

**SOME NONLINEAR NETWORKS CAPABLE OF LEARNING
A SPATIAL PATTERN OF ARBITRARY COMPLEXITY***

BY STEPHEN GROSSBERG

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

Communicated by C. C. Lin, December 15, 1967

(1) *Introduction.*—This note describes some nonlinear networks which can learn a spatial pattern, in “black and white,” of arbitrary size and complexity. These networks are a special case of a collection of learning machines \mathfrak{N} which were introduced in reference 1, where a machine capable of learning a list of “letters” or “events” was described. We list in heuristic terminology some of the properties which arise in the learning of patterns:

(a) *“Practice makes perfect”*: Given a “black and white” pattern of arbitrary size and complexity, a nonlinear network \mathfrak{N} can be found which learns this pattern to any prescribed degree of accuracy.

(b) *An isolated machine never forgets*: If the pattern is learned to a fixed degree of accuracy by \mathfrak{N} , then \mathfrak{N} will remember the pattern to at least this degree of accuracy until a new pattern is imposed upon \mathfrak{N} .

(c) *Overt practice is unnecessary*: \mathfrak{N} remembers the pattern without practicing it overtly.

(d) *Contour enhancement*: If \mathfrak{N} learns the pattern to a “moderate” degree of accuracy, then \mathfrak{N} ’s memory of the pattern spontaneously improves after practices ceases. As a result, when \mathfrak{N} recalls the pattern, its contours are enhanced in the sense that “darks get darker” and “lights get lighter.”

(e) *A new pattern can always be learned*: Even if \mathfrak{N} knows one pattern to an arbitrary degree of accuracy, this pattern can be replaced by any other pattern by a sufficient amount of practice.

(2) *The Machine.*—The nonlinear network which describes \mathfrak{N} is defined as follows for any fixed number $n \geq 1$ of states and any reaction time $\tau \geq 0$.

$$\dot{x}_i(t) = -\alpha x_i(t) + \beta \sum_{m=1}^n x_m(t - \tau) y_{mi}(t) + I_i(t), \quad (1)$$

$$y_{jk}(t) = z_{jk}(t) [\sum_{m=1}^n z_{jm}(t)]^{-1}, \quad (2) \quad (*)$$

and

$$\dot{z}_{jk}(t) = -u z_{jk}(t) + \beta x_j(t - \tau) x_k(t), \quad (3)$$

where $i, j, k = 1, 2, \dots, n$. (*) describes the following process.

Let G be a graph with vertices $V = \{v_i: i = 1, 2, \dots, n\}$ and directed edges $E = \{e_{jk}: j, k = 1, 2, \dots, n\}$. Each v_i is drawn as a point and e_{jk} is drawn as an arrow facing from v_j to v_k . $x_i(t)$ describes a process going on at v_i , and $y_{jk}(t)$ describes a process going on at the *arrowhead* N_{jk} of e_{jk} . Equation (1) has the following interpretation. At time $t - \tau$, each v_m emits a signal of size $\beta x_m(t - \tau)$ into e_{mi} . This signal travels along e_{mi} at finite velocity until it reaches N_{mi} at time t . The signal thereupon activates the process $y_{mi}(t)$, and a quantity $\beta x_m(t - \tau) y_{mi}(t)$ is instantaneously transmitted from N_{mi} to v_i , and thereby

changes the rate of growth $\dot{x}_i(t)$ of $x_i(t)$. Since this is true for every $m = 1, 2, \dots, n$, the total signal received by v_i from all v_m at time t is $\beta \sum_{m=1}^n x_m(t - \tau) y_{mi}(t)$. $x_i(t)$ also spontaneously decays at the rate $-\alpha x_i(t)$. $I_i(t)$ is the input signal to v_i created by the pattern.

$y_{jk}(t)$ in (2) is the ratio of functions $z_{jm}(t)$ which, as (3) shows, cross-correlate the signal $\beta x_j(t - \tau)$ received by N_{jm} from v_j at time t with the value $x_m(t)$ of the contiguous vertex v_m at time t .

These equations can be derived from simple psychological postulates and have a suggestive neural interpretation.² They are studied mathematically in reference 3, and are extended to more realistic neural equations in reference 4, which, for example, contain the Hartline-Ratliff equation⁵ as a special case. The "contour enhancement" in property (d) above will thereupon be seen as an extension of contour enhancement as it is usually discussed in terms of lateral inhibition.

(3) *Spatial Patterns*.—For purposes of learning a spatial pattern, arrange the vertices v_i in a rectangular grid. Not all inputs $I_i(t)$ in (1) represent spatial patterns. For example, the pattern "A" does not depend on the absolute "blackness" of its lines, but only on their relative blackness as compared to the surround. A *pattern* is therefore defined as an input $I_i(t)$ of the form

$$I_i(t) = \theta_i I(t), \quad i = 1, 2, \dots, n, \tag{4}$$

where the θ_i 's are arbitrary, but fixed, nonnegative numbers whose sum can be taken equal to 1 without loss of generality. The pattern "A" is the same whether or not we view it in steady light or flickering light. $I(t)$ can therefore oscillate quite wildly without changing the pattern described by the θ_i 's. In fact the following theorem holds, which describes the way in which the probabilities $y_{jk}(t) = z_{jk}(t) [\sum_{m=1}^n z_{jm}(t)]^{-1}$ and the correspondingly defined probabilities $X_k(t) = x_k(t) [\sum_{m=1}^n x_m(t)]^{-1}$ learn an arbitrary pattern. Other facts and generalizations concerning this learning process are contained in reference 3.

