

Кириченко А.А.

Нейропакеты

– современный

интеллектуальный

инструмент

исследователя

2013

Кириченко А.А.

«Нейропакеты – современный интеллектуальный инструмент исследователя», 2013.
Сетевое электронное издание учебного пособия. 297 страниц, 450 рисунков, формат PDF.

ISBN 978-5-9904911-1-3

Всё, что связано с использованием нейронных сетей получило название нейросетевых технологий (нейрокомпьютинг). Они не требуют программирования, а предусматривают работу по обучению нейронной сети на специально подобранных примерах. На этапе обучения формируются основные отношения между входными параметрами и оформляются в незримые таблицы (образы), которые впоследствии используются при решении задач на сети.

Нейрокомпьютинг предоставляет единую методологию решения очень широкого круга практически интересных задач, как правило - ускоряющих и удешевляющих разработку приложений. В число таких задач входят прогнозирование цен, оценка кредитоспособности, оптическое распознавание, обработка изображений, диагностика, лингвистический анализ, и др.

Использование нейросетей для решения перечисленных задач предусматривает выполнение типовой последовательности действий с помощью нейрокомпьютеров, или нейропакетов.

В книге рассматриваются возможности нейропакетов, использование их для моделирования разнообразных нейронных систем, настройка и обучение универсальных нейропакетов на решение задач с использованием трёх наиболее доступных пакетов: BrainMaker, Deductor Academic, IBM SPSS Neural Network.

Книга предназначена для учащихся старших классов, студентов, бакалавров, магистров, аспирантов.

Кириченко А.А. профессор кафедры архитектуры программных систем Федерального государственного автономного образовательного учреждения высшего профессионального образования "Национальный исследовательский университет "Высшая школа экономики" при правительстве РФ".

ISBN 978-5-9904911-1-3

© Кириченко А.А., 2013

Содержание.

Лекция 1. Назначение и возможности нейросетей. Задачи факультатива.....	3
Лекция 2. НейроЭВМ и нейропакеты.	10
Лекция 3. Концепция нейросетевого исследования.	24
Лекция 4. Deductor - флагманский продукт компании BaseGroup Labs.	38
Лекция 5. Deductor. Нейронная сеть. Настройка и обучение.	46
Лекция 6. Самоорганизующиеся карты Кохонена.	73
Лекция 7. Нейросетевое исследование в статистическом пакете SPSS.....	95
Практикум.	117
ПЗ1. Нейросетевая декомпозиция решаемой задачи.	117
ПЗ 2. Работа с нейропакетом BrainMaker.....	128
ПЗ 3. Подготовка табличных данных.....	154
ПЗ 4. Ирисы Фишера.	178
ПЗ 5 Прогнозирование стоимости дома в пакете Deductor.....	181
ПЗ 6. Автоматическая классификация данных об ирисах Фишера с помощью карт Кохонена.	201
ПЗ 7. Использование радиальной основной функции для классификации телекоммуникационных клиентов.	212
Приложение 1. Описание пакета “BrainMaker 3.11 Pro”.....	228
Приложение 2. Использование Excel для подготовки данных к исследованию.	262
Приложение 3. Поиск информации в Интернет.	284
Приложение 4. Исходные данные для задачи «Ирисы Фишера».	291
Литература.	296

Лекция 1. Назначение и возможности нейросетей. Задачи факультатива.

Введение.

Проведение научного исследования чаще всего заключается в выявлении скрытых правил и закономерностей в наборах данных, формулировке гипотез и выявлении типовых структур. Для этого приходится использовать различные методы обнаружения (добычи) знаний: абстрагирование, ассоциативное объединение, классификацию, кластеризацию, анализ временных рядов, прогнозирование и др.

Человеческий разум не приспособлен для восприятия больших массивов разнородной информации. В среднем человек не способен улавливать более двух-трех взаимосвязей даже в небольших выборках.

Для расширения аналитических возможностей человека можно использовать методы традиционной статистики, эвристические решающие устройства на основе экспертных систем, семантический дифференциал, теорию решения изобретательских задач (ТРИЗ), нейронные сети.

Традиционная статистика решает аналогичные задачи, но она оперирует усредненными характеристиками выборки, которые часто являются фиктивными величинами, например - средней платежеспособностью клиента, в то время, как необходимо уметь прогнозировать состоятельность и намерения конкретного клиента с учётом функции риска или функции потерь.

