



**СООБЩЕНИЯ  
ОБЪЕДИНЕННОГО  
ИНСТИТУТА  
ЯДЕРНЫХ  
ИССЛЕДОВАНИЙ  
ДУБНА**

P17-87-145

**В.А.Загребнов, Я.Л. ван Хеммен\***

**ЗАПОМИНАНИЕ  
ЭКСТЕНСИВНО БОЛЬШОГО ЧИСЛА  
ВЗВЕШЕННЫХ ОБРАЗОВ  
НАСЫЩЕННЫМИ НЕЙРОННЫМИ СЕТЯМИ**

---

\* Университет, Гейдельберг, ФРГ

## I. Введение

Формально нейронная сеть – это множество связанных между собой синапсами нейронов, находящихся в одном из двух возможных состояний, которое имеет такую динамику, что запоминаемые образы являются для неё устойчивыми аттракторами. Одной из главных проблем теории нейронных сетей является изучение возможных способов запоминания или информации и описание механизма, который позволяет как находить (вспоминать), так и забывать хранящуюся в них информацию.

Обычно в качестве вполне реального приближения рассматривают сеть из  $N$  полностью связанных нейронов. Согласно работе [1] нейрон может находиться только в одном из двух состояний (возбужденном и невозбужденном) и, следовательно, может быть ассоциирован с изинговским спином  $S(i)$ ,  $1 \leq i \leq N$ , причем  $S = +1$  соответствует возбуждению, а  $S = -1$  соответствует покою. На этом языке образ есть некоторая конфигурация изинговских спинов.

Общепринято считать [2-7], что все существенные особенности временного поведения сети описываются динамикой Монте – Карло, которая определяется гамильтонианом

$$H_N = -\frac{1}{2} \sum_{i \neq j}^N J_{ij} S(i) S(j) \quad (I.I)$$

Это означает, что динамика сети сводится к "движению вниз с холмов ландшафта", соответствующего (свободной) энергии для  $H_N$ , а асимптотическая устойчивость предела определяется в рамках равновесной статистической механики, задаваемой (I.I). Ниже критической температуры  $T_c$  эргодичность в системе (I.I) нарушается и хранящиеся образы соответствуют притягивающим множествам (равновесные состояния или эргодические компоненты) в фазовом пространстве изинговского спинового стекла, которое мы рассмотрим ниже.

Следуя Хеббу [8], будем считать, что память сосредоточена в синапсах, или точнее, – в распределении величин электрических потенциалов синапсов, которые определяют значения констант обменного взаимодействия в гамильтониане Изинга (I.I). При соответствующем выборе взаимодействий сеть работает как нечувствительная к отказам ассоциативная память, см., например, [2-6]. Дополнительные образы можно запомнить, если соответствующим образом модифицировать  $J_{ij}$ . Для облегчения моделирования обычно предполагают, что образы  $\{E_{i\alpha} : 1 \leq i \leq N\}$  для  $1 \leq \alpha \leq q$  являются "случайными". Это означает,

что  $\varepsilon_{i\alpha} = \pm 1$  являются независимыми, одинаково распределенными случайными величинами, которые принимают значения  $\pm 1$  с равной вероятностью, а заданный образ является некоторой реализацией этой конечной последовательности.

Используемая для описания ассоциативной памяти модель Хопфилда /2/ соответствует следующему выбору обменного взаимодействия:

$$J_{ij} = \frac{1}{N} \sum_{\alpha=1}^q \varepsilon_{i\alpha} \varepsilon_{j\alpha} \quad (1.2)$$

В недавней работе /9/ Амит, Гутфренд и Сомполинский изучали эту сеть вблизи насыщения, т.е. когда  $q = \alpha N$  и  $\alpha > 0$ . С помощью довольно изобретательного анализа в рамках теории среднего поля они показали, что при нулевой температуре ( $T = 0$ ) система способна эффективно восстанавливать информацию, если  $\alpha < \alpha_c = 0,14$ . Однако при  $\alpha = \alpha_c$  количество восстанавливаемых образов (скачком) сокращается, и при  $\alpha > \alpha_c$  эффективное восстановление информации невозможно. Следовательно, система ведет себя таким образом, как если бы при  $\alpha = \alpha_c$  происходил фазовый переход первого рода, после которого вся информация более или менее забывается. С точки зрения физиологии такая модель не внушает доверия. Поэтому предлагались и другие схемы работы памяти (как в режиме запоминания, так и в режиме воспоминания), но до сих пор их статус остается очень удивительным. Например, Паризи /10/ предложил модель памяти, "которая забывает", однако память, работающая по этой схеме, забывает все, за исключением самых последних образов.

В настоящей работе содержатся следующие результаты. Во-первых, предложен более простой и ясный путь, позволяющий получить основные результаты работы /9/, специальное внимание уделено анализу математического статуса используемого формализма. Во-вторых, предложено обобщение модели Хопфилда, такое, что число хранящихся образов в системе размера  $N$  также может быть равно  $N$  (случай полного насыщения), причем каждый образ имеет вес  $\varepsilon_{\nu}$ . Если взвешивание образов соответствует временному порядку их поступления в память и  $\varepsilon_{\nu} \rightarrow 0$  при  $\nu \rightarrow \infty$ , что в этом случае естественно, то при соответствующем выборе весов система запоминает экстенсивно большое число образов, однако с ростом  $\nu$  их четкость постепенно уменьшается.

В разделе 2 дано описание модели и проведен её анализ в рамках теории среднего поля. В следующем разделе метод симметричной реплики использован для вычисления свободной энергии. Соответствующая этой модели ёмкость памяти при  $T = 0$  найдена в разделе 4. Обсуждению полученных результатов посвящен заключительный раздел.

## 2. Модель Хопфилда для взвешенных образов

Мы будем рассматривать гамильтониан (1.1), в котором константы обменного взаимодействия имеют вид

$$J_{ij} = \frac{1}{N} \sum_{n=1}^N \varepsilon_n \varepsilon_{in} \varepsilon_{jn} \quad (2.1)$$

каждый образ имеет вес  $\varepsilon_n$ . Эти веса произвольны, за исключением общего требования:  $0 \leq \varepsilon_n \leq 1$ . Специальный выбор весов, соответствующий модели Хопфилда, следующий:  $\varepsilon_n = 1$  для  $1 \leq n \leq \alpha N$  и  $\varepsilon_n = 0$  для  $n > \alpha N$ .

