

Introdução as Redes Neurais Artificiais

Florianópolis, maio de 2002.

“Nunca se achou que o degrau da escada se destinasse a alguém permanecer em cima dele, mas sim que se destina a sustentar o pé de um homem pelo tempo suficiente para que coloque o outro um pouco mais alto”

Huxley

Sumário

SUMÁRIO	3
LISTA DE FIGURAS	4
LISTA DE TABELAS.....	5
1 O SISTEMA NERVOSO HUMANO	6
1.1 Introdução.....	6
1.2 Cérebro e Conduta	6
1.3 O Neurônio.....	7
2 AS REDES NEURAIS ARTIFICIAIS	9
2.1 Introdução.....	9
2.2 Histórico	10
2.3 Aplicações.....	11
2.4 O Neurônio Artificial	12
2.5 Arquiteturas.....	15
2.6 Aprendizado.....	17
2.6.1 Supervisionado	18
2.6.2 Não supervisionado	19
2.6.3 Velocidade de aprendizado.....	19
2.6.4 Algoritmos de aprendizado.....	20
2.7 Redes Perceptron	21
2.8 A lei de aprendizado do perceptron.....	22
2.9 Limitações: O problema do OU-EXCLUSIVO	24
2.10 Redes Multilayer Perceptron	25
2.11 Algoritmo de treinamento das redes MLP.....	26
3 REFERÊNCIAS BIBLIOGRÁFICAS	31

Lista de Figuras

Figura 1.1 - O Arco reflexo [MINb 00]	7
Figura 1.2 - Estrutura do Neurônio [KAN 98].....	8
Figura 2.1 - O neurônio de McCulloch e implementações de algumas funções booleanas [KOV 96].....	13
Figura 2.2 - O Neurônio artificial [TAFb 96]	14
Figura 2.3 - Funções de transferência [KOV 96].....	15
Figura 2.4 - Rede neural artificial	15
Figura 2.5 - RNA de uma única camada	16
Figura 2.6 - RNA multicamada	16
Figura 2.7 - RNA feedforward ou acíclica	17
Figura 2.8 - RNA feedback ou cíclica.....	17
Figura 2.9 - O perceptron elementar de Roseblatt [BIS 95].....	21
Figura 2.10 - A unidade de processamento do perceptron	22
Figura 2.11 - Perceptron de duas entradas	24
Figura 2.12 - Plano que representa as combinações possíveis do XOR.....	25
Figura 2.13 - Uma rede MLP	26
Figura 2.14 - Rede MLP com os acoplamentos retrógrados para os ajustes sinápticos [PAT 95]	28
Figura 2.15 - Fluxo de treinamento de uma MLP com <i>backpropagation</i>	29

Lista de Tabelas

Tabela 2.1 - Tabela verdade do ou-exclusivo.....	24
--	----

1 O SISTEMA NERVOSO HUMANO

1.1 Introdução

O propósito principal da Neurociência é compreender como o encéfalo produz a acentuada individualidade da ação humana. A compreensão da conduta humana necessita de um estudo neurobiológico que parte da mente e atravessa o mundo molecular, ou seja, como se relacionam as moléculas responsáveis pelas atividades das células nervosas para resultar nos complexos processos mentais.

O encéfalo é uma rede de mais de 100.000 milhões (100 bilhões) de células nervosas delimitadas, que se interconectam em sistemas que produzem nossa percepção do mundo exterior, concentram nossa atenção e controlam o mecanismo da ação. Por tanto, o primeiro passo para conhecer a mente é entender como os neurônios se organizam em vias de comunicação e como as células nervosas individualizadas do encéfalo se comunicam mediante transmissão sináptica. Por fim, seria necessário estudar como alterações em genes individuais afetam a comunicação entre células nervosas e como alterações na comunicação alteram a conduta [KAN 98].

Esses estudos, desde a embriogênese até a neurofisiologia são objeto de estudo da Neurociência que utilizou, inicialmente, a Anatomia macroscópica via disseções anatômicas de órgãos e sistemas (clássicos gregos, Idade Média), depois a Anatomia microscópica (Histologia e Citologia, onde se encontram os célebres estudos de Ramón y Cajal - Espanha), a Biofísica e, recentemente, a Biologia Molecular como ferramenta de estudo. Assim, a Neurociência surgiu no último século como resultado de estudos do sistema nervoso realizados por várias disciplinas clássicas. Atualmente, novas técnicas aportam os meios para vincular diretamente a dinâmica molecular de células nervosas individuais com representações de atos perceptivos e motores do encéfalo e para relacionar estes mecanismos com a conduta observável. As novas técnicas de neuroimagem, por exemplo, permitem observar o encéfalo humano em ação (identificar as regiões específicas que se associam com o pensamento e o sentimento). A neurociência, com sua capacidade de interligar a Biologia Molecular e os estudos cognitivos possibilitou que se comece a explorar a Biologia do potencial humano, de modo que, possamos entender o que determina o que somos (por exemplo, estudos recentes sobre neurotransmissores, associam o comportamento mais agitado, a busca por esportes de riscos, a pessoas com maior quantidade de receptores adrenérgicos em suas terminações nervosas; os problemas de depressão ligados diretamente a quantidade de mediadores serotoninérgicos na corrente sanguínea, e assim por diante.

1.2 Cérebro e Conduta

Talvez a última fronteira da ciência seja entender as bases biológicas da consciência e dos processos mentais pelos que percebemos, atuamos, aprendemos e recordamos.

A tarefa da Neurociência é aportar explicações da conduta em termos de atividades do encéfalo, explicar como atuam os milhões de células nervosas individuais no encéfalo para produzir a conduta e como, por sua vez, estas células estão influenciadas pelo meio ambiente, incluindo a conduta de outros indivíduos. Para responder perguntas como – “Estão localizados os processos mentais em regiões específicas do encéfalo ou representam uma propriedade coletiva e emergente do encéfalo em sua totalidade?” “Vários processos mentais podem localizar-se em diferentes regiões cefálicas?” “Que regras relacionam a Anatomia e a fisiologia de uma região com sua função específica na percepção, no pensamento ou no movimento?” é necessário entender como está estruturado o Sistema Nervoso Central em todos aspectos desde sua embriogênese até sua histofisiologia, biologia molecular e também incluir possíveis alterações teratógênicas.

O estudo integral destes aspectos não pode ser abordado em uma síntese introdutória do Sistema Nervoso, mas que, no entanto é necessária para a construção da estrutura do presente estudo, ou seja, a construção de redes neuronais. Ressalta-se que estes estudos devem realizados pois o simples modelo de organização da estrutura das redes neuronais durante a embriogênese [HAR 87], por exemplo, pode ser útil na estruturação de uma rede neuronal artificial [CHU 92].

1.3 O Neurônio

Todos os animais, inclusive o homem, obtêm informação sobre o seu entorno através de vários receptores sensoriais. A informação conseguida pelos receptores se transforma no encéfalo em percepções ou ordens para o movimento. Respostas tão notáveis são conseguidas somente com a utilização de células nervosas e as conexões estabelecidas entre elas. O comentário realizado neste parágrafo refere-se ao que em fisiologia se denomina de “Arco Reflexo”. Neste caso, como é uma resposta elaborada e interpretada pelos centros nervosos superiores, diz-se que se trata de um Arco Reflexo Central. Já um reflexo de sobrevivência (saltar durante um susto, retirar a mão de uma superfície quente, reflexo patelar, e outras coisas do gênero) são respostas imediatas, sem interpretação detalhada e coordenadas pela medula espinhal, sendo denominadas de Arco Reflexo Periférico. O esquema a seguir ilustra um arco-reflexo de forma simplificada.

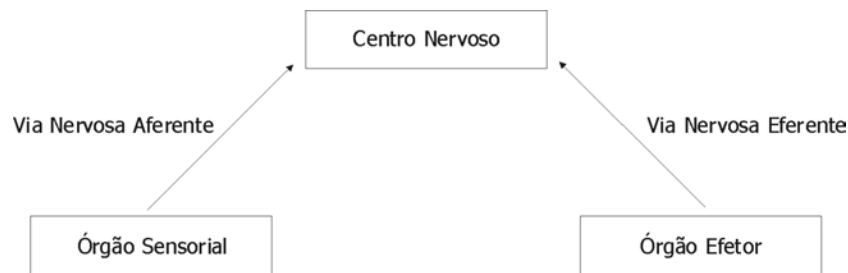


Figura 1.1 - O Arco reflexo [MINb 00]

As unidades básicas do encéfalo, as células nervosas, são muito simples. O encéfalo é capaz de gerar comportamentos tremendamente complexos porque tem uma grande quantidade de células nervosas que se comunicam entre si mediante interconexões específicas. As células nervosas apesar de sua grande quantidade compartilham muitas características. Um dos descobrimentos mais importantes para a compreensão do encéfalo foi que o potencial de ação para produzir condutas complexas não depende, em grande parte, da variedade das células nervosas, mas sim de seu número e de suas conexões específicas entre si e com os receptores sensoriais e os músculos.