Методы математической статистики, эвристические решающие устройства, семантический дифференциал, ТРИЗ, так же, как и статистика относятся к дискретным методам. Для человека же несвойственно использовать при решении жизненных проблем дискретные методы.

Естественным для человека является использование основных принципов мозга — ассоциативное мышление, использование принципов обучения (самообучения) и адаптации, связей «если - то», «посылка - следствие», лежащих в основе распознавания, движения, управления, принятия решений.

Поэтому из различных способов расширения аналитических возможностей человека наиболее эффективными при исследовании задач, не имеющих общепризнанного алгоритма решения, является использование нейронных сетей.

Всё, что связано с использованием нейронных сетей получило название нейросетевых технологий.

Нейросетевые технологии не требуют программирования, а предусматривают работу по обучению нейронной сети на специально подобранных примерах.

Основной функцией обучения нейросети, воспроизводящей работу мозга и ассоциативное мышление является узнавание, умение определять сходство и различия.

На этапе обучения формируются основные отношения между входными параметрами и оформляются в незримые таблицы (образы), которые впоследствии будут использоваться при решении задач на сети.

Какие задачи решают нейросети.

Нейросети наиболее приспособлены к решению широкого круга задач, так или иначе связанных с обработкой образов. Вот список типичных математических постановок задач для нейросетей:

- Аппроксимация функций по набору точек (регрессия).
- Узнавание, классификация данных по заданному набору классов.
- Кластеризация данных с выявлением заранее неизвестных классов-прототипов.
- Сжатие информации.
- Восстановление утраченных данных.
- Ассоциативное преобразование информации.

Этот список можно было бы продолжить и дальше. За этими задачами просматривается некий единый прототип, позволяющий при известной доле воображения сводить их друг к другу. Чаще всего - это типичные примеры некорректных задач, т.е. задач не имеющих единственного решения, алгоритмы которых неизвестны или не имеют строгого обоснования.

Нейрокомпьютинг предоставляет единую методологию решения очень широкого круга практически интересных задач. Это, как правило, ускоряет и удешевляет разработку приложений.

В число таких задач входят прогнозирование цен, оценка кредитоспособности, оптическое распознавание (например, подписи), обработка изображений, диагностика, лингвистический анализ, и др.

Наверное, в каждой предметной области при ближайшем рассмотрении можно найти постановки нейросетевых задач. Список областей, где решение такого рода задач имеет практическое значение уже сейчас, очень широк:

- Экономика и бизнес: предсказание рынков, автоматический дилинг, оценка риска невозврата кредитов, предсказание банкротств, оценка стоимости недвижимости, выявление пере- и недооцененных компаний, автоматическое рейтингование, оптимизация портфелей, оптимизация товарных и денежных потоков, автоматическое считывание чеков и форм, безопасность транзакций по пластиковым карточкам.

- Политические технологии: анализ и обобщение социологических опросов, предсказание динамики рейтингов, выявление значимых факторов, объективная кластеризация электората, визуализация социальной динамики населения.
- Безопасность и охранные системы: системы идентификации личности, распознавание голоса, лиц в толпе, распознавание автомобильных номеров, анализ аэрокосмических снимков, мониторинг информационных потоков, обнаружение подделок.
- Ввод и обработка информации: Обработка рукописных чеков, распознавание подписей, отпечатков пальцев и голоса. Ввод в компьютер финансовых и налоговых документов.

Нейросети - это не что иное, как новый инструмент анализа данных. И лучше других им может воспользоваться именно специалист в своей предметной области.

Основные трудности на пути еще более широкого распространения нейротехнологий - в неумении широкого круга профессионалов формулировать свои проблемы в терминах, допускающих простое нейросетевое решение, т.е. неумение проводить нейросетевую декомпозицию решаемой задачи.

Основные классы задач, решаемых с помощью нейропакетов.

1. Узнавание (Классификация)
2. Кластеризация
3. Регрессия (прогнозирование, предсказание)
4. Понижение размерности

В задачах регрессии целью является оценка значения числовой (принимающей непрерывный диапазон значений) выходной переменной по значениям входных переменных.

В задачах анализа временных рядов целью является **прогноз** будущих значений переменной, зависящей от времени, на основе предыдущих значений ее и/или других переменных.