Так же, как и в работе /9/, мы вначале используем метод реплики и определим устойчивость некоторого конечного числа образов, которые мы обозначим  $\nu$ , интегрируя по оставшимся, которые мы обозначим  $\mu$ . Для этой и других целей удобно ввести в гамильтониан внешнее поле, которое выделяет  $\nu$ -образы:

$$H_{ext} = - \sum_{\nu} h_{\nu} \sum_{i=1}^N \varepsilon_{i\nu} S(i) \quad (2.2)$$

В рамках метода реплики начинают обычно /11/ с вычисления функции

$$\phi_N(n) = \frac{1}{N} \ln \langle Z_N^n \rangle \quad (2.3)$$

для положительных целых  $n$ , затем переходят к термодинамическому пределу  $N \rightarrow \infty$ , который равен  $\phi(n)$ , и, наконец, строят продолжение функции  $\phi(n)$  (обычно для симметричной реплики) в окрестность  $n = 0$ . Тогда производная  $\phi'(0)$  даёт нам  $-\beta f(\beta)$ , где  $f(\beta)$  — плотность свободной энергии для обратной температуры  $\beta$ , а  $Z_N = \text{Tr} \exp(-\beta H_N)$  — статистическая сумма по всем изинговским конфигурациям. Угловые скобки в (2.3) обозначают усреднение по всем возможным образам  $\{\varepsilon_{in}\}$ . Поскольку мы вначале интегрируем по  $\mu$ -образам, то вместо  $\langle Z_N^n \rangle$  сосредоточимся на вычислении

$$\left\langle \exp \left\{ \frac{\beta}{2N} \sum_{\mu, \rho} \varepsilon_{\mu} \left( \sum_{i=1}^N \varepsilon_{i\mu} S_{\rho}(i) \right)^2 - \frac{1}{2} \beta n \sum_{\mu} \varepsilon_{\mu} \right\} \right\rangle \quad (2.4)$$

Здесь  $\rho : 1 \leq \rho \leq n$  нумерует  $n$  реплик. Далее, вплоть до формулы (3.5), мы будем всюду опускать член  $(-\frac{1}{2} \beta n \sum_{\mu} \varepsilon_{\mu})$ . С помощью соотношения

$$\exp \left( \frac{1}{2} \lambda a^2 \right) = \int_{-\infty}^{+\infty} \frac{dz}{\sqrt{2\pi}} \exp \left( -\frac{1}{2} z^2 + \sqrt{\lambda} a z \right) \quad (2.5)$$

мы линеаризуем квадрат в экспоненте (2.4):

$$\left\langle \prod_{\mu, \rho} \frac{dm_{\mu\rho}}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \sum_{\mu, \rho} m_{\mu\rho}^2 + \sum_{\mu, \rho} m_{\mu\rho} \left[ \left( \frac{\beta \varepsilon_{\mu}}{N} \right)^{1/2} \sum_{i=1}^N \xi_{i\mu} S_{\rho}(i) \right] \right\} \right\rangle, \quad (2.6)$$

и преобразуем среднее по  $\{\xi_{i\mu}\}$  к виду:

$$\int \prod_{\mu, \rho} \frac{dm_{\mu\rho}}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \sum_{\mu, \rho} m_{\mu\rho}^2 + \sum_{i, \mu} \ln \left( \text{ch} \left[ \left( \frac{\beta \varepsilon_{\mu}}{N} \right)^{1/2} \sum_{\rho=1}^n m_{\mu\rho} S_{\rho}(i) \right] \right) \right\}. \quad (2.7)$$

Если предположить, что член в фигурных скобках в (2.7) "мал" и заменить  $\ln \text{ch } x$  для малых  $x$  на  $\frac{1}{2} x^2$ , то получим

$$\int \prod_{\mu, \rho} \frac{dm_{\mu\rho}}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \sum_{\mu, \rho} m_{\mu\rho}^2 + \sum_{i, \mu} \sum_{\rho, \sigma} \frac{\beta \varepsilon_{\mu}}{2N} m_{\mu\rho} m_{\mu\sigma} S_{\rho}(i) S_{\sigma}(i) \right\}, \quad (2.8)$$

где  $1 \leq \rho, \sigma \leq n$ . Для каждого фиксированного  $\mu$  интегрирование по  $m_{\mu\rho}$  может быть выполнено точно, и мы получаем

$$\det(Q_{\mu})^{-1/2}, \quad (2.9)$$

где  $Q_{\mu}$  - симметричная  $n \times n$  матрица, элементы которой имеют вид

$$(Q_{\mu})_{\rho\sigma} = \delta_{\rho\sigma} - \beta \varepsilon_{\mu} \frac{1}{N} \sum_{i=1}^N S_{\rho}(i) S_{\sigma}(i), \quad (2.10)$$

здесь  $\delta_{\rho\sigma}$  - символ Кронекера. Таким образом, выражение (2.8) принимает вид

$$\prod_{\mu} \det(Q_{\mu})^{-1/2} = \exp \left\{ -\frac{1}{2} \sum_{\mu} \text{Tr}(\ln Q_{\mu}) \right\}. \quad (2.11)$$

Сделаем теперь следующее отступление /12/. Для того чтобы (2.9) имело смысл, матрица  $Q_{\mu}$  должна быть положительно определенной, т.е. её диагональные элементы должны быть положительными. Однако  $(Q_{\mu})_{\rho\rho} = 1 - \beta \varepsilon_{\mu}$  отрицательно для достаточно больших  $\beta$ . Причину этого легко усмотреть, если вернуться к переходу от (2.7) к (2.8). В то время, как выражение (2.7) хорошо определено для всех  $\beta$  (поскольку  $\ln \text{ch } x \sim |x|$  при  $|x| \rightarrow \infty$ ), выражение (2.8) таковым не является. Положение выглядит безвыходным, так как мы действительно нуждаемся в квадратичной аппроксимации  $\ln \text{ch } x$  - в противном случае невозможны аналитические вычисления. Ниже мы покажем, однако, что положительная определенность восстанавливается в пределе  $n \rightarrow 0$ .