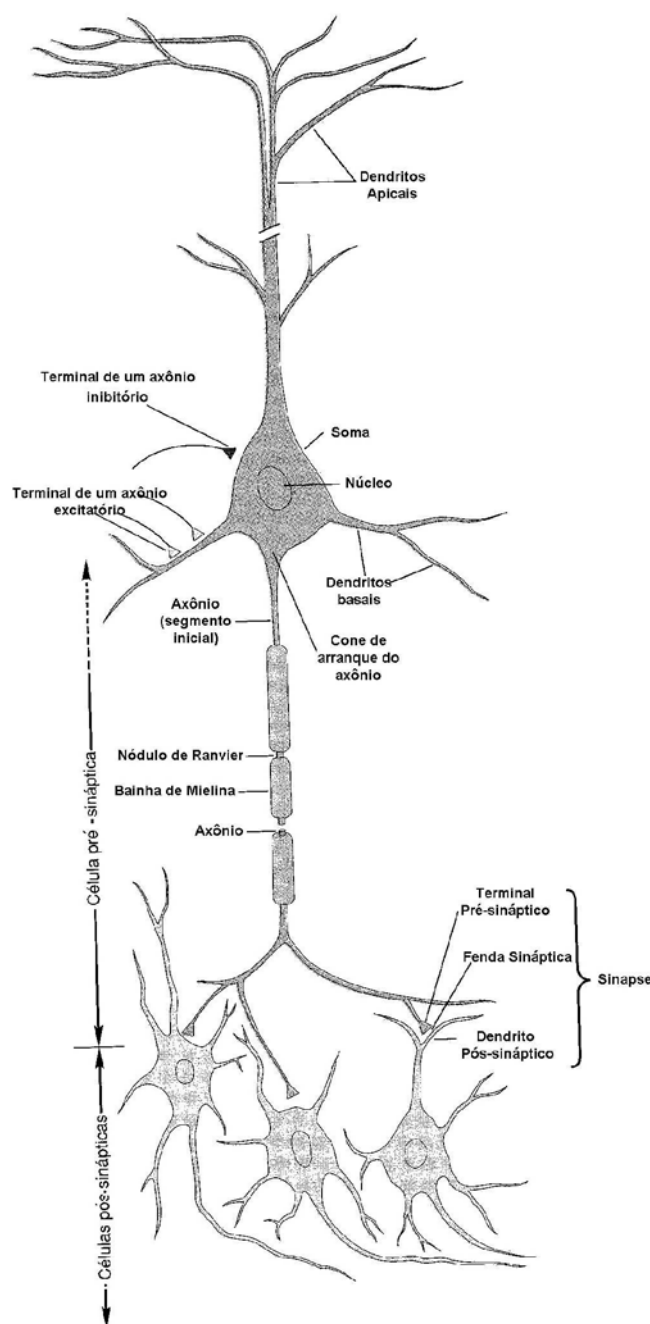


Figura 1.2 - Estrutura do Neurônio [KAN 98]

SINAPSE é a comunicação estabelecida entre um neurônio com outro(s) neurônio(s) ou com outros tecidos. A condução de um estímulo elétrico pela membrana celular de um neurônio é unidirecional, sendo assim, em uma comunicação entre um neurônio e outra célula qualquer, sempre se observa a distinção entre o neurônio que faz a sinapse e seu elemento subsequente. Assim, o elemento que fica antes da comunicação (sinapse) celular é denominado pré-sináptico e o que fica depois de pós-sináptico. O espaço entre o elemento pré-sináptico e o pós-sináptico é denominado de fenda sináptica e é onde são liberados os mediadores químicos inibidores ou excitadores de membrana.

Os temas superficialmente discutidos até o presente momento de forma simplificada servem de base para a compreensão do estudo apresentado a seguir sobre **redes neuronais artificiais**.

2 AS REDES NEURAIS ARTIFICIAIS

2.1 Introdução

A tecnologia das Redes Neurais Artificiais (RNA's) visa solucionar problemas de reconhecimento de padrões que geralmente são baseados em um conjunto de informações previamente conhecido. Geralmente os conjuntos de dados são divididos em conjunto de treinamento e conjunto de teste. Atualmente, pesquisadores em RNA's estão buscando uma compreensão das capacidades da natureza humana, as quais possibilitam que as pessoas construam soluções para problemas que não sejam resolvidos através de métodos tradicionais.

As redes neurais artificiais visam na sua maioria solucionar problemas de inteligência artificial, modelando sistemas através de circuitos (conexões) que possam simular o sistema nervoso humano, abrangendo a capacidade que o mesmo possui de aprender e agir perante as mais adversas situações apresentadas, bem como adquirir conhecimento através da experiência e da observação.

Segundo o pesquisador da Universidade de Helsinki, Teuvo Kohonen, uma rede neural artificial tem a seguinte definição: "uma rede massivamente paralela de elementos interconectados e suas organizações hierárquicas que estão preparadas para iterar com objetos do mundo real do mesmo modo que um sistema nervoso biológico faz".

A complexidade das estruturas elementares das Redes Neurais Biológicas é muito maior do que a dos modelos matemáticos usados nas Redes Neurais Artificiais, demonstrando as dificuldades encontradas para se tentar imitar o funcionamento do sistema nervoso humano. O sistema nervoso é formado por bilhões de células nervosas, enquanto que uma rede neural artificial possui de dezenas a no máximo milhares de unidades de processamento (neurônios).

Uma rede neural artificial pode ser vista como um conjunto de várias unidades interconectadas (similar à estrutura do cérebro), denominadas de neurônios artificiais, cada qual contendo uma pequena porção local de memória. Estes conceitos foram baseados e fundamentados nos estudos realizados nas células nervosas naturais. Portanto, busca-se aproximar ao máximo o funcionamento das redes neurais artificiais das redes neurais biológicas, na tentativa de buscar a desenvoltura com que o cérebro humano desempenha suas funções.

Alguns modelos de redes neurais artificiais possuem muitos neurônios conectados numa estrutura de pesos de conexão e com facilidade de adaptação, proporcionando uma estrutura paralela. A estrutura paralela é desejável pois se algum(s) neurônio(s) falhar (em), os efeitos na rede como um todo não será significativo para o desempenho do sistema se outro caminho de conexão entre os neurônios puder burlar a falha, surgindo então a tolerância à falha.

A princípio, as redes neurais podem calcular qualquer função computável que é realizada em um computador digital, ou seja, possuem a capacidade de modelar relações lineares e não lineares.

Principais características das RNA's [BAR 99]:

- capacidade de "aprender" através de exemplos e de generalizar este aprendizado de forma a reconhecer elementos similares, que não foram apresentados no conjunto de exemplos (treinamento);
- bom desempenho em tarefas pouco ou mal definidas, onde falta o conhecimento explícito de como resolvê-las, o aprendizado se dá através de exemplos;
- robustez à presença de informações falsas ou ausentes, escolha dos elementos no próprio conjunto de treinamento (integridade do conjunto de treinamento);

- no contexto de classificação de padrões, uma rede neural pode fornecer informações sobre quais padrões selecionar em função do grau de confiança apresentado (confiabilidade do conjunto de treinamento);
- tolerância à falha.

2.2 Histórico

As Redes Neurais Artificiais surgiram na década de 40, mais precisamente em 1943, quando o neurofisiologista Warren McCulloch e o matemático Walter Pitts, da Universidade de Illinois, fizeram uma analogia entre as células nervosas e o processo eletrônico num artigo publicado no *Bulletin of Mathematical Biophysics* com o título: *A Logical Calculus of the Ideas Immanent in Nervous Activity*.

Em 1949, o biólogo e psicólogo Donald Hebb, que estudava o comportamento dos animais, escreveu um livro chamado *The Organization of Behavior*, que reforçava as teorias de que o condicionamento psicológico estava presente em qualquer parte dos animais, pelo fato de que esta é uma propriedade de neurônios individuais. As idéias de Hebb não eram pioneiras, mas ele propôs um princípio de aprendizado em sistemas nervosos complexos, ou seja, uma lei que descreve o funcionamento quantitativo da sinapse e do processo de treinamento humano.

Desde, então, vários outros pesquisadores, entusiasmados com as novas descobertas, voltaram-se para esta linha de pesquisa.

Em 1951, Marvin Minsky, co-fundador do Laboratório de Inteligência Artificial do MIT, construiu o SNARC, o primeiro simulador de cadeia neural. O SNARC trabalhava com êxito e podia ajustar seus pesos sinápticos automaticamente. Ele nunca chegou a executar alguma função de processamento de informação interessante, servindo somente de fator motivador para idéias que surgiram posteriormente.

Em 1956, na Primeira Conferência Internacional de Inteligência Artificial, foi apresentado um modelo de rede neural artificial pelo pesquisador da IBM Nathaniel Rochester. Seu modelo consistia numa simulação de centenas de neurônios interconectados através de um sistema que verificaria como a rede responderia aos estímulos ambientais.

Já em 1959, Frank Rosenblatt na Universidade de Cornell, criou uma rede de múltiplos neurônios do tipo discriminadores lineares e a batizou de rede perceptron. Rosenblatt baseou-se nas linhas de pensamento de McCulloch para desenvolver o seu modelo matemático de sinapse humana. Devido as suas complexas pesquisas e inúmeras contribuições técnicas, muitos o consideram como fundador da neurocomputação.

No final da década de 50, Minsky e Seymour Papert lançaram em uma obra chamada *Perceptron*, a qual demonstrava que o modelo apresentado por Rosenblatt não era muito promissor, devido ao uso de técnicas empíricas, das grandes dificuldades da matemática envolvida e dos poucos recursos computacionais disponíveis na época. A publicação de Minsky e Papert acabou esfriando as pesquisas e praticamente todo o investimento financeiro nesta área foi cancelado.

Enquanto Rosenblatt trabalhava no perceptron, Bernard Widrow da Universidade de Stanford com a ajuda de alguns estudantes desenvolveu um novo modelo de processamento de redes neurais chamado de Adaline (ADaptive LINear Element), a qual se destacava pela sua poderosa lei de aprendizado. O princípio de treinamento para as redes Adalines ficou conhecido como a Regra Delta, que foi mais tarde generalizada para redes com modelos neurais mais sofisticados. Mais tarde, Widrow criou a Madaline, que era uma generalização multidimensional do adaline.

Nos anos seguintes, muitos artigos foram publicados, e várias previsões exageradas e pouco confiáveis para a época foram anunciadas [TAFb 96]. A maioria destas suposições falava de computadores com um poder de raciocínio e/ou processamento igual ou superior ao do

cérebro humano. Desta forma, a credibilidade de futuros estudos das RNA's foram fortemente comprometidos.

No início da década de 80, muitos pesquisadores publicaram inúmeras propostas para a exploração de desenvolvimento e pesquisa em redes neurais. Foi quando o administrador de programas da DARPA (Defense Advanced Research Projects Agency) Ira Skurnick resolveu dar atenção às proposições da neurocomputação, contrariando todos os preceitos, fundando em 1983 as pesquisas em neurocomputação da DARPA. Este fato acabou abrindo novos horizontes para a neurocomputação.