Как правило, прогнозируемая переменная является числовой, поэтому прогнозирование временных рядов - это частный случай регрессии. Однако такое ограничение часто в пакет не закладывается, так что в нем можно прогнозировать и временные ряды номинальных (т.е. классифицирующих) переменных.

В задаче узнавания сеть должна отнести каждое наблюдение к одному из нескольких классов (или, в более общем случае, оценить вероятность принадлежности наблюдения к каждому из классов).

Задачи классификации (кластеризации) заключаются в группировке похожих образов в классы (кластеры). В нейропакетах они решаются с помощью сети СОКК (самоорганизующаяся карта Кохонена).

Самоорганизующаяся карта Кохонена (СОКК или сеть Кохонена) принципиально отличается от всех других типов сетей. В то время как все остальные сети предназначены для задач с управляемым обучением (т.е. обучением с учителем), сети Кохонена главным образом рассчитаны на неуправляемое обучение (без учителя).

В алгоритмах неуправляемого обучения значения весов и/или порогов меняются на основании только входных обучающих данных (выходные значения не требуются и, если они присутствуют, игнорируются).

Наиболее общая возможность *понижения размерности* - это анализ главных компонент, который часто в нейропакетах реализован в виде оператора «Что – если?» - он может помочь извлечь небольшое число самых основных компонент из исходных данных довольно большой размерности, сохранив при этом неизменной структуру результатов.

Нейросетевая технология исследований.

Использование нейросетей для решения перечисленных задач предусматривает выполнение типовой последовательности действий:

1. Получение исходных данных.
2. Отбор входных данных и понижение размерности.
3. Оцифровка данных (шкалирование, преобразование номинальных значений, и др.)
4. Выбор подходящей архитектуры сети.
5. Обучение нейронной сети.
6. Тестирование нейронной сети.
7. Получение готового решения

Что собой представляет нейрокомпьютинг.

Исторически электронные вычислительные машины (компьютеры) развивались в виде трёх разновидностей: Цифровые, Аналоговые и Адаптивные (самонастраивающиеся, дообучающиеся).

Выпуск аналоговых ВМ полностью прекратился к 1980 году. Развитие двух оставшихся видов ВМ иллюстрируется следующим графиком:

При конструировании ВМ основное внимание уделялось цифровым вычислительным машинам.

Адаптивные ВМ представляли собой альтернативную, полностью параллельную архитектуру обработки информации -- «по образу и подобию» биологических нервных систем – *нейрокомпьютинг*. Датой рождения этой науки принято считать 1943 год, в котором появилась статья МакКаллока и Питтса о вычислениях в сетях формальных нейронов, хотя работы по созданию адаптивных систем велись ещё в 30-е годы прошлого столетия.

Вначале развития нейрокомпьютинга [1] шло в направлении создания и использования нейросетевых программных продуктов, рынок которых включал в себя

нейропакеты общего назначения, системы разработки нейроприложений, готовые решения на основе нейросетей, нейроконсалтинг. Впоследствии к этому рынку добавилась работа по созданию нейроЭВМ. Характеризовать секторы рынка нейросетевых программных продуктов можно следующим образом:

Нейропакеты общего назначения нацелены на решение информационных задач в диалоговом режиме - при непосредственном участии пользователя. Они не применимы в условиях потоковой обработки данных. Кроме того, они не приспособлены для разработки сложных систем обработки данных, состоящих из многих блоков, содержащих, скажем, сотни нейросетей,

К ним относятся такие пакеты, как BrainMaker Professional, NeuroForecaster, Лора-IQSOO, Штутгартский симулятор для UNIX-машин.

Коммерческие пакеты отличаются от свободно распространяемых большим набором средств импорта и предобработки данных, дополнительными возможностями по анализу значимости входов и оптимизации структуры сети. Стоимость коммерческих эмуляторов - масштаба \$1000. Как правило, такие пакеты (BrainMaker Professional, NeuroForecaster, Лора-IQ300) имеют собственный встроенный блок предобработки данных, хотя иногда для этой цели удобнее использовать стандартные программные средства типа электронных таблиц.