Зависящее от спиновых переменных выражение (2.11) представляет собой шумы, которые вносят  $\mu$ -образы. Вклад этих шумов необходимо

учесть в гамильтониане для  $\nu$ -образов. С точностью до оговорок, сделанных выше, это можно сделать строго. Из соотношений (2.3) и (2.11) получаем

$$\phi_N(n) = \frac{1}{N} \ln \left\langle \text{Tr}_{S_{\nu}} \exp \left[ N \left\{ \frac{1}{2} \beta \sum_{\nu, \rho} \varepsilon_{\nu} \left( \frac{1}{N} \sum_{i=1}^N \xi_{i\nu} S_{\rho}(i) \right)^2 + \beta \sum_{\nu, \rho} h_{\nu} \left( \frac{1}{N} \sum_{i=1}^N \xi_{i\nu} S_{\rho}(i) \right) - \frac{1}{2N} \sum_{\mu} \text{Tr}(\ln Q_{\mu}) \right\} \right] \right\rangle. \quad (2.12)$$

Первый след является суммой по всем  $2^{nN}$  изинговским спиновым конфигурациям  $n$  реплик. Определим два типа параметров порядка:

$$m_{\nu\rho} = \frac{1}{N} \sum_{i=1}^N \xi_{i\nu} S_{\rho}(i), \quad 1 \leq \rho \leq n; \quad (2.13)$$

$$q_{\rho\sigma} = \frac{1}{N} \sum_{i=1}^N S_{\rho}(i) S_{\sigma}(i), \quad 1 \leq \rho < \sigma \leq n.$$

Тогда выражение в фигурных скобках в (2.12) можно представить в виде

$$F(\underline{m}, \underline{q}) = \beta \sum_{\nu, \rho} \left( \frac{1}{2} \varepsilon_{\nu} m_{\nu\rho}^2 + h_{\nu} m_{\nu\rho} \right) - \frac{1}{2N} \sum_{\mu} \text{Tr} \ln Q_{\mu}(\underline{q}), \quad (2.14)$$

откуда следует, что более удобными переменными являются не спины  $S_{\rho}(i)$ ,  $1 \leq i \leq N$  и  $1 \leq \rho \leq n$ , а новые переменные  $m_{\nu\rho}$  и  $q_{\rho\sigma}$ . Чтобы сделать этот переход, нам необходимо знать соответствующий "якобиан"  $\mathcal{D}_N(\underline{m}, \underline{q})$ . Как показано в работах /13-15/, в пределе  $N \rightarrow \infty$  такое преобразование переменных возможно, причем с вероятностью единица (по отношению к априорной вероятностной мере) имеем

$$\mathcal{D}_N(\underline{m}, \underline{q}) = \exp \left\{ -N c^*(\underline{m}, \underline{q}) \right\}, \quad (2.15)$$

где

$$c^*(\underline{m}, \underline{q}) = \sup_{\underline{x}, \underline{y}} \{ \underline{m} \cdot \underline{x} + \underline{q} \cdot \underline{y} - c(\underline{x}, \underline{y}) \} \quad (2.16)$$

является преобразованием Лежандра /16/ (строго) выпуклой функции

$$c(\underline{x}, \underline{y}) = \left\langle \ln \text{Tr} \exp \left\{ \sum_{\nu, \rho} x_{\nu} \xi_{\nu} S_{\rho} + \sum_{(\rho, \sigma)} y_{\rho\sigma} S_{\rho} S_{\sigma} \right\} \right\rangle. \quad (2.17)$$

Второе суммирование в (2.17) осуществляется только по парам  $(\rho, \sigma)$ . След относится к  $n$  изинговским спином  $S_{\rho}$ ,  $1 \leq \rho \leq n$ , а в усреднении, обозначенном угловыми скобками, каждая переменная  $\xi_{\nu}$  встречается только один раз, причем их число конечно. Кроме того, уже в формулировке (2.15) неявно содержится утверждение о том /15/, что при  $N \rightarrow \infty$  правая часть (2.15) не зависит от случайной конфигурации

$\underline{\mu}$  - образам. По этой причине угловые скобки в выражении (2.12) можно опустить<sup>1)</sup>.

Если обозначить пару переменных  $(\underline{m}, \underline{q})$  через  $\underline{\mu}$ , то из (2.12)-(2.17) получаем выражение

$$\phi(n) = \lim_{N \rightarrow \infty} \frac{1}{N} \ln \int d\underline{\mu} \exp\{N[F(\underline{\mu}) - c^*(\underline{\mu})]\}, \quad (2.18)$$

которое с помощью стандартных аргументов (Лапласа) можно представить в виде

$$\phi(n) = \sup_{\underline{\mu}} \{F(\underline{\mu}) - c^*(\underline{\mu})\}. \quad (2.19)$$

Для достаточно малых  $\beta$  это выражение является точным, однако оно становится совершенно формальным, если  $\beta \max \epsilon_{\mu} > 1$ .

Если мы, несмотря на это, продолжим выкладки, то получим (см. приложение), что

$$\phi(n) = \max_{\underline{\mu}} \{F(\underline{\mu}) - \underline{\mu} \cdot \nabla F(\underline{\mu}) + c(\nabla F(\underline{\mu}))\}, \quad (2.20)$$

где  $\underline{\mu}$  удовлетворяет уравнению самосогласования (уравнение для неподвижной точки) вида

$$\underline{\mu} = \nabla c(\nabla F(\underline{\mu})). \quad (2.21)$$

Соотношение (2.20) имеет то же смысл, что и (2.19).