O físico e biólogo de reputação mundial John Hopfield também se interessou pela neurocomputação e escreveu vários artigos em 1982 que levaram vários cientistas a se unirem nesta nova área emergente. Hopfield reascendeu as pesquisas em neurocomputação, criticando fortemente as teorias apresentadas por Minsky e Papert na década de 50.

Este campo de pesquisa explodiu mesmo em 1986, quando o professor de psicologia da Universidade de Stanford, David E. Rumelhart, e seu colega James L. McClelland, professor de psicologia da Universidade de CarnegieMellon, publicaram o livro *Parallel Distributed Processing: Explorations in the Microstructure of Cognition (vol.1: Foundations, vol.2: Psychological and Biological Models)*. Nesse livro, eles apresentam um modelo matemático e computacional que propicia o treinamento supervisionado dos neurônios artificiais. Surgia, então, o algoritmo *backpropagation*, um algoritmo de otimização global sem restrições.

Em 1987 ocorreu a Primeira Conferência de Redes Neurais. Também foi formada a Sociedade Internacional de Redes Neurais (*International Neural Networks Society - INNS*) juntamente com o *INNS Journal* em 1989, do *Neural Computation* e do *IEEE Transactions on Neural Networks* em 1990.

A partir destes acontecimentos, muitas instituições formaram institutos de pesquisa e programas de educação em neurocomputação.

2.3 Aplicações

Um dos principais objetivos da pesquisa sobre redes neurais artificiais na computação é desenvolver modelos matemáticos das estruturas neurais, não necessariamente baseadas na biologia, que podem efetuar diversas funções. Na maior parte dos casos, os modelos neurais são compostos por conjuntos de elementos não lineares que operam em paralelo e que são classificados de acordo com modelos/padrões relacionados à biologia. Quando um método é criado visando utilizar aspectos de redes neurais artificiais, começam com o desenvolvimento de um neurônio artificial ou computacional baseado no entendimento de estruturas neurais biológicas, seguidas do aprendizado de mecanismos voltados para um determinado conjunto de aplicações e o treinamento do suposto sistema. Segue-se mais detalhadamente as seguintes fases:

- estudo do problema;
- desenvolvimento de modelos neurais motivados por neurônios biológicos;
- modelos de estruturas e conexões sinápticas;
- escolha de um algoritmo de aprendizado (um método de ajuste de pesos ou forças de conexões internodais);
- construção de um conjunto de treinamento;
- o treinamento propriamente dito;
- fase de testes;
- utilização da rede.

As diferenças entre as aplicações, os algoritmos de aprendizagem e as estruturas de interconexões entre os neurônios levam os pesquisadores a desenvolver diferentes modelos (arquiteturas) de redes neurais. Do ponto de vista estrutural, a arquitetura de redes neurais pode ser classificada como estática, dinâmica ou fuzzy, podendo ter uma ou múltiplas camadas. Além disso, diferenças computacionais surgem devido a forma como são feitas as conexões entre os neurônios. Estas conexões podem ser *feed forward*, *backward*, lateralmente conectadas, topologicamente ordenadas ou híbridas. As aplicações de redes neurais podem ser classificadas em diversas classes como:

- reconhecimento e classificação de padrões;
- processamento de imagem;
- visão computacional;
- identificação e controle de sistemas;
- processamento de sinais;
- robótica;
- filtros contra ruídos eletrônicos;
- análise do mercado financeiro;
- controle de processos.

Cabe ressaltar que em uma determinada aplicação de um sistema, que faz o uso das redes neurais artificiais, não precisa necessariamente ser classificada em apenas uma das citadas acima.

2.4 O Neurônio Artificial

O primeiro modelo matemático para uma rede neural, proposto por McCulloch e Pitts, era simples diante das informações disponíveis naquela época sobre o funcionamento elétrico de uma célula nervosa (Figura 2.1). Era um dispositivo binário, sendo que a saída do neurônio poderia ser pulso ou não pulso (ativo ou não), e as várias entradas tinham um ganho arbitrário, podendo ser excitatórias ou inibitórias. Para se determinar a saída do neurônio, calculava-se a soma ponderada das entradas com os respectivos ganhos como fatores de ponderação, excitatórios ou inibitórios. Se o resultado atingisse um certo limiar, a saída do neurônio era pulso (ativo), caso contrário, não pulso (não ativo).

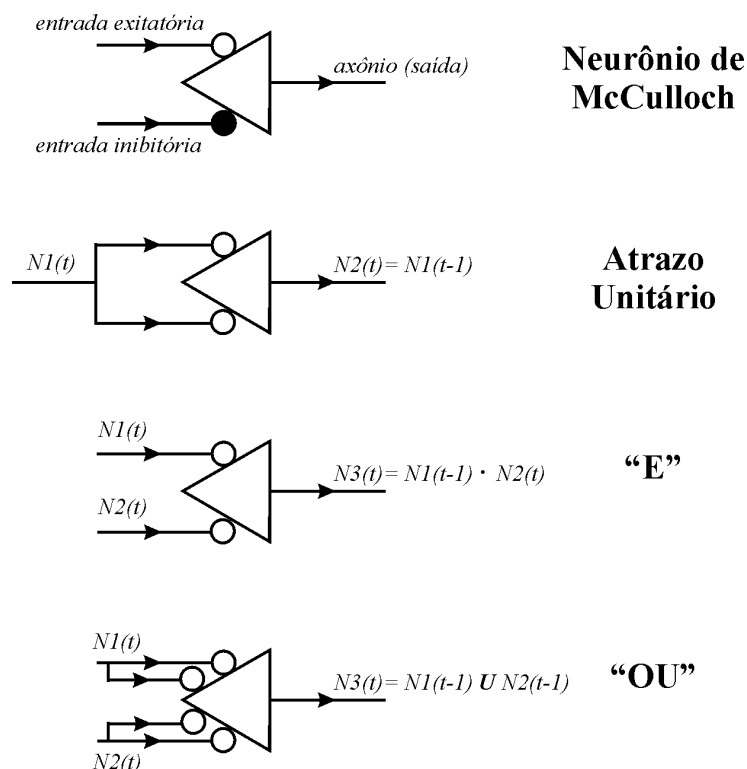


Figura 2.1 - O neurônio de McCulloch e implementações de algumas funções booleanas [KOV 96]

Assim como o neurônio biológico, o neurônio artificial possui um ou mais sinais de entrada e apenas um sinal de saída. As informações podem ser recebidas através de sensores ou de outros neurônios artificiais que fazem parte da Rede Neural Artificial (RNA). Estes sinais são processados e enviados para a saída. Os sinais de entrada (estímulos) devem chegar até o neurônio simultaneamente, isto é, todas as informações devem chegar ao núcleo do neurônio artificial ao mesmo tempo.

O processamento paralelo em computadores seqüenciais (por exemplo, os microcomputadores atuais) pode ser paradoxal, mas não o é, ocorre de fato. A simulação de um ambiente paralelo é possível, e é desta forma que ocorre esse tipo de processamento para as redes neurais. O modelo matemático simula o paralelismo da rede neural através de um algoritmo [TAF 96].

Um dos atributos de grande importância do neurônio artificial é o peso. Os pesos, também conhecidos por pesos sinápticos, são representados pela letra **w** (*weight*) e representam o grau de importância que determinada entrada possui em relação àquele determinado neurônio.

O valor do peso é alterado em função da intensidade do sinal de entrada, e dessa forma, o peso muda o seu valor representativo para a rede (processo de aprendizagem). Deduz-se que, quanto mais estimulada for uma entrada, mais estimulado será o peso correspondente, e quanto mais for estimulado um peso, mais significativo e influente o mesmo será para o resultado do sinal de saída do respectivo neurônio.

Matematicamente, os pesos são vistos como um vetor de valores $[w_1, w_2, \dots, w_n]$ para um neurônio, ou uma matriz de pesos, coleção de vetores, para um conjunto de neurônios.

O sinal de excitação do neurônio é resultante do somatório do produto dos sinais de entrada, representados por um vetor $[x_1, x_2, \dots, x_n]$, pelo vetor de pesos do neurônio $(\sum_{i=0}^n x_i w_i -$ o valor correspondente a $x_0 w_0$ será explicado adiante e corresponde ao viés, representando um estímulo inicial a rede). Após esta operação, os sinais de entrada passam a ser chamados de entradas ponderadas.

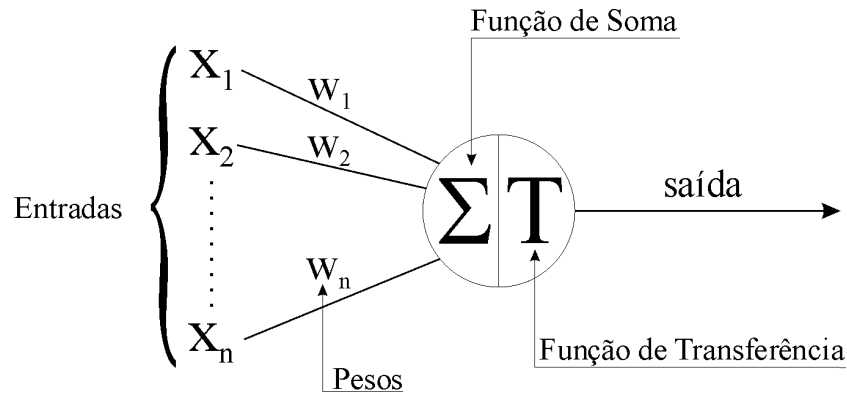


Figura 2.2 - O Neurônio artificial [TAFb 96]

A próxima tarefa a ser tomada pelo neurônio, é a de verificar se o valor resultante da soma entre o produto dos sinais de entrada pelos respectivos pesos atingiu ou não um valor predeterminado, chamado de limiar. Se o valor calculado atingiu o limiar, o mesmo é repassado adiante através da saída. Caso contrário, se o valor não atingiu o limiar, o sinal não será transferido. Esse processo de verificação é chamado de função de transferência, que também é conhecido como limiar lógico.

