta)$ is an ESS $\rightarrow \hat{x}$ is the unique NE																		
Coarse Correlated Equilibrium (CCE, no-regret set)	$(p_{ij}) \in CCE \leftrightarrow \sum_{i,j} a_{i'j} p_{ij} \leq \sum_{i,j} a_{ij} p_{ij}$ and $\sum_{i,j} b_{ij'} p_{ij} \leq \sum_{i,j} b_{ij} p_{ij}$																		
Correlated Equilibrium (CE)	$(p_{ij}) \in CE \leftrightarrow \sum_k a_{i'k} p_{ik} \leq \sum_k a_{ik} p_{ik}$ and $\sum_l b_{lj'} p_{lj} \leq \sum_l b_{lj} p_{lj}$																		
	$\{\text{strict NE}\} \subset \{\text{ESS}\} \subset \{\text{NE}\}$ $\{\text{NE}\} \subset \{\text{CE}\} \subset \{\text{CCE}\}$																		
<p>Replicator Dynamics (one population)</p>	<table border="1"> <tr> <td>Dynamics</td> <td> <math>\dot{x}_i = x_i((Ax)_i - x \cdot Ax), i = 1, \dots, n</math>  <math>\frac{d}{dt} \frac{x_i}{x_j} = \frac{x_i}{x_j}((Ax)_i - (Ax)_j)</math> </td> </tr> <tr> <td>Convergence</td> <td> <p><b>Every <math>n \times n</math> replicator game has a Nash Equilibrium</b></p> <p>NE is an equilibrium of Replicator dynamics  <math>\hat{x}</math> is omega-limit of an orbit <math>x(t), \hat{x} \in \operatorname{int}(\Delta)</math>. Then <math>\hat{x}</math> is NE.  <math>\hat{x}</math> is lyapunov stable. Then <math>\hat{x}</math> is NE.</p> <p>ESS is asymptotically stable equilibrium  <math>\hat{x} \in \operatorname{int}(\Delta)</math> is ESS, then <math>\hat{x}</math> globally attracts all initial points <math>x \in \operatorname{int}(\Delta)</math></p> </td> </tr> </table>	Dynamics	$\dot{x}_i = x_i((Ax)_i - x \cdot Ax), i = 1, \dots, n$ $\frac{d}{dt} \frac{x_i}{x_j} = \frac{x_i}{x_j}((Ax)_i - (Ax)_j)$	Convergence	<p><b>Every <math>n \times n</math> replicator game has a Nash Equilibrium</b></p> <p>NE is an equilibrium of Replicator dynamics  <math>\hat{x}</math> is omega-limit of an orbit <math>x(t), \hat{x} \in \operatorname{int}(\Delta)</math>. Then <math>\hat{x}</math> is NE.  <math>\hat{x}</math> is lyapunov stable. Then <math>\hat{x}</math> is NE.</p> <p>ESS is asymptotically stable equilibrium  <math>\hat{x} \in \operatorname{int}(\Delta)</math> is ESS, then <math>\hat{x}</math> globally attracts all initial points <math>x \in \operatorname{int}(\Delta)</math></p>														
Dynamics	$\dot{x}_i = x_i((Ax)_i - x \cdot Ax), i = 1, \dots, n$ $\frac{d}{dt} \frac{x_i}{x_j} = \frac{x_i}{x_j}((Ax)_i - (Ax)_j)$																		
Convergence	<p><b>Every <math>n \times n</math> replicator game has a Nash Equilibrium</b></p> <p>NE is an equilibrium of Replicator dynamics  <math>\hat{x}</math> is omega-limit of an orbit <math>x(t), \hat{x} \in \operatorname{int}(\Delta)</math>. Then <math>\hat{x}</math> is NE.  <math>\hat{x}</math> is lyapunov stable. Then <math>\hat{x}</math> is NE.</p> <p>ESS is asymptotically stable equilibrium  <math>\hat{x} \in \operatorname{int}(\Delta)</math> is ESS, then <math>\hat{x}</math> globally attracts all initial points <math>x \in \operatorname{int}(\Delta)</math></p>																		