Матричные элементы  $(Q_{\underline{\mu}}(q))_{\rho\sigma}$  для  $\rho < \sigma$  имеют вид

$$(Q_{\underline{\mu}})_{\rho\sigma} = \delta_{\rho\sigma} (1 - \beta \epsilon_{\mu}) - \beta \epsilon_{\mu} q_{\rho\sigma} = (Q_{\underline{\mu}})_{\sigma\rho}. \quad (2.22)$$

Используя соотношение

$$\frac{\partial}{\partial Q_{\rho\sigma}} \tau \ln Q = 2 (Q^{-1})_{\rho\sigma}, \quad (2.23)$$

легко проверить, что

<sup>1)</sup> Такая возможность не является следствием свойства самоусреднения плотности свободной энергии, как это утверждали авторы работы [9]. Заметим, что усреднение необходимо провести под знаком логарифма. Отметим также, что параметры порядка  $z_{\rho\sigma}$ , введенные этими авторами, не имеют физического смысла, без них можно обойтись.

$$\nabla F(\underline{\mu}) = \left( \begin{array}{c} \beta (\epsilon_{\nu} m_{\nu\rho} + h_{\nu}) \\ \frac{1}{N} \sum_{\mu} \beta \epsilon_{\mu} (Q_{\mu}^{-1})_{\rho\sigma} \end{array} \right), \quad \underline{\mu} = (\underline{m}, \underline{q}). \quad (2.24)$$

Теперь мы можем продолжить функцию (2.20) на  $n = 0$ , предполагая, что реплика симметрична.

### 3. Симметричная реплика

Это означает, что реплики равны:  $m_{\nu\rho} = m_{\nu}$  и  $q_{\rho\sigma} = q$  ( $\rho \neq \sigma$ ). Такое требование совместимо с уравнением (2.21) для неподвижной точки. Далее, представим матрицу  $Q_{\underline{\mu}}(q)$  в виде

$$Q_{\underline{\mu}}(q) = (1 - \beta + \beta \epsilon_{\mu} q) \mathbb{1} - \beta \epsilon_{\mu} q n \frac{1}{\sqrt{n}} \mathbb{1} \langle \frac{1}{\sqrt{n}} \mathbb{1} | \equiv a(n) - n b(n) P, \quad (3.1)$$

где  $\mathbb{1}$  - единичная матрица,  $\frac{1}{\sqrt{n}} \mathbb{1}$  - единичный вектор  $(1, \dots, 1) \in \mathbb{R}^n$ , а  $P = P^2$  - проекционный оператор. Поскольку обратную матрицу можно представить в виде  $Q^{-1} = c - d \cdot P$ , то её матричные элементы легко найти, и

$$(Q^{-1})_{\rho\sigma} = -b(n) [a(n)(n b(n) - a(n))]^{-1}. \quad (3.2)$$

Следовательно,

$$\frac{1}{N} \sum_{\mu} \beta \epsilon_{\mu} (Q_{\mu}^{-1})_{\rho\sigma} \equiv \beta^2 q \tau(n), \quad (3.3)$$

где  $\tau(n)$  таково, что явно не зависит от  $\beta$ . Тогда из (2.24) получаем

$$\nabla F(\underline{\mu}) = \left( \begin{array}{c} \beta (\epsilon_{\nu} m_{\nu} + h_{\nu}) \\ \beta^2 q \tau(n) \end{array} \right). \quad (3.4)$$

Используя (3.1), можно непосредственно убедиться, что невырожденные собственные значения матрицы  $Q_{\underline{\mu}}(q)$  равны  $1 - \beta \epsilon_{\mu} (1 - q) + \beta \epsilon_{\mu} q n$ , а  $(n-1)$  -кратно вырожденные - равны  $1 - \beta \epsilon_{\mu} (1 - q)$ . В пределе  $n \rightarrow 0$  остается только величина  $1 - \beta \epsilon_{\mu} (1 - q)$ , которая должна быть положительна, см. (3.7) ниже. В следующем разделе мы убедимся, что этот предел действительно положителен, см. также [9].

С помощью (2.17), (2.20) и (3.1)-(3.4), а кроме того, добавляя константу, которую мы опускали, начиная с формулы (2.4), получаем

$$\phi(n) = -\frac{1}{2} \beta n \left( \frac{1}{N} \sum_{\mu} \epsilon_{\mu} \right) - \frac{1}{2} \beta n \sum_{\nu} m_{\nu}^2 - \frac{1}{2N} \sum_{\mu} [\ln(1 - \beta \epsilon_{\mu} (1 - q)) + \beta \epsilon_{\mu} q n] + (n-1) \times$$

$$\times \ln(1 - \beta \epsilon_{\mu} (1 - q)) - \frac{1}{2} n (n-1) \tau(n) (\beta q)^2 + \langle \ln \tau_{\rho\sigma} \exp\{\beta \sum_{\nu, \rho} (\epsilon_{\nu} m_{\nu} + h_{\nu}) \epsilon_{\nu} S_{\rho} + \frac{1}{2} \beta^2 \sum_{\rho \neq \sigma} q \tau(n) S_{\rho} S_{\sigma}\} \rangle. \quad (3.5)$$

Если воспользоваться представлением (2.5) и провести линеаризацию, то след можно вычислить и представить последний член в (3.5) в виде

$$-\frac{1}{2} n \beta^2 q \tau (n) + n \ln 2 + \left\langle \ln \int_{-\infty}^{+\infty} \frac{d\tilde{x}}{\sqrt{2\pi}} e^{-\frac{1}{2} \tilde{x}^2} \text{ch}^n \left\{ \beta \left[ (\varepsilon \underline{m} + \underline{h}) \cdot \underline{\xi} + \tilde{x} \sqrt{q \tau (n)} \right] \right\} \right\rangle, \quad (3.6)$$

где  $\varepsilon = \text{diag}(\varepsilon_\nu)$  - диагональная матрица. Если выбрать теперь в качестве продолжения  $\phi(n)$  её "очевидное" продолжение как вещественной функции, то получим

$$-\beta f(\beta) = \lim_{n \rightarrow 0} n^{-1} \phi(n) = -\frac{1}{2N} \sum_{\mu} \varepsilon_{\mu} - \frac{1}{2} \beta \sum_{\nu} \varepsilon_{\nu} m_{\nu}^2 - \frac{1}{2N} \sum_{\mu} \left[ \ln (1 - \beta \varepsilon_{\mu} (1-q)) - \beta \varepsilon_{\mu} q (1 - \beta \varepsilon_{\mu} (1-q))^{-1} \right] - \frac{1}{2} \beta^2 q (1-q) \tau + \left\langle \int_{-\infty}^{+\infty} \frac{d\tilde{x}}{\sqrt{2\pi}} e^{-\tilde{x}^2/2} \ln \left[ 2 \text{ch} \left\{ \beta \left[ (\varepsilon \underline{m} + \underline{h}) \cdot \underline{\xi} + \tilde{x} \sqrt{q \tau} \right] \right\} \right] \right\rangle, \quad (3.7)$$