A resposta final da rede ou das camadas subjacentes está diretamente ligada com o resultado obtido pela função de transferência. Por isso, deve-se dar a devida atenção a este processo. A lógica neural expõe, que a intensidade dos sinais de entrada, dispara, ou não, o sinal do neurônio, fazendo com que este estimule o neurônio seguinte [TAFb 96].

Além da função de transferência, há a função de ativação, a qual antecede a mesma e tem como função, suceder um nível de ativação dentro do próprio neurônio, ou seja, o neurônio, através desta função, decidirá o que fazer com o resultado da soma ponderada das entradas (ativar ou não). Essa decisão tem efeito somente ao respectivo neurônio artificial.

Em alguns modelos simples de redes neurais artificiais, a função de ativação pode ser a própria função de soma das entradas ponderadas do neurônio. Já em modelos mais elaborados, a função de ativação pode possuir um processamento atribuído, o qual pode ser, por exemplo, o uso de um valor prévio de saída como uma entrada para o próprio neurônio, servindo de auto-excitação para o mesmo [TAFb 96].

O valor de saída do neurônio será produzido após a chamada da função de ativação, seguido pela função de transferência.

Em alguns casos, o neurônio artificial pode não ter efeito no neurônio seguinte se o valor de ativação não ultrapassar um certo valor mínimo. Este fator é resultante das características sigma ou ríspidas que a função de transferência tem como propriedade. Devido a esse fator, há vários tipos de funções de transferência (Figura 2.3).

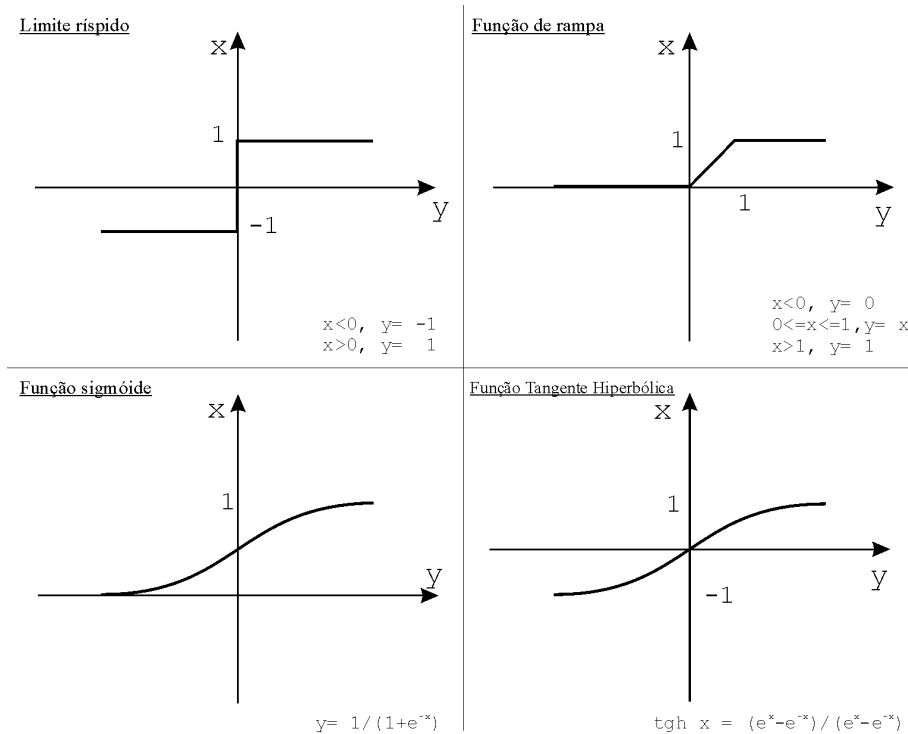


Figura 2.3 - Funções de transferência [KOV 96]

Assim como nas redes neurais biológicas, o conjunto de vários neurônios artificiais interconectados, formam as redes neurais artificiais.

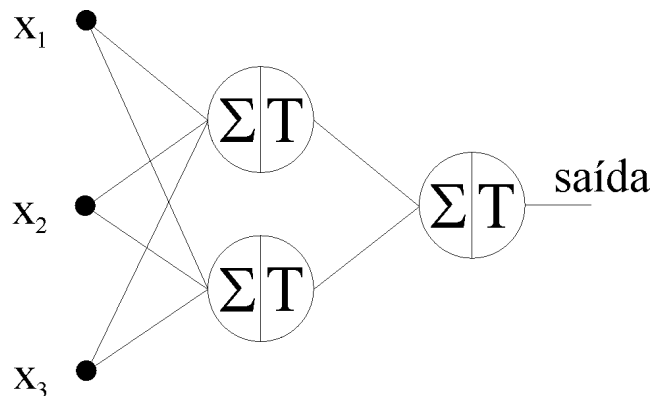


Figura 2.4 - Rede neural artificial

2.5 Arquiteturas

Um outro detalhe importante a ser considerado é a maneira como os neurônios artificiais podem ser agrupados. Este agrupamento se sucede no cérebro humano de maneira que as informações possam ser processadas de forma dinâmica ou interativa. Biologicamente, as redes neurais são organizadas e construídas de forma tridimensional por componentes microscópicos. Há uma forte restrição no número de camadas que a rede pode conter, limitando consideravelmente o tipo e o escopo da implementação da mesma em silício.

Uma rede neural pode ter uma ou várias camadas. As redes que possuem uma única camada são as redes que possuem um nó entre uma entrada e uma saída da rede (Figura 2.5). Esse tipo de rede é indicado para a solução de problemas linearmente separáveis. Já as redes

multicamadas possuem mais de uma camada entre as já existentes camadas de entrada e saída (Figura 2.6).

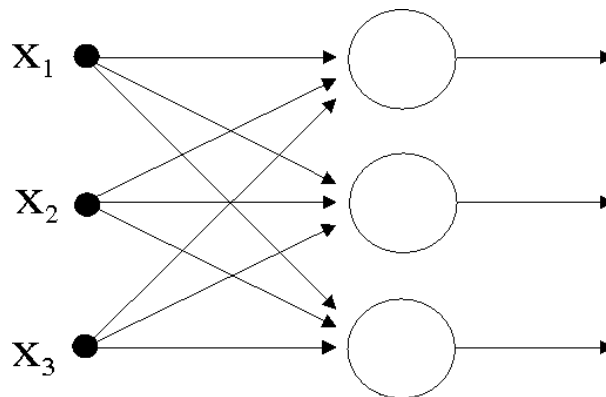


Figura 2.5 - RNA de uma única camada

As redes neurais artificiais multicamadas possuem as chamadas camadas escondidas (hidden), que também são chamadas de intermediárias ou ocultas. Esse número de camadas pode ser indeterminado, e estão situadas entre a camada de entrada e a camada de saída da rede neural [CAR 98].

As camadas ocultas são constituídas por neurônios artificiais, da mesma forma com que as camadas externas (entrada e saída) são compostas, e tendo como característica diferenciada o não contato com o mundo externo (Figura 2.6). Os sinais são passados para os outros neurônios obedecendo às funções de transferência que cada neurônio possui [NAS 94].

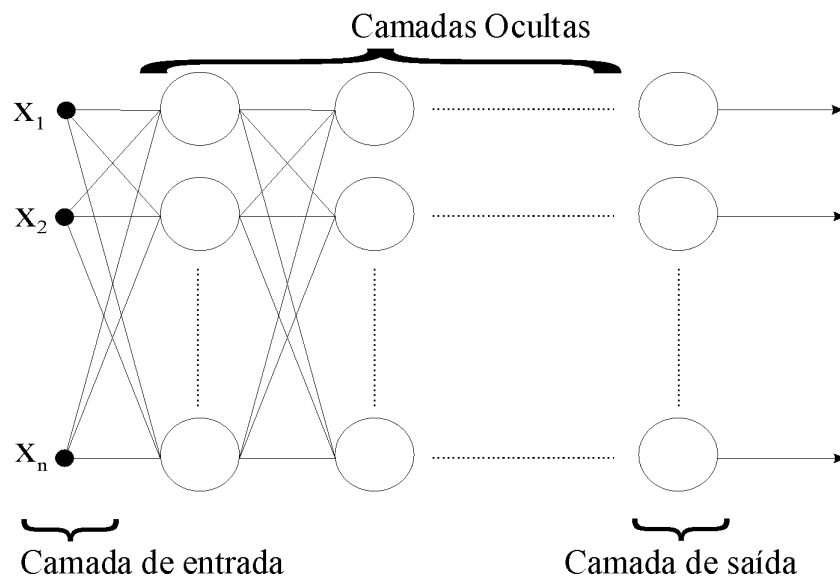


Figura 2.6 - RNA multicamada

Os nodos que compõe a rede neural artificial podem ter conexões do tipo:

- *feedforward* ou acíclicas (Figura 2.7) – a saída de um neurônio na i -ésima camada da rede não pode ser usada como entrada de nodos em camadas de índice menor ou igual a i [CAR 98]. Uma aplicação típica para as redes neurais artificiais *feedforward* é de desenvolver modelos não-lineares que também são usados para o reconhecimento e classificação de padrões. Uma rede *feedforward* pode ser vista como uma ferramenta que realiza a análise de regressão não linear [NAS 94].

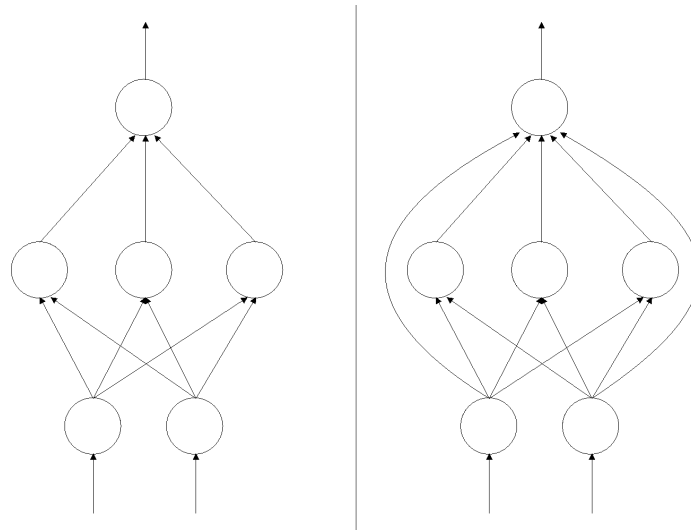


Figura 2.7 - RNA feedforward ou acíclica

- conexões feedback ou cíclica (Figura 2.8) – a saída de algum neurônio na i -ésima camada da rede é usada como entrada de nodos em camadas de índice menor ou igual a i . Se todas as ligações entre os neurônios forem cíclicas, a rede é chamada autoassociativa; estas redes associam um padrão de entrada com ele mesmo, e são particularmente úteis para a recuperação ou regeneração de um padrão de entrada [CAR 98].