Так, нейро-продукты группы нейрокомпьютинга ФИАН встраиваются непосредственно в Microsoft Excel в качестве специализированных функций обработки данных. При этом всю предобработку данных и визуализацию результатов можно проводить стандартными средствами Excel, который, кроме того, имеет богатый и расширяемый набор конверторов для импорта и экспорта данных. Пакет прикладных программ Excel Neural Package, использовал в качестве функций активизации гиперболический тангенс, а в качестве алгоритма обучения — алгоритм Rprop.

Удобным инструментом разработки сложных нейросистем является MATLAB с прилагающимся к нему нейросетевым инструментарием, органично вписавшимся в матричную идеологию этой системы. MATLAB предоставляет удобную среду для синтеза нейросетевых методик с прочими методами обработки данных (wavelet-анализ, статистика, финансовый анализ и т.д.).

Разработанные в системе MATLAB приложения могут быть затем перетранслированы в C++.

Полностью переведен на русский язык нейросетевой программный продукт *STATISTICA Neural Networks*. Он является средой анализа нейросетевых моделей и соответствует самым современным инструментам.

Инструменты разработки нейроприложений - главное, что отличает этот класс программного обеспечения - способность генерировать "отчуждаемые" нейросетевые продукты, т.е. генерировать программный код, использующий обученные нейросети для обработки данных. Такой код может быть встроен в качестве подсистемы в любые сколь угодно сложные информационные комплексы.

Примерами подобных систем, способных генерировать исходные тексты программ являются NeuralWorks Professional II Plus (стоимостью от \$3000) фирмы NeuralWare и отечественный Neural Bench (нейро-верстак). Последний интересен, кроме прочего, тем, что может генерировать коды на многих языках, включая Java.

Готовые решения на основе нейросетей - это - конечный результат. Здесь нейросети спрятаны от пользователя в недрах готовых автоматизированных комплексов, предназначенных для решения конкретных производственных задач.

Например, такой продукт, как Falcon встраивается в банковскую автоматизированную систему обслуживания платежей по пластиковым карточкам.

В другом случае такая система может выполнять функции автоматизированной системы управления заводом или реактором. Конечного пользователя, как правило, не интересует способ достижения результата, ему важно лишь качество продукта

Поскольку многие такие готовые решения обладают уникальными возможностями (пока специалисты по нейрокомпьютерингу еще в дефиците) и обеспечивают реальные конкурентные преимущества, их цена может быть довольно высока - гораздо выше, чем стоимость нейроЭВМ (neuro-hardware).

Тем не менее, возможность перевода алгоритмов на нейросетевую основу является очень ценной, поскольку в случае реального применения алгоритмов на практике, нейросетевой подход позволяет получать высокоэффективные параллельные архитектуры при аппаратной реализации алгоритмов в устройствах.

Нейросетевой консалтинг используется вместо того, чтобы продавать готовые программы либо инструменты для их разработки. Некоторые задачи, например такие, как предсказание рыночных временных рядов, являются настолько сложными, что доступны лишь настоящим профессионалам. Не каждая компания может позволить себе издержки, ассоциируемые с передовыми научными разработками (например, постоянное участие в международных конференциях). Поэтому приобретают популярность фирмы, единственной продукцией которых являются исследования рынков. При большом числе клиентов цена таких предсказаний может быть весьма умеренной.

Примером здесь может служить Prediction Company, основанная в 1991 году физиками Дойном Фармером и Норманом Паккардом. Продукция компании пользуется большим успехом среди Швейцарских банков, скупающих прогнозы "на корню" для игры на фондовых и валютных рынках.

Фирма Richard Borst, торгующая недвижимостью, применяет предельно дешевый ("университетский") нейропакет для уточнения оценки выставляемых на продажу домов и квартир. Как свидетельствуют старожилы фирмы, внедрение нейропакета (стоимостью всего \$300) увеличило оборот фирмы в Нью-Йорке и Пенсильвании на 6%.

В статье «Нейронная сеть – оружие финансиста» Андрей Масалович пишет: «Мой приятель Джим - журналист-аналитик, работающий в "информационной империи" McGraw-Hill. Весь его офис стоит менее \$2000, помещается в кармане и состоит из портативного palmtop-компьютера, факс-модема и нейропакета Brain Maker. Ежедневно Джим подключается к необъятным базам данных McGraw-Hill и при помощи специальной утилиты DataMaker "просеивает" через свой нейропакет мегабайты финансовой, экономической и прочей информации. После нескольких минут яростного перемалывания тысяч и тысяч разнородных цифровых параметров, нейропакет выдает

(кстати, в виде изящных таблиц Excel) прогноз ряда макроэкономических индикаторов - на завтра, на неделю и на месяц вперед. Попытка проделать такой объем аналитической работы вручную (даже с использованием вспомогательных программ корреляционного анализа), наверное, привела бы беднягу на больничную койку в первый же день».