здесь подразумевается, что  $N$  велико и

$$\tau = \lim_{n \rightarrow 0} \tau(n) = \frac{1}{N} \sum_{\mu} \varepsilon_{\mu}^2 \left[ 1 - \beta \varepsilon_{\mu} (1-q) \right]^{-2} \geq 0. \quad (3.8)$$

Далее, необходимо выбрать такое решение уравнений для неподвижной точки

$$\underline{m} = \left\langle \underline{\xi} \text{th} \left\{ \beta \left[ \varepsilon \underline{m} + \underline{h} \right] \cdot \underline{\xi} + \tilde{x} \sqrt{q \tau} \right\} \right\rangle, \quad (3.9a)$$

$$q = \left\langle \text{th}^2 \left\{ \beta \left[ \varepsilon \underline{m} + \underline{h} \right] \cdot \underline{\xi} + \tilde{x} \sqrt{q \tau} \right\} \right\rangle, \quad (3.9b)$$

которое обеспечивает максимум правой части (3.7). Двойные угловые скобки в (3.9) обозначают среднее и по конечному числу переменных  $\{\xi_{\nu}\}$ , и по гауссовскому распределению для переменной  $\tilde{x}$ . Если  $\varepsilon_{\mu} = 1$  для  $1 \leq \mu \leq \alpha N$  и  $\varepsilon_{\mu} = 0$  для  $\mu > \alpha N$ , тогда (3.7)-(3.9) воспроизводят результаты работы /9/ (если всюду положить  $\underline{h} = 0$ ). Из (3.8) и (3.9) следует, что  $\tau$  необходимо интерпретировать как константу, перенормирующую параметр порядка  $q$ . Это следствие шума, который порождает хвост из "бесконечного числа" ( $N \rightarrow \infty$ ) образов.

#### 4. Ёмкость памяти

Для заданного набора весовых факторов  $\varepsilon_{\mu}$  ёмкостью памяти называется (максимальное) число образов, которые не теряют своей стабильности полностью. В этом разделе мы исследуем ёмкость памяти при  $T = 0$  для сети с весовыми факторами  $\varepsilon_{\mu} = \mu^{-\alpha}$ ,  $\alpha > 0$ . При соответствующем выборе  $\alpha$  имеется постепенная потеря чётности об-

разов по мере роста индекса  $\mu$ . Легко проверить, что для  $\alpha > 0$  среднее  $N^{-1} \sum_{\mu} \mu^{-\alpha}$  сходится (при  $N \rightarrow \infty$ ) к нулю (сходимость  $\mu^{-\alpha}$  по Чезаро). Поэтому член  $N^{-1} \sum_{\mu} T \tau (\ln Q_{\mu})$ , соответствующий фоновому шуму, также стремится (при  $N \rightarrow \infty$ ) к нулю. Это не означает, однако, что им можно полностью пренебречь, поскольку вес  $\varepsilon_{\mu}$  образа с номером  $\mu$  также стремится к нулю при  $\mu \rightarrow \infty$ . Имеется, однако, особенность, связанная с предельными переходами, которые мы хотим изучать.

Мы будем рассматривать два предела  $N \rightarrow \infty$  и  $\beta \rightarrow \infty$ , которые необходимо совершать в определенном порядке. Вначале мы должны совершить переход к термодинамическому пределу  $N \rightarrow \infty$ , а затем к пределу нулевой температуры  $\beta \rightarrow \infty$ . Поскольку параметр  $q$  при этом будет стремиться к единице, мы начнем с уравнения (3.9a). В пределе  $\beta \rightarrow \infty$  функция  $\text{th}\{\dots\}$  стремится к  $\text{sgn}\{\dots\}$ , а интеграл по гауссовскому распределению совпадает, с точностью до члена порядка  $T$ , с интегралом ошибок:

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} d\tilde{x} e^{-\tilde{x}^2/2} \text{th}\{\dots\} \approx \sqrt{\frac{2}{\pi}} \int_0^{\varepsilon \underline{m} \cdot \underline{\xi} / \sqrt{q \tau}} d\tilde{x} e^{-\tilde{x}^2/2} = \text{erf}(\varepsilon \underline{m} \cdot \underline{\xi} / \sqrt{q \tau}), \quad (4.1)$$

где

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x dy e^{-y^2}, \quad (4.2)$$

так что уравнение (3.9a) принимает вид

$$\underline{m} = \left\langle \underline{\xi} \text{erf}(\varepsilon \underline{m} \cdot \underline{\xi} / \sqrt{q \tau}) \right\rangle_{\underline{\xi}}. \quad (4.3)$$

Рассмотрим поведение одного выделенного образа с номером  $\nu$ . Тогда  $\underline{m}$  имеет только одну ненулевую компоненту  $m$  и

$$m = \left\langle \xi_{\nu} \text{erf}(m \varepsilon_{\nu} \xi_{\nu} / \sqrt{q \tau}) \right\rangle = \text{erf}(m \varepsilon_{\nu} / \sqrt{q \tau}). \quad (4.4)$$

Как следует из (2.13), чем ближе  $m$  к единице, тем лучше восстановление образа. Поскольку веса  $\varepsilon_{\nu}$  стремятся к нулю в смысле Чезаро, параметр  $\tau$  сходится к нулю при  $N \rightarrow \infty$ , а величина  $m$  (для фиксированного  $\nu$ ) стремится к единице. Аналогично для конечной группы фиксированных индексов  $\{\nu\}$  и конечной температуры получаем

$$\underline{m} = \left\langle \underline{\xi} \text{th} \left\{ \beta \varepsilon \underline{m} \cdot \underline{\xi} \right\} \right\rangle_{\underline{\xi}}, \quad (4.5)$$

а в случае единственной ненулевой компоненты

$$m = th \beta \varepsilon, m. \quad (4.6)$$

Следовательно, чем больше  $\nu$ , тем меньше критическая температура  $T_c(\nu)$ , ниже которой появляется ветвь ненулевых значений  $m$ , соответствующих образу с индексом  $\nu$  (компонента в направлении образа  $\nu$ ).