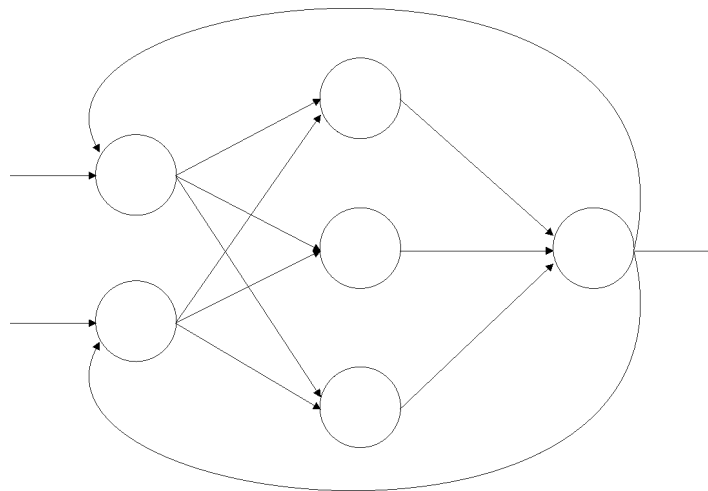


Figura 2.8 - RNA feedback ou cíclica

2.6 Aprendizado

Para o aprendizado das redes neurais, foram propostos diversos métodos de treinamento, sendo estes subdivididos em dois paradigmas principais: o aprendizado supervisionado e o não supervisionado. Para estes modelos existem vantagens e desvantagens que serão expostas a seguir. As RNA's possuem a capacidade de aprender por exemplos, determinando a intensidade de conexões entre os neurônios pertencentes à rede. Logo, um conjunto de procedimentos definidos para ajustar os parâmetros de uma RNA, a fim que a mesma possa aprender uma determinada função, é chamado de algoritmo de aprendizado. A designação de uma RNA, na resolução de um determinado problema, passa inicialmente por um

processo de aprendizagem, onde a rede procura extrair informações relevantes de padrões de informação apresentados a ela, modelando uma representação própria.

2.6.1 Supervisionado

A vasta maioria das redes neurais artificiais tem utilizado o treinamento supervisionado. Deste modo, a saída atual da rede neural é comparada com a saída desejada. Os pesos terão os seus valores iniciais setados aleatoriamente, e serão ajustados, através do algoritmo de aprendizagem, pela rede na próxima iteração ou ciclo.

O ajuste sináptico é dependente do valor esperado e do sinal atual de saída. Desta maneira, o método de aprendizado tenta minimizar o fluxo corrente de erros de todos os elementos em processamento. Esta redução global de erros trabalha modificando continuamente os pesos até que a rede alcance uma certa precisão.

Com o aprendizado supervisionado, as redes neurais artificiais devem ser treinadas antes de serem usadas. O treinamento consiste da apresentação dos sinais de entrada e saída à rede. Estes dados são freqüentemente referenciados ao conjunto de treinamento. A fase de treinamento pode consumir uma grande fatia de tempo. Em alguns sistemas protótipos, com um inadequado poder de processamento, o aprendizado pode levar semanas. O treinamento é considerado completo quando a rede neural alcança um certo nível de performance. Este nível significa que a rede alcançou uma precisão estatística conforme as produções de saída necessárias para uma dada seqüência de entradas. Quando não há mais a necessidade de aprendizado, os pesos são praticamente “congelados” para a aplicação. Alguns tipos de redes neurais permitem um treinamento contínuo, com uma taxa muito baixa de aprendizado, enquanto a mesma está em operação. Este processo ajuda a rede a adaptar-se gradualmente as condições de mudança.

O conjunto de treinamento precisa ser suficientemente grande para conter as informações necessárias para que a rede aprenda os moldes e as relações importantes. Se a rede é treinada somente com um exemplo em um determinado tempo, todos os pesos serão setados meticulosamente para este fato, os quais poderiam sofrer alterações drásticas no aprendizado de um próximo fato. Conforme um resultado, o sistema precisa aprender com todos os fatos em conjunto, provendo posteriormente o melhor ajuste dos pesos para todo o conjunto de fatos.

A maneira com que os sinais de entrada são representados, ou codificados, determina o maior componente constituinte para o sucesso de instrução da rede. Normalmente, as redes neurais artificiais somente manipulam, ou trabalham, com dados numéricos como entrada. Por este motivo, os dados do mundo exterior, devem ser tratados e convertidos para que se possa alimentar a rede. Esta captura de estímulos do mundo real pode ser feita através de vários tipos de dispositivos, tais como: câmeras de vídeo, diversos tipos de sensores, microfones, etc.

Várias técnicas de condicionamento já estão disponíveis para serem aplicadas a implementações de redes neurais artificiais, viabilizando e principalmente facilitando para que o desenvolvedor da rede encontre o melhor formato para os dados, e uma arquitetura adequada para a rede objetivando uma determinada aplicação.

Após o treinamento supervisionado, é importante analisar o que a rede pode realizar com os dados que ainda não foram apresentados à mesma. Se o resultado de saída do sistema não for razoável para este novo conjunto de dados (chamado conjunto de teste), presume-se que o treinamento da rede ainda não foi suficiente.

Esta avaliação é crítica para assegurar que a rede simplesmente não memorizou um dado conjunto de dados, mas sim aprendeu os modelos/padrões gerais envolvidos na aplicação (generalização). É importante ressaltar que às vezes o problema da generalização é devido à má qualidade dos dados usados para o treinamento e não um problema da rede.

2.6.2 Não supervisionado

O aprendizado não supervisionado é a grande promessa para o futuro, visto que implica que a rede aprenda se a necessidade de um conjunto de treinamento.

Estas redes não suportam influências externas para ajustar os seus pesos sinápticos, pois há um monitoramento de desempenho interno da mesma, analisando as regularidades e/ou tendências dos sinais de entrada, e conseqüentemente adaptando-se automaticamente as necessidades da rede.

Possuindo características de autonivelção, sem um suposto auxílio para determinar se o aprendizado converge ou não para o caminho certo, a rede possui mecanismos, mais precisamente, informações, de como se organizar. Esta propriedade e percepção da rede devem-se a topologia e as regras de aprendizado adotados pela rede neural artificial.

Uma rede com o algoritmo de aprendizado não supervisionado deve ter enfatizado a cooperação entre as camadas de unidades de processamento. A competição entre estas unidades é a base de aprendizado da rede. Normalmente, quando a competição pelo aprendizado ocorre de fato, somente os pesos pertencentes à unidade de processamento vencedora são ajustados.

2.6.3 Velocidade de aprendizado

A velocidade com que uma rede neural artificial aprende, depende de vários fatores. A baixa taxa de aprendizado resulta num tempo relativamente grande para a realização do aprendizado da rede, a fim de produzir um resultado adequado para o sistema em questão. Já com uma alta taxa de aprendizado, a rede pode não ser capaz de fazer uma possível discriminação fina em relação aos sistemas que aprendem de forma mais lenta.

Geralmente, vários fatores, além do tempo de aprendizado, precisam ser considerados quando se aborda a tarefa de aprendizado e treinamento da rede.

Alguns fatores que estão fortemente acoplados ao processo de aprendizado de uma RNA [DAC 92]:

- a complexidade da rede;
- o número de camadas (tamanho);
- o paradigma de seleção;
- a arquitetura adotada;
- algoritmo de aprendizado;
- as regras empregadas;
- a precisão desejada;

Todos estes fatores contribuem e alteram o tempo de treinamento da rede. A mudança de qualquer um destes fatores pode estender o tempo de treinamento para uma razão não muito significativa, ou resultando em uma precisão não satisfatória.

A maioria dos algoritmos de aprendizado possuem alguma provisão para a taxa de aprendizado ou em muitos casos, uma constante de tempo. Normalmente, este valor está compreendido num intervalo entre zero e um $[0, 1]$. Se a taxa de aprendizado exceder o valor máximo, o algoritmo de aprendizado irá corrigir os pesos da rede.

Pequenos valores da taxa de aprendizado não retificarão os erros tão rapidamente, mas se pequenos passos são tomados na correção de erros, há uma grande possibilidade de se alcançar uma boa convergência de aprendizado da rede.

2.6.4 Algoritmos de aprendizado

Muitas das leis de aprendizado estão em uso, e a maioria delas são apenas variações da mais difundida lei, que é a lei de Hebb. As pesquisas em torno das funções de aprendizado continuam, e busca-se aproximar cada vez mais estes modelos dos tão sonhados e perfeitos padrões biológicos.

Embora o homem esteja muito além de descobrir realmente como ocorre o processamento biológico, e o aprendizado seja algo extremamente complexo, simplificações e morfologias matemáticas continuam em desenvolvimento.