<p>Replicator Dynamics two player</p>	<table border="1"> <tr> <td>Dynamics</td> <td> <math>P_A(x, y) = x \cdot Ay, BR_A(y) = \operatorname{argmax}_{x \in \Delta_A} x \cdot Ay</math>  <math>P_B(x, y) = x \cdot By, BR_B(x) = \operatorname{argmax}_{y \in \Delta_B} x \cdot By</math>  <math>\dot{x}_i = x_i((Ay)_i - x \cdot Ay)</math>  <math>\dot{y}_j = y_j((x^{tr} B)_j - x \cdot By)</math> </td> </tr> <tr> <td>Convergence</td> <td>Each bimatrix game (A;B) has a Nash equilibrium</td> </tr> </table>	Dynamics	$P_A(x, y) = x \cdot Ay, BR_A(y) = \operatorname{argmax}_{x \in \Delta_A} x \cdot Ay$ $P_B(x, y) = x \cdot By, BR_B(x) = \operatorname{argmax}_{y \in \Delta_B} x \cdot By$ $\dot{x}_i = x_i((Ay)_i - x \cdot Ay)$ $\dot{y}_j = y_j((x^{tr} B)_j - x \cdot By)$	Convergence	Each bimatrix game (A;B) has a Nash equilibrium														
Dynamics	$P_A(x, y) = x \cdot Ay, BR_A(y) = \operatorname{argmax}_{x \in \Delta_A} x \cdot Ay$ $P_B(x, y) = x \cdot By, BR_B(x) = \operatorname{argmax}_{y \in \Delta_B} x \cdot By$ $\dot{x}_i = x_i((Ay)_i - x \cdot Ay)$ $\dot{y}_j = y_j((x^{tr} B)_j - x \cdot By)$																		
Convergence	Each bimatrix game (A;B) has a Nash equilibrium																		
<p>Best response Dynamics</p>	<table border="1"> <tr> <td>Dynamics</td> <td> <math>\dot{x} \in BR(x) - x;</math>  <math>\dot{x} = BR(x) - x</math>, at differentiable <math>x</math>  <math>x \rightarrow BR(x)</math> is upper semi-continuous                      Characteristic curve is continuous, almost everywhere differentiable                 </td> </tr> <tr> <td>Convergence</td> <td>Assume that (A,B) is a <b>zero-sum game</b>. Best Response dynamics converge to the set of <b>Nash equilibria</b>.</td> </tr> </table>	Dynamics	$\dot{x} \in BR(x) - x;$ $\dot{x} = BR(x) - x$ , at differentiable $x$ $x \rightarrow BR(x)$ is upper semi-continuous Characteristic curve is continuous, almost everywhere differentiable	Convergence	Assume that (A,B) is a <b>zero-sum game</b> . Best Response dynamics converge to the set of <b>Nash equilibria</b> .														
Dynamics	$\dot{x} \in BR(x) - x;$ $\dot{x} = BR(x) - x$ , at differentiable $x$ $x \rightarrow BR(x)$ is upper semi-continuous Characteristic curve is continuous, almost everywhere differentiable																		
Convergence	Assume that (A,B) is a <b>zero-sum game</b> . Best Response dynamics converge to the set of <b>Nash equilibria</b> .																		
<p>Fictional play dynamics</p>	<table border="1"> <tr> <td>Dynamics</td> <td> <math>p(s) = \frac{1}{s} \int_0^s x(u) du; q(s) = \frac{1}{s} \int_0^s y(u) du; s = e^t</math>  <math>\dot{p}(s) = \frac{1}{s} x(s) - \frac{1}{s} p(s)</math> and <math>\dot{q}(s) = \frac{1}{s} y(s) - \frac{1}{s} q(s)</math>  <math>x(s) \in BR_A(q(s))</math> and <math>y(s) \in BR_B(p(s))</math> for <math>s \geq 1</math>.  <math>\dot{p}(s) \in \frac{1}{s} (BR_A(q(s)) - p(s)); \dot{q}(s) \in \frac{1}{s} (BR_B(p(s)) - q(s))</math>  <math>\dot{p}(t) = (BR_A(q(t)) - p(t)); \dot{q}(t) = (BR_B(p(t)) - q(t)).</math> </td> </tr> </table>	Dynamics	$p(s) = \frac{1}{s} \int_0^s x(u) du; q(s) = \frac{1}{s} \int_0^s y(u) du; s = e^t$ $\dot{p}(s) = \frac{1}{s} x(s) - \frac{1}{s} p(s)$ and $\dot{q}(s) = \frac{1}{s} y(s) - \frac{1}{s} q(s)$ $x(s) \in BR_A(q(s))$ and $y(s) \in BR_B(p(s))$ for $s \geq 1$ . $\dot{p}(s) \in \frac{1}{s} (BR_A(q(s)) - p(s)); \dot{q}(s) \in \frac{1}{s} (BR_B(p(s)) - q(s))$ $\dot{p}(t) = (BR_A(q(t)) - p(t)); \dot{q}(t) = (BR_B(p(t)) - q(t)).$																