Обратимся теперь к уравнению (3.9б). В термодинамическом пределе оно сводится к равенству:

$$q = \langle th^2 \{ \beta \varepsilon m \cdot \underline{\underline{\varepsilon}} \} \rangle_{\underline{\underline{\varepsilon}}}, \quad (4.7)$$

тогда при  $\beta \rightarrow \infty$  правая часть сходится к единице экспоненциально быстро, и  $C(\beta) = \beta(1-q)$  стремится к нулю. Как и выше,  $\varepsilon = \text{diag}(\varepsilon_\nu)$ .

Что произойдет для очень больших, но конечных  $N$ ? Какое количество образов система может хранить при  $T = 0$  в этом случае? Фактор шума  $\tau$  в этом случае отличен от нуля. Его можно оценить, используя (3.8) и рассуждения, приведенные выше:

$$\tau = N^{-1} \sum_{\mu} [\varepsilon_{\mu}^{-1} - C(\beta)]^{-2} \xrightarrow{\beta \rightarrow \infty} N^{-1} \sum_{\mu} \mu^{-2\alpha}. \quad (4.8)$$

Это выражение, в свою очередь, можно оценить следующим образом:

$$\tau = N^{-1} \int_1^N d\mu \mu^{-2\alpha} = (1-2\alpha)^{-1} [N^{-2\alpha} - N^{-1}], \quad (4.9)$$

если  $\alpha \neq 1/2$ , и, соответственно,

$$\tau = N^{-1} \int_1^N d\mu \mu^{-1} = N^{-1} \ln N, \quad (4.10)$$

если  $\alpha = 1/2$ . Мы будем различать три случая: (а)  $0 < \alpha < 1/2$ ; (б)  $\alpha = 1/2$ ; (в)  $\alpha > 1/2$ . В первом случае из (4.9) получаем

$$\tau = (1-2\alpha)^{-1} N^{-2\alpha}. \quad (4.11a)$$

Во втором случае из (4.10) следует

$$\tau = N^{-1} \ln N, \quad (4.11б)$$

в то время как в последнем случае имеем

$$\tau = (2\alpha - 1)^{-1} N^{-1}. \quad (4.11в)$$

Теперь эти оценки для  $\tau$  позволяют нам вернуться к выражению (4.4).

Функция  $\psi(x) = \text{erf}(\eta x)$  является выпуклой ( $\eta > 0$ ) и монотонно возрастающей по  $x \geq 0$ . Она начинается в нуле и стремится к единице. Её изменения напоминают поведение  $th(\eta x)$ , за исключением наклона в нуле, который равен  $2\eta/\sqrt{\pi}$ , см. (4.2). Следовательно, уравнение

$$m = \text{erf}(\eta m) \quad (4.12)$$

имеет единственное решение  $m = 0$ , если  $2\eta/\sqrt{\pi} < 1$ , т.е.  $\eta < \eta_c = \sqrt{\pi}/2$ . Для  $\eta > \eta_c$  существует, кроме того, и нетривиальное решение, которое непрерывно стремится к нулю, когда  $\eta \rightarrow \eta_c$  сверху.

Из (4.4) и (4.12) следует, что существует такое критическое значение  $\nu_c$ , что образы с  $\nu > \nu_c$  не запоминаются. Имеем

$$\frac{\varepsilon_\nu}{\sqrt{2\tau}} \Big|_{\nu=\nu_c} = \frac{1}{2} \sqrt{\pi} \iff \sqrt{\tau} (\nu_c)^\alpha = \sqrt{\frac{2}{\pi}}. \quad (4.13)$$

Тогда в случае (а) получаем

$$(1-2\alpha)^{-1/2} \frac{(\nu_c)^\alpha}{N} = \sqrt{\frac{2}{\pi}}, \quad (4.14)$$

так что величина  $\nu_c$  является экстенсивной:

$$\nu_c = \left[ (1-2\alpha) \frac{2}{\pi} \right]^{1/2\alpha} N. \quad (4.15a)$$

В случае (б) имеем

$$\nu_c = \frac{2}{\pi} N / \ln N, \quad (4.15б)$$

в то время, как в случае (в), когда  $2\alpha > 1$ , эта величина равна

$$\nu_c = \left[ (2\alpha - 1) \frac{2}{\pi} \right]^{1/2\alpha} N^{1/2\alpha}. \quad (4.15в)$$

Следовательно, в последних двух случаях  $\nu_c$  не является экстенсивной.

Оптимальную величину коэффициента  $\alpha = \left[ (1-2\alpha) \frac{2}{\pi} \right]^{1/2\alpha}$  в (4.15a) можно найти, варьируя  $\alpha$ . Это дает  $\alpha_{\max} = 0,103$  для  $\alpha = 0,280$ . Тем не менее для этой величины показателя  $\alpha$  значение  $\nu_{\max}$  необходимо уменьшить до  $0,01 N$ , чтобы величина ошибки в воспроизведении информации не превышала 0,5%. Наоборот, можно

фиксировать величину ошибки, скажем 0,5%, и найти величину  $\alpha$ , которая делает соответствующее значение  $\nu$  максимальным, см. (4.4). Эта процедура дает  $\nu_{\max} = 0,013 N$  для  $\alpha = 0,366$ , что не является существенным улучшением.

Смысл оптимизации, приведенной выше, — прост. Если  $\alpha = 0$ , то память находится в состоянии полной неразберихи /9/, т.е. никаких образов не хранит. С другой стороны, если  $\alpha \geq 1/2$ ,  $\nu \rightarrow \infty$ , то образы легко теряются в фоновом шуме и хранение экстенсивного их числа также становится невозможным. Оптимальное значение  $\alpha$  лежит между этими крайними значениями.