Abaixo, são apresentadas algumas das mais conhecidas e difundidas leis de aprendizado [DAC 92]:

- **Princípio de aprendizado de Hebb** – a primeira e indiscutivelmente a mais conhecida regra de aprendizado, foi apresentada pelo biólogo e psicólogo Donald Hebb. A descrição da mesma foi exposta em seu famoso livro *The Organization of Behavior* em 1949. A regra parte do seguinte pressuposto. Se um neurônio recebe uma entrada, proveniente de outro neurônio, e ambos estão ativos, isto é, possuem o mesmo sinal, os pesos entre os neurônios precisam ser excitados;
- **A Lei de Hopfield** – é praticamente similar ao princípio de aprendizado de Hebb com apenas uma exceção: a lei de Hopfield especifica a magnitude da excitação ou inibição. Se a saída desejada e o sinal de entrada estão ambos ativos ou inativos, os pesos são incrementados pela taxa de aprendizado, caso contrário, decrementados;
- **A Regra Delta de Widrow** – esta regra é uma variação um pouco além do princípio de aprendizado proposto por Hebb. A regra delta apresentada por Widrow é a mais comumente usada nos dias atuais. Esta regra se baseia na simples idéia da contínua modificação da intensidade e a importância das conexões de entrada; reduzindo consideravelmente a diferença entre o valor de saída desejado e o atual valor de saída da unidade de processamento, no caso, o neurônio artificial. A regra delta altera os pesos sinápticos de modo que minimize o erro quadrático da rede, trabalhando da seguinte forma: o erro calculado na saída é transformado pela derivação da função de transferência e conseqüentemente usado para ajustar os pesos de entrada da camada prévia da rede, ou seja, o erro é retro-propagado às camadas anteriores, sendo uma de cada vez. O processo de retro-propagação dos erros da rede continua até que a primeira camada da rede seja alcançada. Os tipos de redes chamadas de *feedforward* e *backpropagation* derivam seus nomes dos métodos adotados pelo processamento do erro. Quando se usa a regra delta, é importante assegurar que o conjunto de dados de entrada está disposto de forma aleatória ou gerado de forma randômica. Uma vez mal organizado, este conjunto de treinamento pode conduzir a rede a não convergência da precisão desejada, impossibilitando o aprendizado do problema em questão.
- **A Lei de aprendizado de Teuvo Kohonen** – desenvolvida por Teuvo Kohonen, a mesma foi inspirada nos sistemas biológicos, onde os elementos competem entre si por uma oportunidade de aprender, ou atualizar/ajustar seus respectivos pesos. A unidade de processamento que possuir o melhor sinal de saída será considerada o mais apto, e conseqüentemente passarão a ter a capacidade e privilégio de inibir os ajustes sinápticos de seus concorrentes e excitar seus vizinhos. Somente a unidade apta e seus respectivos vizinhos terão permissão para ajustar seus pesos. A abrangência e a possível área que uma unidade vizinha pertence está relacionada ao período de treinamento da rede. O paradigma atual usa o seguinte procedimento: é formada uma grande área de vizinhança e a medida com que ocorre o processo de treinamento, há uma seleção e conseqüentemente um estreitamento da mesma.

2.7 Redes Perceptron

As redes neurais artificiais com função de ativação foram inicialmente estudadas por Rosenblatt em meados de 1958, as quais foram chamadas por ele de Perceptrons. O entusiasmo de Rosenblatt levou-o a construir suas redes em hardware, inclusive usando um algoritmo de aprendizado.

Estas redes foram aplicadas para a classificação de problemas que geralmente possuíam como fonte de alimentação imagens binárias de caracteres ou simplesmente moldes de informações [BIS 95]. O perceptron em sua origem era uma simulação computacional da retina, a qual demonstrou como o sistema nervoso visual reconhece padrões [TAF 96].

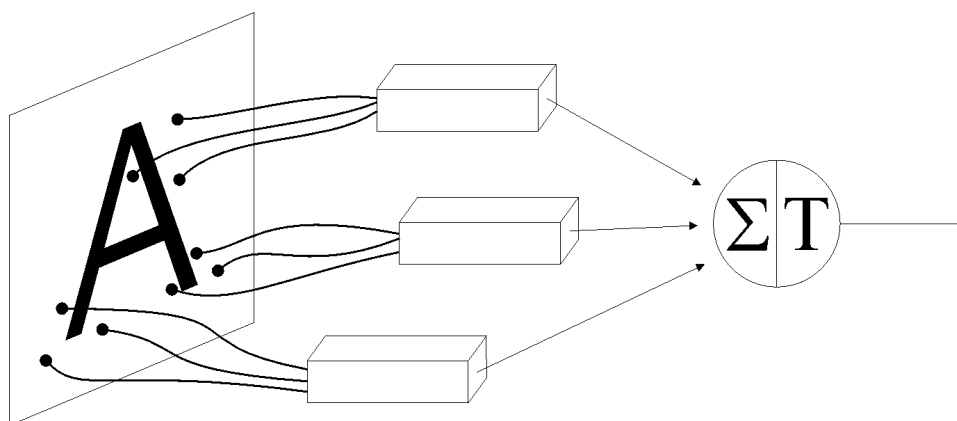


Figura 2.9 - O perceptron elementar de Rosenblatt [BIS 95]

Enquanto Rosenblatt estava desenvolvendo o perceptron, Widrow e seus colegas, estavam trabalhando em uma linha de pesquisa similar a de Rosenblatt; mais conhecida como ADALINE. Como já exposto, o termo **AD**aptive **LI**near **E**lement refere-se a uma única unidade de processamento com um limiar não linear.

Como as redes neurais artificiais de uma única camada possuem uma certa limitação, Rosenblatt resolveu então usar um número fixo de neurônios para transformar e tratar os dados provindos do mundo exterior. Estas unidades de processamento podem ser chamadas de função base de um discriminador limiar [BIS 95].

Rosenblatt propunha resolver problemas como a implementação das funções booleanas **E** e **OU** de duas variáveis, sendo que a escolha dos ganhos para este caso parecia ser trivial. Entretanto, para a implementação de uma função discriminatória arbitrária, a escolha não é tão simples e muito menos trivial, e dependendo do número de variáveis envolvidas, sem a existência de algum método, beira o impossível [KOV 96].

Inspirado também pelas idéias de McCulloch, Rosenblatt compôs a rede perceptron por uma camada de entrada, onde cada elemento pertencente à camada de entrada fazia a distribuição do sinal que ele recebia para todas as unidades de processamento. Os neurônios eram essencialmente compostos por unidades sigma e de funções de transferência, sendo que estas, eram responsáveis pela soma ponderada dos sinais oriundos das conexões com os dados de entrada. Foi adicionada a camada de entrada um elemento especial chamado viés, o qual possui um sinal de valor sempre um. A conexão entre o viés e a unidade sigma tem peso w_0 , que por sua vez é ajustado da mesma maneira com que os demais pesos o são.

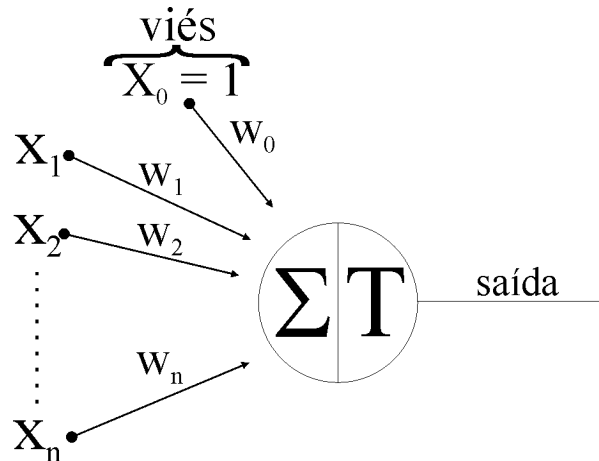


Figura 2.10 - A unidade de processamento do perceptron

O nível de ativação de uma rede perceptron é dado pela soma ponderada dos pesos sinápticos com os valores de entrada, $\sum x_i \cdot w_i$.

Estas redes usam uma função de transferência do tipo *hard-limiter* (limite ríspido), onde a ativação do limiar resulta num valor de saída 1, ou, -1 caso contrário. Dados os valores de entrada x_i , os pesos w_i , e um limiar t , o perceptron computa os valores de saída da seguinte maneira:

$$\begin{cases} \mathbf{1} & \text{se } \sum x_i w_i \geq t \\ -\mathbf{1} & \text{se } \sum x_i w_i < t \end{cases} \quad (2.1)$$

As redes perceptron usam como configuração, o treinamento supervisionado. O perceptron altera os seus pesos, visando reduzir o erro.

2.8 A lei de aprendizado do perceptron

Sendo o parâmetro c uma referência a taxa de aprendizado na medida em que reflete a taxa com que os ganhos são alterados em função dos erros, e d sendo o valor de saída esperado, o ajuste sináptico dos pesos no i -ésimo componente do vetor de entrada, Δw_i , é dado por:

$$\Delta w_i = c \cdot (d - \text{sign}(\sum x_i w_i)) \cdot x_i \quad (2.2)$$

A expressão $\text{sign}(\sum x_i w_i)$ é o valor de saída do perceptron, que pode assumir os valores $+1$ ou -1 . A diferença entre o valor desejado d e a saída atual, poderá ser 0 , 2 ou -2 . Logo, para cada componente do vetor de entrada, conclui-se que [LUG 98]:

- se, a saída desejada e a atual possuírem valores iguais, não haverá alteração alguma no peso;
- se o valor atual da saída for -1 e a saída desejada for 1 , o peso deverá ser ajustado na i -ésima linha da rede por $2 \cdot c \cdot x_i$;
- se o valor atual de saída for 1 e a saída desejada for -1 , o peso deverá ser ajustado por $-2 \cdot c \cdot x_i$.

Logo,

$$\begin{cases} \Delta w_i = 0, & \text{para } y_l^d = y_1 \\ \Delta w_i = 2 \cdot c \cdot x_i, & \text{para } y_l^d \neq y_1 \end{cases} \quad (2.3)$$

$$\Downarrow$$

$$\Delta w_i = c \cdot (1 - y_l^d \cdot y_1) \cdot y_l^d \cdot x_{i,l}$$

O procedimento mostrado acima tem como principal função, ajustar o conjunto de pesos da rede, a fim de minimizar o erro médio do conjunto de treinamento.

Como os perceptrons são utilizados em problemas de classificação, eles possuem a capacidade de aprender e classificar os dados de entrada em grupos ou classes.