Dynamics	$p(s) = \frac{1}{s} \int_0^s x(u) du; q(s) = \frac{1}{s} \int_0^s y(u) du; s = e^t$ $\dot{p}(s) = \frac{1}{s} x(s) - \frac{1}{s} p(s)$ and $\dot{q}(s) = \frac{1}{s} y(s) - \frac{1}{s} q(s)$ $x(s) \in BR_A(q(s))$ and $y(s) \in BR_B(p(s))$ for $s \geq 1$ . $\dot{p}(s) \in \frac{1}{s} (BR_A(q(s)) - p(s)); \dot{q}(s) \in \frac{1}{s} (BR_B(p(s)) - q(s))$ $\dot{p}(t) = (BR_A(q(t)) - p(t)); \dot{q}(t) = (BR_B(p(t)) - q(t)).$																		

	Convergence	<p><b>Fictitious play converges to the no-regret set CCE</b></p> <p>Assume that (A,B) is a zero-sum game. Fictitious Play dynamics converge to the set of <b>Nash equilibria</b>.</p> <p><b>Fiction Play orbits Pareto dominates Nash payoff</b></p> <p>Time averages of Replicator Dynamics converge to pseudo-orbits of Fictitious Play</p>						
Reinforcement Learning	Dynamics	<table border="1"> <tr> <td data-bbox="558 259 714 442">Model update</td> <td data-bbox="717 259 974 442"> <p>Cross Learning</p> <p>Erev-Roth Cumulative payoff matching (CPM)</p> <p>Arthur model</p> </td> <td data-bbox="977 259 1425 442"> <p><math>\theta^{t+1} = (1 - \vartheta u^t)\theta^t + \vartheta u^t x^t, \quad t \geq 1</math></p> <p><math>\theta^{t+1} = \theta^t + u^t x^t, \quad t \geq 1</math></p> <p><math>\theta^{t+1} = (\theta^t + u^t x^t) \frac{C(t+1)}{Ct + u^t}, \quad t \geq 1</math></p> </td> </tr> <tr> <td data-bbox="558 446 714 608">Action choose</td> <td data-bbox="717 446 974 608"> <p>Proportional</p> <p>Greedy</p> <p>Softmax</p> </td> <td data-bbox="977 446 1425 608"> <p><math>p(t) = Q^t /  Q^t </math></p> <p><math>p(Q) = (1 - \epsilon)BR_t(Q) + \epsilon(1/n, \dots, 1/n)</math></p> <p><math>\text{softmax}_T(Q) = \frac{1}{\sum_i \exp(Q_i/T)} (\exp(Q_1/T), \dots, \exp(Q_n/T))</math></p> </td> </tr> </table>	Model update	<p>Cross Learning</p> <p>Erev-Roth Cumulative payoff matching (CPM)</p> <p>Arthur model</p>	<p><math>\theta^{t+1} = (1 - \vartheta u^t)\theta^t + \vartheta u^t x^t, \quad t \geq 1</math></p> <p><math>\theta^{t+1} = \theta^t + u^t x^t, \quad t \geq 1</math></p> <p><math>\theta^{t+1} = (\theta^t + u^t x^t) \frac{C(t+1)}{Ct + u^t}, \quad t \geq 1</math></p>	Action choose	<p>Proportional</p> <p>Greedy</p> <p>Softmax</p>	<p><math>p(t) = Q^t /  Q^t </math></p> <p><math>p(Q) = (1 - \epsilon)BR_t(Q) + \epsilon(1/n, \dots, 1/n)</math></p> <p><math>\text{softmax}_T(Q) = \frac{1}{\sum_i \exp(Q_i/T)} (\exp(Q_1/T), \dots, \exp(Q_n/T))</math></p>
Model update	<p>Cross Learning</p> <p>Erev-Roth Cumulative payoff matching (CPM)</p> <p>Arthur model</p>	<p><math>\theta^{t+1} = (1 - \vartheta u^t)\theta^t + \vartheta u^t x^t, \quad t \geq 1</math></p> <p><math>\theta^{t+1} = \theta^t + u^t x^t, \quad t \geq 1</math></p> <p><math>\theta^{t+1} = (\theta^t + u^t x^t) \frac{C(t+1)}{Ct + u^t}, \quad t \geq 1</math></p>						