THEOREM 1. *Suppose $u > 2(\alpha - \beta) > 0$ and $\beta > 0$. Let n be any fixed number of states and let τ be any fixed nonnegative reaction time. Let $I_i(t) = \theta_i I(t)$ be any pattern with $I(t)$ nonnegative, continuous, and bounded, and such that positive constants k and T_0 exist for which*

$$\int_0^t e^{\alpha v} I(v) dv \geq k e^{\alpha t}, \quad t \geq T_0. \tag{5}$$

Then for arbitrary nonnegative and continuous initial data in (), the limits $Q_i = \lim_{t \rightarrow \infty} X_i(t)$ and $P_{jk} = \lim_{t \rightarrow \infty} y_{jk}(t)$ exist, and obey the equations*

$$P_{ji} = Q_i = \theta_i, \quad i, j = 1, 2, \dots, n. \tag{6}$$

Equation (6) says that the probability $X_i(t)$ of v_i and the correlations $y_{ji}(t)$ of all N_{ji} touching v_i learn the relative weight θ_i of the pattern, just so long as the absolute intensity $I(t)$ of the pattern is not "too small" in the sense of (5). $I(t)$ can in fact oscillate very wildly without violating (5). A pattern can therefore be learned to arbitrary accuracy if only it is presented sufficiently often. In

order to learn ever more subtle gradations of shading in the pattern, it suffices to take the number n of vertices in the rectangular grid ever larger.

Equation (5) requires that $I(t)$ take on positive values at arbitrarily large values of t . We now describe what happens if a "truncated" pattern $I_i^{(w)}(t) = \theta_i I^{(w)}(t)$ is presented, where $I^{(w)}(t) = I(t)$, $0 \leq t < w$, and $I^{(w)}(t) = 0$, $t \geq w$. That is, \mathfrak{N} is exposed to the pattern only in the time interval $[0, w]$.

THEOREM 2. *Suppose $u > 2(\alpha - \beta) > 0$ and $\beta > 0$. Let $n \geq 2$ (to avoid trivialities) and $\tau \geq 0$. Let $I_i(t) = 0, t \geq w$, for all $i = 1, 2, \dots, n$. Then for arbitrary nonnegative and continuous data in $[w - \tau, w]$, the limits Q_i and P_{ji} exist and lie in the interval $[m_i(w), M_i(w)]$, where*

$$m_i(w) = \min \{ X_i(w), y_{ki}(w) : k = 1, 2, \dots, n \}$$

and

$$M_i(w) = \max \{ X_i(w), y_{ki}(w) : k = 1, 2, \dots, n \}.$$

Denoting the functions of (*) which are exposed to $I_i^{(w)}(t)$ by superscripts " (w) " (for example, $X_i(t)$ becomes $X_i^{(w)}(t)$), we find the following corollary.

COROLLARY 1.

$$\lim_{w \rightarrow \infty} \lim_{t \rightarrow \infty} X_i^{(w)}(t) = \lim_{w \rightarrow \infty} \lim_{t \rightarrow \infty} y_{ji}^{(w)}(t) = \theta_i, \quad i, j = 1, 2, \dots, n. \tag{7}$$

Proof: By Theorem 1, $\lim_{w \rightarrow \infty} m_i(w) = \lim_{t \rightarrow \infty} M_i(w) = \theta_i$.

These theorems say that if the pattern is exposed to \mathfrak{N} during $[0, w]$ and if w is taken sufficiently large, then \mathfrak{N} will learn the pattern to an arbitrary degree of accuracy and will remember the pattern to at least this degree of accuracy thereafter. \mathfrak{N} does this without "practicing overtly" because the outputs $x_i(t)$ from \mathfrak{N} decay exponentially to 0 for $t \geq w$ whenever $\alpha > \beta > 0$.

Contour enhancement occurs in \mathfrak{N} because of the following corollary, which describes the "envelope"

$$Y_i(t) = \max \{ y_{ki}(t) : k = 1, 2, \dots, n \}$$

and

$$y_i(t) = \min \{ y_{ki}(t) : k = 1, 2, \dots, n \}$$

of correlations whose arrowheads N_{ki} touch v_i .