Supondo-se uma rede perceptron, teremos a seguinte unidade sigma (Σ):

$$\Sigma = x_0 \cdot w_0 + x_1 \cdot w_1 + x_2 \cdot w_2 + \dots + x_n \cdot w_n \quad (2.4)$$

onde, o vetor $[x_0, x_1, x_2, \dots, x_n]$ são os sinais de entrada, e o vetor $[w_0, w_1, w_2, \dots, w_n]$ são os pesos respectivos ao vetor de entrada.

Teremos então:

$$\Sigma = \sum_{i=0}^n x_i \cdot w_i \quad (2.5)$$

Se a camada de entrada possuir dois elementos, x_1 e x_2 , a unidade sigma da rede será representada por:

$$\Sigma = x_0 \cdot w_0 + x_1 \cdot w_1 + x_2 \cdot w_2 \quad (2.6)$$

Como o viés, é representado por $x_0 = 1$, teremos:

$$\Sigma = w_0 + x_1 \cdot w_1 + x_2 \cdot w_2 \quad (2.7)$$

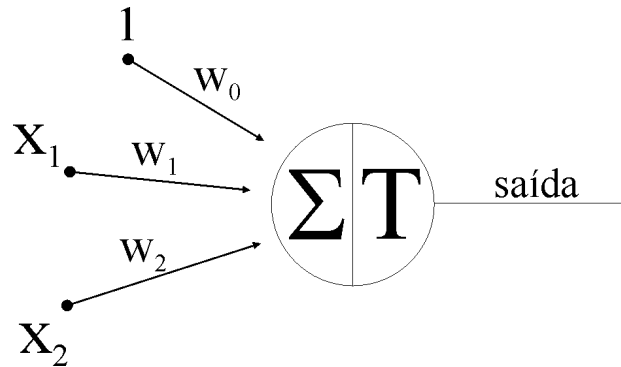


Figura 2.11 - Perceptron de duas entradas

Na equação (2.1), usando o limiar $t = 0$, tem-se que a saída é 1 se, $\sum_{i=0}^n x_i w_i \leq 0$, caso contrário, teremos -1.

A separação entre duas classes é chamada de superfície de decisão. Como só existem duas entradas, a superfície de decisão é uma reta. Se existem mais do que duas entradas, haverá então um hiperplano [TAF 96].

As redes perceptrons foram inicialmente elogiadas, entretanto, Nils Nilson em 1965, Minsky e outros pesquisadores, analisaram as limitações das redes perceptrons. Eles demonstraram que os perceptrons não podiam resolver uma certa classe de problemas, normalmente chamados de problemas linearmente não separáveis [LUG 98].

Estes problemas foram solucionados a partir da década de 80, onde houve o surgimento de outras técnicas de aprendizado, como por exemplo, o *backpropagation*.

2.9 Limitações: O problema do OU-EXCLUSIVO

Um dos problemas que o perceptron não seria capaz de resolver era o do ou-exclusivo. Foi baseado neste exemplo que Minsky e Papert mostraram à comunidade científica que o modelo de Rosenblatt não era tão eficiente e promissor.

Tabela 2.1 - Tabela verdade do ou-exclusivo

x_1	x_2	saída
1	1	0
1	0	1
0	1	1
0	0	0

Considerando uma rede perceptron com duas entradas $[x_1, x_2]$, dois pesos $[w_1, w_2]$, e um limiar t , a rede, para aprender com estes fatos, deveria encontrar os pesos designados para satisfazer a tabela verdade e as seguintes premissas [LUG 98]:

- para a linha 1 da tabela verdade: $w_1 \cdot 1 + w_2 \cdot 1 < t$

- para a linha 2 da tabela verdade: $w_1 \cdot 1 + 0 > t$
- para a linha 3 da tabela verdade: $0 + w_2 \cdot 1 > t$
- para a linha 4 da tabela verdade: $0 + 0 < t$

As premissas apresentadas, baseadas nos pesos $[w_1, w_2]$ e no limiar t , não possuem solução. Logo, o perceptron de uma única camada é incapaz de resolver este tipo de problema.

O motivo pelo qual torna o problema do ou-exclusivo impossível para as redes do tipo perceptron é que as duas classes que precisam ser distinguidas não são linearmente separáveis.

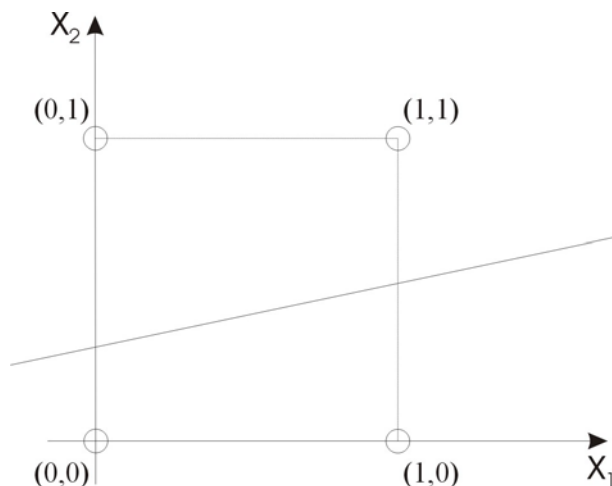


Figura 2.12 - Plano que representa as combinações possíveis do XOR

Percebe-se que é impossível plotar uma linha reta que separe em duas dimensões os pontos $\{(0,0), (1,1)\}$ de $\{(0,1), (1,0)\}$.

Cada parâmetro dos dados de entrada corresponde a uma dimensão, com cada valor de entrada definindo um ponto no espaço [LUG 98].

2.10 Redes Multilayer Perceptron

Os problemas não linearmente separáveis podem ser resolvidos através das redes com uma ou mais camadas intermediárias. A alteração da arquitetura da rede, como a inserção de camadas ocultas e/ou o número de neurônios, a princípio, não parece ser problema, pois um dos principais agravantes passa a ser o algoritmo de treinamento para as redes multicamadas. Fator este que, devido à inexistência ou desconhecimento, causou uma atenuação nas pesquisas em redes neurais artificiais em meados da década de 70. Uma das alternativas adotadas é dividir a rede em um conjunto de subredes, sendo uma subrede para cada camada, com um treinamento independente. Este método de subdivisão, muitas vezes, ou não é possível ou é muito complicado. Outra possibilidade seria realizar um treinamento completo, isto é, de uma só vez. O problema encontrado para este segundo método está em como realizar o treinamento dos nodos que pertencem à camada intermediária, visto que é extremamente complicado determinar que tipo de resposta desejada estes teriam, ou seja, como determinar o erro. A aplicabilidade deste método está restrita a definição do erro nos nodos pertencentes às camadas intermediárias da rede. Se for utilizada uma função do tipo limiar, a avaliação do erro será complexa, visto que, os nodos das camadas intermediárias e de saída não terão como saber a margem de erro ou a diferença entre as respostas de seus nodos com relação às respostas desejadas. Uma das soluções para o problema apresentado seria a utilização de uma função de ativação não linear, a qual resolve o mesmo em parte, visto que a utilização deste tipo de função em redes multicamada resultaria na equivalência de uma rede de uma única camada [CAR 98].

Adotou-se então treinar as redes com mais de uma camada através de métodos baseados no gradiente descendente. Métodos baseados no gradiente descendente precisam ter a função de ativação contínua, diferenciável e não decrescente. A função adotada precisa informar os erros que a rede cometeu para as camadas anteriores, com uma boa precisão. Logo a função que mais se adapta a estas características é a função do tipo sigmóide [CAR 98].

O processamento atribuído a cada neurônio pertencente à rede é resultante da combinação do processamento realizado pelos neurônios da camada anterior, que por sua vez estão atribuídos a este nodo da próxima camada. A medida com que cada camada intermediária da rede se aproxima da camada de saída há uma delimitação do espaço de decisão dos dados que está recebendo. Para uma rede com duas camadas intermediárias, teríamos a primeira camada oculta, delimitando o espaço de padrões de treinamento através das “retas traçadas” pelos neurônios. A segunda camada forma regiões convexas, onde o número de lados que compõe tal região é determinado pela quantidade de unidades conectadas a este neurônio, que por sua vez combina as retas que surgiram da camada anterior. Cada neurônio da camada de saída forma regiões, provenientes das combinações das regiões convexas [CAR 98]. Conclui-se que cada neurônio que compõe uma rede Multilayer Perceptron contribui para a detecção de características dos dados apresentados.

A determinação do número de camadas a ser utilizada influi de forma crucial no aprendizado da rede. O uso de um grande número de camadas intermediárias não é recomendado, visto que o erro ocorrido em uma camada é propagado a camadas anteriores da rede. A determinação do número de neurônios que pertence a camadas intermediárias é definida de forma empírica, e normalmente depende da distribuição dos padrões de treinamento e validação da rede. Um uso excessivo de neurônios levará a rede a decorar o conjunto de treinamento, ao invés de extrair as características gerais (generalizar). Ao processo de memorização do conjunto de treinamento, dá-se o nome de *overfitting*. Um número razoavelmente pequeno de neurônios levará a rede a aumentar o tempo de treinamento, dificultando a determinação da representação ótima do problema proposto. Neste caso, alguns neurônios poderão ficar sobrecarregados, pois estes precisam lidar com um número elevado de restrições a serem analisadas.

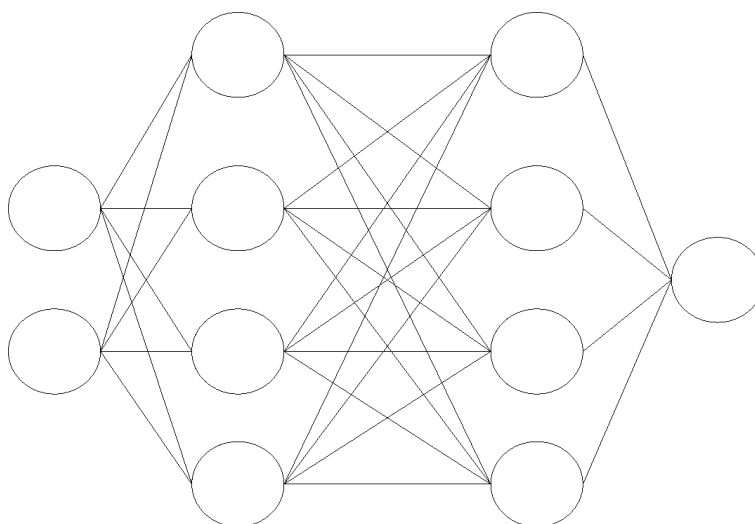


Figura 2.13 - Uma rede MLP

2.11 Algoritmo de treinamento das redes MLP

O algoritmo de aprendizado mais conhecido para a realização do treinamento das redes neurais multicamadas é o *backpropagation*. Cabe ressaltar que existem vários algoritmos de

aprendizado para as redes MLP, e estes normalmente possuem um aprendizado supervisionado. Pode-se ainda classificar os mesmos em dois grupos:

- estáticos;
- dinâmicos.