Зачем эта дисциплина нужна студентам исследовательского университета?

Статус национального исследовательского университета обязывает иметь среди преподаваемых дисциплин методики проведения научных исследований – разнообразные, современные, новые, может быть даже – необычные.

К таким методикам относятся, например, построенные на использовании нейроЭВМ, нейронных сетей ЭВМ, нейроматематики, в простейших случаях (и наиболее доступных) – на нейропакетах.

Такие методики позволяют решать интеллектуальные задачи в различных областях человеческих знаний. Например, специалисты – политологи (политобозреватели) заняты решением задачи, связанной с просмотром огромного количества материалов (новостей) из разных источников за последние сутки. Читать эти новости можно непрерывно. Но нужно среди них отобрать самые важные. Для этой цели разработаны нейропакеты, просматривающие огромный поток сообщений с электронной скоростью и отбирающие только те, которые представляют интерес для данного специалиста.

Поскольку основные задачи, решаемые с помощью нейропакетов, связаны с поиском скрытых закономерностей, классификацией, прогнозированием, сокращением размерности, область их применения охватывает все направления человеческих знаний – и технических, и гуманитарных.

Изучение возможностей нейропакетов, моделирование разнообразных нейронных систем, разработка программ для автоматизации выполнения интеллектуальных операций, настройка и обучение универсальных нейропакетов на решение определённых задач – вот далеко не полный перечень задач, стоящих перед бакалаврами и магистрами национального исследовательского университета.

В результате освоения дисциплины студент должен:

Знать

- Отличительные особенности задач, эффективно решаемых на нейросистемах, принципы обучения нейросетей. Устройство и принцип действия нейрокомпьютеров, нейронных компьютерных сетей, нейропакетов.

Уметь

- создавать нейронные сети и обучающие выборки для них, обучать нейропакеты решению задач классификации, прогнозирования, снижения размерности данных

Иметь навыки

- (приобрести опыт) проведения аналитических исследований на доступных нейропакетах.

Для освоения учебной дисциплины, студенты должны владеть следующими знаниями и компетенциями:

- Знания основ математической статистики и теории вероятностей;
- Знания в объёме курса «Информатика» бакалаврской подготовки;
- Знание английского языка (умение читать и переводить).

Тематический план учебной дисциплины:

Критерии оценки знаний, навыков

- На текущем контроле студент должен продемонстрировать навыки использования пакета BrainMaker при проведении исследований.
- На итоговом контроле студент должен продемонстрировать навыки самостоятельного поиска исходных данных в Интернет и применимого в пакете Deductor или SPSS метода анализа для решения поставленной задачи, интерпретации и представления результатов анализа, формулировки выводов на основе проведенного анализа данных.

Тематика заданий текущего контроля

- Настройка пакета BrainMaker на распознавание десятичных цифр, текста, пространственного расположения объекта;
- Настройка интерфейса пакета BrainMaker;
- Формирование обучающей выборки (текст, графика, звук);
- Прогнозирование в IBM SPSS Neural Networks;
- Классификация и узнавание в нейронных сетях;
- Технология снижения размерности данных. Задача политобозревателя;
- Моделирование ассоциативной памяти в нейронных сетях;
- Формирование отчетности в нейропакетах.

Программные средства

Для успешного освоения дисциплины, студент использует статистические пакеты BrainMaker, Deductor и IBM SPSS Neural Networks

Лекция 2. НейроЭВМ и нейропакеты.

НейроЭВМ.

Универсальные ЭВМ неэффективны при решении трудноформализуемых задач – задач, алгоритм решения которых неизвестен. Они неспособны обучаться решению какой-либо задачи «на примерах». А человек и животные такой способностью обладают! Разница объясняется, видимо, тем, что принципиально отличаются методы решения задач, характерные для живых организмов и ЭВМ.