Action choose	<p>Proportional</p> <p>Greedy</p> <p>Softmax</p>	<p><math>p(t) = Q^t /  Q^t </math></p> <p><math>p(Q) = (1 - \epsilon)BR_t(Q) + \epsilon(1/n, \dots, 1/n)</math></p> <p><math>\text{softmax}_T(Q) = \frac{1}{\sum_i \exp(Q_i/T)} (\exp(Q_1/T), \dots, \exp(Q_n/T))</math></p>						
Q-learning	Dynamics	<p><math>Q^{t+h}(s) = Q^t(s) + ah \cdot (u_A^t + \gamma \max_j Q_j^t(s) - Q^t(s) \cdot a(t)) a(t)</math></p> <p><math>\frac{dx_i}{dt} = x_i \tau \left( \frac{dQ_i}{dt} - \sum_j \frac{dQ_j}{dt} x_j \right)</math></p> <p><math>\frac{dx_i}{dt} = x_i \tau \alpha \left( r_i - \sum_j x_j r_j + (1/\tau) \sum_j x_j \log(x_j/x_i) \right)</math></p> <p><math>\frac{dx_i}{dt} = x_i \tau \alpha \left( (Ay)_i - x \cdot Ay + (1/\tau) \left[ -\log x_i + \sum_j x_j \log x_j \right] \right)</math></p> <p><math>\frac{dy_i}{dt} = y_i \tau \alpha \left( (Bx)_i - y \cdot Bx + (1/\tau) \left[ -\log y_i + \sum_j y_j \log y_j \right] \right)</math></p>						
No-regret Learning	Dynamics	<p><math>\text{SWAP}_A^t(j, k) = \begin{cases} e_k \cdot Ay^i &amp; \text{if } x^i = e_j \\ x^i \cdot Ay^i &amp; \text{if } x^i \neq e_j \end{cases}</math></p> <p><math>\text{DIFF}_A^t(j, k) = \frac{1}{t} \left( \sum_{i=1}^t [\text{SWAP}_A^i(j, k) - x^i \cdot Ay^i] \right)</math></p> <p><math>\text{REGRET}_A^t(j, k) = \max(\text{DIFF}_A^t(j, k), 0)</math></p> <p><math>p_j^{t+1} = \frac{1}{\mu} \text{REGRET}_A^t(j^*, j)</math> for all <math>j \neq j^*</math></p> <p><math>p_{j^*}^{t+1} = 1 - \sum_{j \neq j^*} p_j^{t+1}</math> when <math>j = j^*</math></p>						
Convergence	<p>(Hart and Mas-Colell). Provided we fix <math>\mu</math> sufficiently large, if player A follows this algorithm then almost surely <math>\text{REGRET}_A^t(j, k) \rightarrow 0</math></p> <p><b>If two player play no-regret learning, the system converges to CE</b></p>							
Blackwell's Approachability	<p>a convex set <math>C \in R^k</math> is approachable for the vector payoff A if for each t and all probabilities <math>\{p^i, q^i\}_{i=1}^{t-1}</math>, there exists a choice <math>p^t</math> so that for each choice of <math>q^t</math> (which player A does not know before choosing <math>p^t</math>), the vectors <math>a_t</math> converge to <math>C</math> as <math>t \rightarrow \infty</math></p> <p>For any closed convex set C the following are equivalent.</p> <ol style="list-style-type: none"> <li>C is approachable for the vector payoff A;</li> <li>for each q there exists p so that <math>A(p, q) \in C</math></li> <li>every half space containing C is approachable.</li> </ol>							
Connections between Learning dynamics	FP vs. RD	<p>The time average of a replicator orbit corresponds to a pseudo-orbit of fictitious play dynamics</p> <p>Exists hyperbolic orbits in FP, then exists a corresponding orbits in RD</p>						
RL vs. FP	<p>Reinforcement learning with choosing -greedy choices is very closely related to the type of dynamics one sees in Best Response dynamics and Fictitious Play.</p>							
RD vs. RL	<p>Softmax Q-learning with <math>\alpha = 1/\tau \rightarrow 0</math>, the dynamics converges to the usual replicator dynamics</p> <p><a href="https://arxiv.org/abs/nlin/0408039">https://arxiv.org/abs/nlin/0408039</a></p>							