### 5. Обсуждение

Для последовательности весов вида  $\epsilon_{\nu} = \nu^{-\alpha}$  существует такое критическое значение  $\nu_c$ , что воспроизведение образов с индексами  $\nu > \nu_c$  становится невозможным: образы теряются в фоновом шуме. Для  $\nu < \nu_c$  корреляция между хранящимися в памяти и воспроизводимыми образами существует, но стремится к нулю (непрерывно), когда  $\nu \rightarrow \nu_c$  снизу. В этом случае образы постепенно теряют четкость вплоть до полного исчезновения в точке  $\nu_c$ . Если  $0 < \alpha < 1/2$ , тогда  $\nu_c$  является экстенсивной величиной, пропорциональной объему системы.

Все эти результаты, и в частности (4.15), легко понять качественно. При нулевой температуре образ с индексом  $\nu$  должен быть устойчивой неподвижной точкой для динамики, которая задается соотношением

$$S(i) := \operatorname{sgn} \left( \sum_j J_{ij} S(j) \right). \quad (5.1)$$

Константа  $J_{ij}$  определена в (2.1) и, следовательно,

$$J_{ij} \sim \sum_n \epsilon_n \xi_{in} \xi_{jn}, \quad (5.2)$$

тогда (с учетом  $N-1 \approx N$ ) имеем

$$\begin{aligned} \sum_j J_{ij} \xi_{j\nu} &\sim \sum_{n,j} \epsilon_n \xi_{jn} \xi_{j\nu} \xi_{in} = \epsilon_{\nu} \xi_{i\nu} N + \\ &+ \sum_{\mu \neq \nu, j} \epsilon_{\mu} \xi_{i\mu} \xi_{j\mu} \xi_{j\nu}. \end{aligned} \quad (5.3)$$

Здесь сумма по  $\mu (\neq \nu)$  и  $j (\neq i)$  является суммой независимых, одинаково распределенных случайных величин, поэтому её значение (по порядку величины) равно квадратному корню из дисперсии, т.е. равно  $(N \sum_{\mu} \epsilon_{\mu}^2)^{1/2}$ . Значение  $\nu_c$  определяется тогда из условия того, что оба члена в правой части (5.3) имеют одинаковую величину:

$$(\epsilon_{\nu} N)^2 \approx N \sum_{\mu} \epsilon_{\mu}^2. \quad (5.4)$$

Следовательно, для  $\nu = \nu_c$  получаем

$$\epsilon_{\nu}^2 \approx N^{-1} \sum_{\mu} \epsilon_{\mu}^2, \quad (5.5)$$

где, по определению,  $\epsilon_{\nu}^2 = 1/\nu^{2\alpha}$ . Это соотношение, с точностью до (существенного) множителя  $2/\alpha$ , который возникает из интеграла ошибок (4.4), воспроизводит, с учетом (4.8)–(4.10), условие (4.15). Без учета этого множителя оптимальная величина константы  $\alpha$  была бы равна максимуму функции  $(1-2\alpha)^{1/2\alpha}$ , т.е.  $e^{-1} = 0,368$ , который достигается в точке  $\alpha = +0$ .

В заключение суммируем полученные результаты. Мы предлагаем простой и строгий вывод термодинамических свойств модели Хопфилда вблизи насыщения. Этот анализ подтверждает основные результаты работы /9/, а также позволяет распространить модель на случай взвешенных образов. При соответствующем выборе  $\epsilon_{\nu}$  имеет место постепенное уменьшение четкости образов с ростом  $\nu$ , что можно интерпретировать как забывание образа с ростом промежутка времени с момента его поступления в память. Однако необходимы дополнительные основания для того, чтобы решить, является ли процесс забывания внутренним свойством памяти либо связан с течением времени или с обоими этими факторами.

Один из авторов (Я.Л. ван Хеммен) благодарен ОИИИ за гостеприимство, способствовавшее написанию настоящей работы. Он также благодарен Немецкому исследовательскому обществу (Бонн) за финансовую поддержку.

### Приложение

Здесь мы хотим показать, что, если  $c(\underline{t})$  является строго выпуклой функцией, а  $c^*(\underline{m})$  — её преобразование Лежандра и функция  $F(\underline{m})$  гладкая, то выражение для

$$\sup_{\underline{m}} \{ F(\underline{m}) - c^*(\underline{m}) \} \quad (П.1)$$

можно представить в виде

$$\max_{\underline{\mu}} \{ F(\underline{\mu}) - \underline{\mu} \cdot \nabla F(\underline{\mu}) + c(\nabla F(\underline{\mu})) \}, \quad (П.2)$$

где  $\underline{\mu}$  удовлетворяет уравнению для неподвижной точки:

$$\underline{\mu} = \nabla c(\nabla F(\underline{\mu})). \quad (П.3)$$

Доказательство сводится к следующему. Заметим, что  $\underline{t} \rightarrow \nabla c(\underline{t})$  — это отображение  $\mathbb{R}^n$  на  $\mathbb{R}^n$ . Обратное этому отображению существует и равно  $\nabla c^*$ , см. /16/. Поскольку  $\sup$  в (П.1) реализуется



для тех  $\underline{m}$ , которые удовлетворяют соотношению  $\nabla F(\underline{m}) = \nabla c^*(\underline{m})$ , то мы немедленно получаем уравнение для неподвижной точки:

$$\underline{m} = \nabla c(\nabla F(\underline{m})). \quad (\text{II.4})$$

Его решение обозначим символом  $\underline{\mu}$ .

Рассмотрим теперь  $c^*(\underline{\mu})$ . По определению,

$$c^*(\underline{\mu}) = \sup_{\underline{t}} \{ \underline{\mu} \cdot \underline{t} - c(\underline{t}) \}. \quad (\text{II.5})$$

Чтобы получить  $\sup$ , мы должны найти такое  $\underline{t}$ , что

$$\nabla c(\underline{t}) = \underline{\mu}. \quad (\text{II.6})$$

Однако  $\underline{\mu}$  удовлетворяет уравнению (II.4). Сравнивая (II.6) с (II.3), получаем, что  $\underline{t} = \nabla F(\underline{\mu})$ . Если подставить это соотношение в (II.5) и вернуться к (II.1), тогда немедленно получим (II.2).