К числу задач, для решения которых ЭВМ не приспособлены, относятся:

- Задачи предсказания (прогнозирования);
- Нахождение скрытых закономерностей, причинно-следственных связей;
- Распознавание изображений, звуковой информации и иных образов (автоматическая классификация);
- Ассоциативная память;
- Сжатие и кодирование информации в реальном масштабе времени;

- И др.

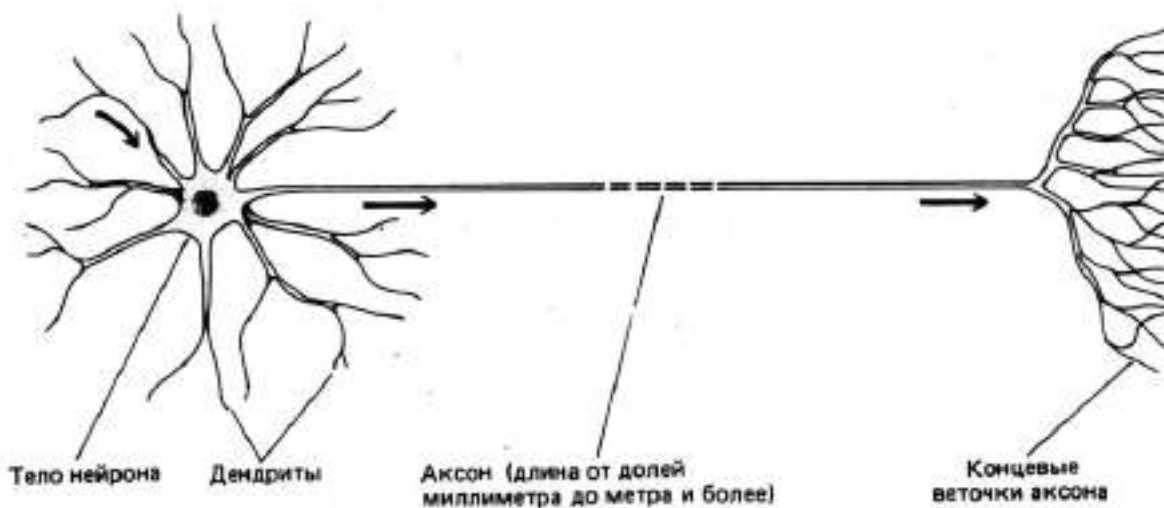
Точные алгоритмы решения таких задач чаще всего отсутствуют. А следовательно, невозможно составить программу для решения такой задачи.

Одновременно с разработкой методов распознавания, реализуемых на универсальных ЭВМ фон-Неймановской архитектуры, ведется поиск методов решения таких задач (в том числе – задач распознавания образов), основанных на иных принципах. Наиболее интересные результаты в этом направлении научных исследований получены при попытках создания распределенных самоорганизующихся систем, например, перцептрона Розенблатта, пандемониума Селфриджа, нейронных ЭВМ (сокращенно – нейроЭВМ), моделирующих работу нервных сетей живых организмов. На их основе в настоящее время разработаны и практически используются параллельные, волновые, матричные и др. системы распознавания.

Элементы, из которых состоит нейросеть, не отличаются разнообразием. Все они имеют одинаковую структуру. Каждый такой элемент суммирует приходящие на него сигналы, и если полученная сумма превышает пороговый уровень, на выходе элемента появляется выходной сигнал. В противном случае элемент не реагирует на входные сигналы. Такие элементы названы *нейронами*, а ЭВМ, построенные на основе таких элементов – нейрокомпьютерами.

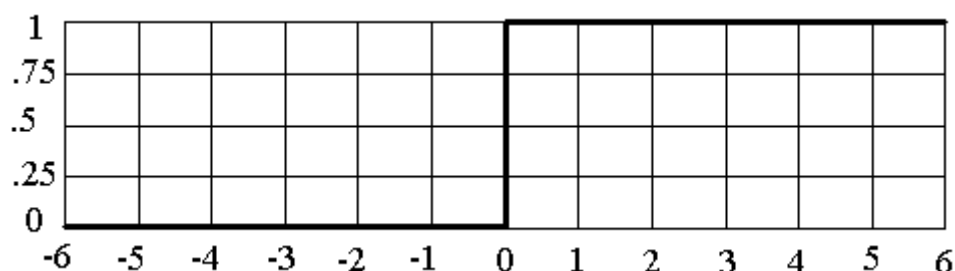
Для нейрокомпьютеров программы не нужны, не нужно и знание того, как решается задача. Их конструкция подсмотрена у природы: элементы, из которых состоит мозг человека и животных; и принципы обработки информации сильно отличаются от элементов и принципов, используемых в ЭВМ.