COROLLARY 2. *For w sufficiently large, one of the following alternatives holds for each $i = 1, 2, \dots, n$:*

(a) $Y_i^{(w)}(t) \geq X_i^{(w)}(t) \geq \theta_i$, $y_i^{(w)}(t) \geq \theta_i$, and $Y_i^{(w)}(t)$ is monotone decreasing for $t \geq w$; or

(b) $\theta_i \geq X_i^{(w)}(t) \geq y_i^{(w)}(t)$, $\theta_i \geq Y_i^{(w)}(t)$, and $y_i^{(w)}(t)$ is monotone increasing for $t \geq w$; or

(c) $Y_i^{(w)}(t) \geq \theta_i \geq y_i^{(w)}(t)$, $Y_i^{(w)}(t) \geq X_i^{(w)}(t) \geq y_i^{(w)}(t)$, $Y_i^{(w)}(t)$ is monotone decreasing, and $y_i^{(w)}(t)$ is monotone increasing for $t \geq w$.

In other words, after a sufficient amount of exposure to the pattern, the envelope of correlations "spontaneously" approaches the pattern probabilities θ_i .

Suppose, for example, that $\theta_i = 0$, which designates a "black" portion of the pattern at state v_i . Then case (a) holds, in which $Y_i^{(w)}(t)$ decreases towards zero. That is, "darks get darker."

To see how \mathfrak{M} recalls a pattern, suppose that a pattern has been practiced over a long time interval $[0, w)$ and that the outputs $x_i(t)$ have decayed nearly to zero in the subsequent interval $[w, W]$. We now show that if even a single speck of light is thereupon shined on the machine at a given vertex (say v_1), then τ time units later the pattern will reappear in all its glory at all the vertices v_i if the reaction time τ is sufficiently large. Since all $x_i(W) \cong 0$, we find by (1) that

$$\dot{x}_1(t) \cong -\alpha x_1(t) + I(t), \quad t \in [W, W + \tau],$$

where $I(t)$ represents the speck of light shined on v_1 . Thus a signal is emitted from v_1 to all points v_i . By Theorem 2, $y_{1i}(t) \cong \theta_i$ for $t \in [W, W + \tau]$, and thus

$$\dot{x}_i(t) \cong -\alpha x_i(t) + \beta x_1(t - \tau)\theta_i. \tag{8}$$

Suppose τ is so large that $x_1(t)$ has a chance to decay back toward zero before it receives the signal which it has created in e_{11} . Then by (8),

$$x_i(t) \cong \beta\theta_i e^{-\alpha t} \int_{W+\tau}^t e^{\alpha v} x_1(v - \tau) dv,$$

for all $i = 1, 2, \dots, n$, and in particular

$$\frac{x_i(t)}{x_j(t)} \cong \frac{\theta_i}{\theta_j}, \quad i, j = 1, 2, \dots, n.$$

The outputs $x_i(t)$ therefore reproduce the relative shadings θ_i of the original pattern.

The very act of recalling the pattern helps to destroy \mathfrak{M} 's memory of it, because the speck of light is itself a pattern of the form $I_i(t) = \theta_i I(t)$ with $\theta_1 = 1$ and all $\theta_j = 0, j \neq 1$. This is not true of all the machines introduced in reference 1. An outstar can, for example, recall as many times as it pleases without destroying its memory. An obvious artifice for preserving \mathfrak{M} 's memory under recall is to postulate that every output from \mathfrak{M} creates a proportional "feedback input" which is returned to \mathfrak{M} through the external medium surrounding \mathfrak{M} , much as we "hear ourselves talk." The outstar does not require this artifice because it has a source vertex v_1 which never receives memory-destroying non-linear feedback from other vertices when it is perturbed by a "speck of light." Although as a graph the pattern-learning machine of this note can be viewed as n outstars connected together, the dynamics of this machine is not simply the sum of the dynamics of connected outstars. It would be highly desirable to be able to recapture the stability of an outstar's memory even in a graph whose vertices are interconnected. Reference 4 shows that (*) must be altered to include lateral inhibitory signals and thresholds to achieve this effect. In other words, lateral inhibition and thresholds "localize" the dynamics of the graph.

Even if \mathfrak{M} knows a given pattern perfectly at time $t = T$, it can relearn any other pattern thereafter. This is because the values of (*)'s variables in the

interval $[T - \tau, T]$ can be viewed as the initial data for (*) in the interval (T, ∞) . Since these values are nonnegative and continuous, and Theorems 1 and 2 hold for all nonnegative and continuous initial data, our contention is proved.

* The preparation of this work was supported in part by the National Science Foundation (NSF GP-7477).

¹ Grossberg, S., "Nonlinear difference-differential equations in prediction and learning theory," these PROCEEDINGS, **58**, 1329 (1967).

² Grossberg, S., "Embedding fields: A new theory of learning with physiological implications," *J. Math. Psych.*, to appear.

³ Grossberg, S., "A prediction theory for some nonlinear functional-differential equations, II. Learning of patterns," *J. Math. Anal. Appl.*, to appear.

⁴ Grossberg, S., "On learning, information, lateral inhibition, and transmitters."

⁵ Ratliff, F., *Mach Bands: Quantitative Studies on Neural Networks in the Retina* (New York: Holden-Day, 1965).