#### Литература

1. McCulloch W.S. and Pitts W.A. Bull. Math. Biophys., 5 (1943), II5.
2. Hopfield J.J. Proc. Nat. Acad. Sci. USA, 79 (1982) 2554; ibid. 81 (1984) 3088.
3. Little W.A. Math. Biosci., 19 (1974) 101; Little W.A. and Shaw G.L. Math. Biosci., 39 (1978) 281.
4. Peretto P. Biol. Cybern., 50 (1984) 51.
5. Toulouse G., Dehaene S. and Changeux J.-P. Proc. Nat. Acad. Sci. USA, 83 (1986) 1695.
6. van Hemmen J.L. and Kühn R. Phys. Rev. Lett., 57 (1986) 913.
7. van Enter A.C.D. and van Hemmen J.L. Phys. Rev., A 29 (1984) 355.
8. Hebb D. The Organization of Behavior (Wiley, New York, 1949).
9. Amit D.J., Gutfreund H. and Sompolinsky H. Phys. Rev. Lett., 55 (1985) 1530 and Ann. Phys. (N.Y.), to appear.
10. Parisi G. J. Phys. A: Math. Gen., 19 (1986) L 617.
11. van Hemmen J.L. and Palmer R.G. J. Phys. A: Math. Gen., 12 (1979) 563.
12. Grensing D., Kühn R. and van Hemmen J.L. to be published.
13. van Hemmen J.L. Phys. Rev. Lett., 49 (1982) 409.
14. van Hemmen J.L., van Enter A.C.D. and Canisius J. Z. Phys., 50 (1983) 311.

15. van Hemmen J.L. In: Heidelberg Colloquium on Spin Glasses, edited by van Hemmen J.L. and Morgenstern I., Lecture Notes in Physics, 192 (Springer, Heidelberg, 1983), pp. 203-233, in particular the Appendix.
16. Roberts A.W. and Varberg D.E. Convex Functions (Academic Press, New York, 1973).

Рукопись поступила в издательский отдел  
9 марта 1987 года.

### НЕТ ЛИ ПРОБЕЛОВ В ВАШЕЙ БИБЛИОТЕКЕ?

Вы можете получить по почте перечисленные ниже книги, если они не были заказаны ранее.

Д9-82-664	Труды совещания по коллективным методам ускорения. Дубна, 1982.	3 р. 30 к.
Д3,4-82-704	Труды IV Международной школы по нейтронной физике. Дубна, 1982.	5 р. 00 к.
Д11-83-511	Труды совещания по системам и методам аналитических вычислений на ЭВМ и их применению в теоретической физике. Дубна, 1982.	2 р. 50 к.
Д7-83-644	Труды Международной школы-семинара по физике гажелых ионов. Алушта, 1983.	6 р. 55 к.
Д2,13-83-689	Труды рабочего совещания по проблемам излучения и детектирования гравитационных волн. Дубна, 1983.	2 р. 00 к.
Д13-84-63	Труды XI Международного симпозиума по ядерной электронике. Братислава, Чехословакия, 1983.	4 р. 50 к.
Д2-84-366	Труды 7 Международного совещания по проблемам квантовой теории поля. Алушта, 1984.	4 р. 30 к.
Д1,2-84-599	Труды VII Международного семинара по проблемам физики высоких энергий. Дубна, 1984.	5 р. 50 к.
Д17-84-850	Труды III Международного симпозиума по избранным проблемам статистической механики. Дубна, 1984. /2 тома/	7 р. 75 к.
Д10,11-84-818	Труды V Международного совещания по проблемам математического моделирования, программированию и математическим методам решения физических задач. Дубна, 1983	3 р. 50 к.
	Труды IX Всесоюзного совещания по ускорителям заряженных частиц. Дубна, 1984 /2 тома/	13 р. 50 к.
Д4-85-851	Труды Международной школы по структуре ядра, Алушта, 1985.	3 р. 75 к.
Д11-85-791	Труды Международного совещания по аналитическим вычислениям на ЭВМ и их применению в теоретической физике. Дубна, 1985.	4 р.
Д13-85-793	Труды XII Международного симпозиума по ядерной электронике. Дубна, 1985.	4 р. 80 к.
Д3,4,17-86-747	Труды V Международной школы по нейтронной физике. Алушта, 1986.	4 р. 50 к.

Заказы на упомянутые книги могут быть направлены по адресу:  
101000 Москва, Главпочтамт, п/я 79  
Издательский отдел Объединенного института ядерных исследований

Загребнов В.А., ван Хеммен Я.Л.

P17-87-145

Запоминание экстенсивно большого числа взвешенных образов насыщенными нейронными сетями

Эффективность модели Хопфилда для нейронных сетей анализируется в случае, когда число запоминаемых образов экстенсивно велико и им приписаны некоторые веса. Если размеры системы равны  $N$ , тогда она может запомнить  $\alpha N$  образов с соответствующими весами. Эти веса можно связать с временным порядком запомнившихся образов либо, при соответствующем выборе, с процессом постепенного забывания экстенсивно большого числа образов. Особое внимание уделено математической структуре модели.

Работа выполнена в Лаборатории теоретической физики ОИЯИ.

Сообщение Объединенного института ядерных исследований. Дубна 1987

Перевод Г.Г.Сандуковской

Zagrebнов V.A., van Hemmen J.L.

P17-87-145

Strong Extensively Many Weighted Patterns in a Saturated Neural Network

The performance of the Hopfield model of a neural network with extensively many weighted patterns is analyzed. If the system size is  $N$ , then  $N$  patterns, each provided with a suitable weight, are stored. The weights may be associated with a temporal order and, if appropriately chosen, they allow a gradual fading out of the extensively many stored patterns. Particular emphasis is put on the underlying mathematical structure.

The investigation has been performed at the Laboratory of Theoretical Physics, JINR.

Communication of the Joint Institute for Nuclear Research. Dubna 1987