Нейрон (т.е. нервная клетка) живого организма представляет собой клетку, состоящую из тела клетки, ядра, находящегося в плазме, коротких отростков (дендритов), одного длинного отростка – аксона. Дендриты и аксоны имеют разветвления (синапсы), с помощью которых нейроны связаны друг с другом. Входами нейрона являются дендриты, выходом – аксон.

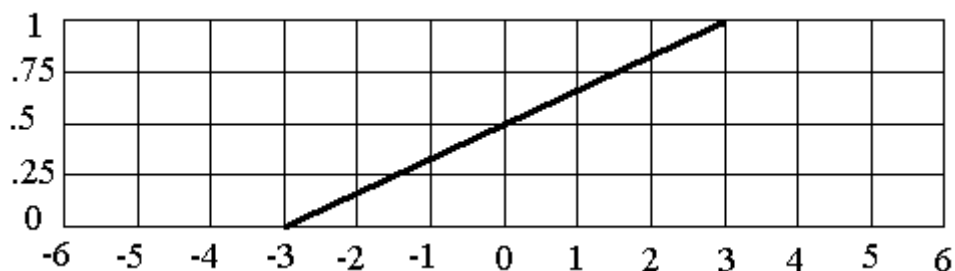


Работает нейрон так: входная информация поступает в него через дендриты и накапливается в теле нейрона. Если сумма поступивших сигналов превышает порог срабатывания, то в аксон выдается импульс, который является выходным сигналом нейрона. Такой нейрон называется *бинарным*.

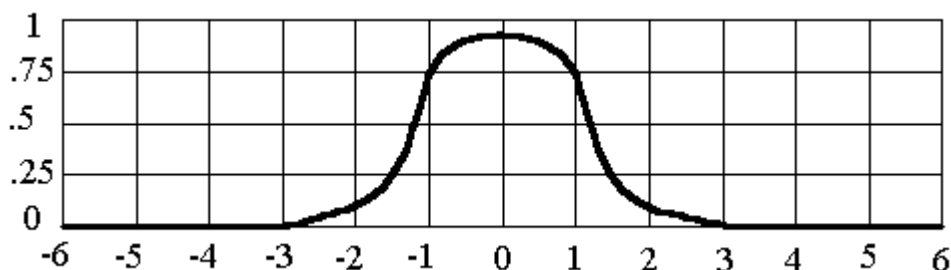
Математическая модель бинарного нейрона представлена на рис.:



Linear – линейная функция возбуждения. Ее графическое представление:



Gaussian – гауссовая. Ее графическое представление:



Алгоритм обучения.

В качестве алгоритма обучения пакет Brain Maker использует метод с обратным распространением ошибки.

Этот метод предусматривает передачу информации об ошибочном срабатывании нейронов в направлении, обратном направлению распространения возбуждения. Метод получил широкое распространение для организации обучения в нейронных сетях различного типа – импульсных, статических, дискретных и аналоговых.

Алгоритм модификации весов синапсов включает ряд шагов:

ШАГ 1. Инициализация весов и порогов.

Устанавливаем малые случайные величины для весов синапсов и порогов нейронов.

ШАГ 2. Предъявление входной реализации и эталона на выходе сети.

Вводятся непрерывные величины компонент входного вектора x_0, x_1, \dots, x_{N-1} и соответствующего эталона выходного вектора y_0, y_1, \dots, y_{M-1} . Если сеть используется как классификатор, тогда все каналы выхода устанавливаются равными 0 за исключением канала, соответствующего требуемому классу, которому присваивается значение 1. Входные реализации предъявляются циклически, пока не стабилизируются веса синапсов.

ШАГ 3. Вычисление действительных выходов.

Используя сигмоидальную нелинейность, у которой

$$\alpha - \Theta = \sum_{i=0}^{N-1} w_{ij} \cdot x_i - \Theta_j,$$

вычисляем y_0, y_1, \dots, y_{M-1} .