Os algoritmos de aprendizado estáticos não alteram a estrutura da rede, mudando somente o valor dos pesos sinápticos. Já os dinâmicos, podem mudar a arquitetura da rede, isto é, podem alterar o número de camadas, o número de neurônios da camada intermediária e o número de conexões da rede.

O método de aprendizado *backpropagation* foi descoberto através de inúmeras linhas pesquisas. Werbos, em 1974, foi um dos primeiros a propor o uso deste método de aprendizado na Universidade de Harvard em sua teste de doutorado "*Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*". Parker em 1985 redescobriu as técnicas utilizadas por Werbos no relatório do MIT, "Learning Logic". Até então, as pesquisas e principalmente os créditos eram dados a Rumelhart e aos outros membros do Grupo de Processamento Paralelo e Distribuído, por popularizar e desenvolver procedimentos que realmente pudessem ser utilizados. Este grupo publicou dois volumes que relatavam vários tipos de arquiteturas de redes neurais artificiais, incluindo um tratamento considerável sobre o procedimento de treinamento da regra delta generalizada, o *backpropagation* propriamente dito e alguns tópicos relacionados as RNA's.

O método de aprendizado *backpropagation* pode ser aplicado a qualquer rede que usufrui de uma função de ativação diferencial e aprendizado supervisionado. Assim como a regra delta, sua otimização é baseada no gradiente descendente, que ajusta os pesos para reduzir o erro da rede. O nome *backpropagation* surgiu do método na qual as correções da rede são realizadas nos pesos das conexões.

Durante a fase de treinamento, os sinais/padrões de entrada são apresentados a rede artificial em uma determinada ordem. Cada padrão de treinamento é propagado adiante, camada após camada, até a produção do sinal/padrão de saída. A saída computada pela rede é então comparada com uma saída desejada. Esta comparação irá gerar um valor que determinará o erro. Este erro será utilizado como uma realimentação para as conexões, que resultará no ajuste dos pesos sinápticos de cada camada num sentido oposto à propagação dos sinais de treinamento. Os acoplamentos retrógrados somente existirão na fase de treinamento, considerando que as conexões adiante (sentido entrada → saída) serão usadas durante a fase de treinamento e uso da rede.

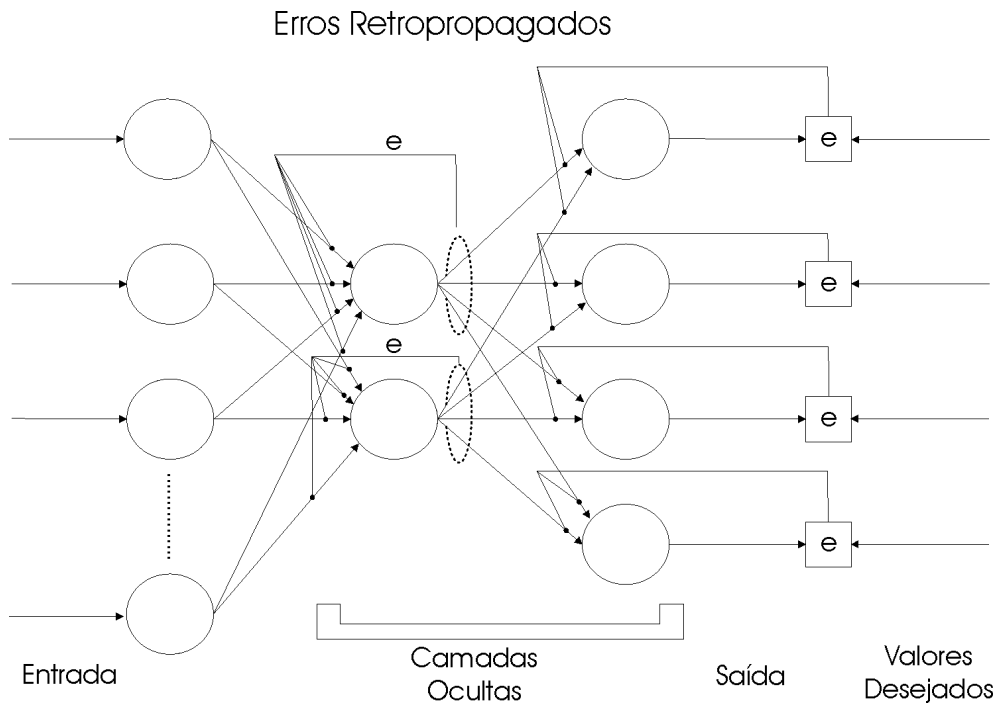


Figura 2.14 - Rede MLP com os acoplamentos retrógrados para os ajustes sinápticos [PAT 95]

Fazendo o uso do *backpropagation*, as camadas ocultas terão os seus pesos ajustados de acordo com as camadas subseqüentes, i.e, com as camadas seguintes. Deste modo, os erros computados na camada de saída serão usados para ajustar os pesos entre a última camada escondida ou oculta com a camada de saída. Assim, o erro calculado de uma camada escondida será usado para ajustar os pesos da camada oculta anterior. Este processo será repetido até que a primeira camada oculta seja ajustada. Desta forma, os erros serão retropropagados camada a camada com as devidas correções. Rotina esta que será realizada de uma maneira repetitiva, ajustando os pesos das respectivas camadas.

O processo é repetido por um número determinado de vezes para cada padrão de dados durante o treinamento até que o erro total da saída tenha convergido a um valor mínimo, ou até que algum limite predeterminado de iterações tenha sido completado.

Pode-se então criar duas fases para o algoritmo de treinamento do *backpropagation*. Cada fase percorre um sentido da rede. A primeira fase, chamada de *forward*, define a saída da rede para um determinado padrão de dados de entrada. A segunda e última está incumbida de utilizar a saída desejada/esperada e a saída fornecida pela última camada da rede para ajustar os pesos sinápticos da rede neural.

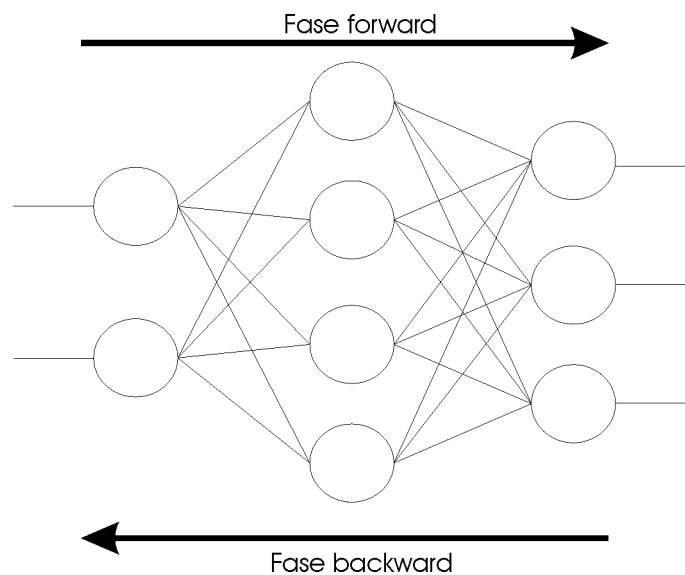


Figura 2.15 - Fluxo de treinamento de uma MLP com *backpropagation*

Segundo Carvalho [CAR 98], podemos definir os passos da seguinte maneira:

Fase forward

1. Os padrões de saída são apresentados a primeira camada c_1 que compõe a rede;
2. Para cada camada c_i a partir da camada de entrada:
 - 2.1. Os sinais de saída do neurônio da camada c_i irão alimentar a entrada da camada c_{i+1} , i.e., serão os sinais de entrada para a próxima camada;
3. Os sinais provenientes da última camada serão comparados com os sinais desejados;

Fase backward

1. Da última camada da rede até a primeira [$n \rightarrow 1$];
 - 1.1. Os neurônios artificiais da camada c_a (camada atual) devem ajustar seus pesos sinápticos de forma a reduzir seus erros;
 - 1.2. O erro de um neurônio das camadas intermediárias, $c_{[2, n-1]}$, por exemplo, c_i , será calculado utilizando os erros retropropagados dos neurônios que pertencem às camadas seguintes conectados a ele, no caso, c_{i+1} , os quais serão ponderados pelos pesos das conexões entre eles;

O backpropagation

1. Inicializar a rede, ou seja, pesos e parâmetros;
2. Repita
 - 2.1. Para cada padrão/dados de treinamento $P[x]$, para todo $x \in [1, n]$, sendo n o número total de amostras do conjunto de treinamento.
 - 2.1.1. Calcular a saída (S) da rede utilizando a fase *forward*;
 - 2.1.2. Comparar a saída (S), calculada no item 2.1.1, com as saídas desejadas;
 - 2.1.3. Realizar atualização dos pesos sinápticos fazendo o uso da fase *backward*;
3. Até o erro ser mínimo ou até x ciclos;

O algoritmo *backpropagation* também é chamado de regra delta generalizada, pois é baseado na regra delta apresentada por Widrow e Hoff.

3 REFERÊNCIAS BIBLIOGRÁFICAS

- [BAR 98] BARROS, Carlos, PAULINO, Wilson. **O Corpo Humano**. Editora Ática, 1998.
- [BAR99] BARONE, Dante Augusto Couto, “**Projeto Revox**”, versão eletrônica em <http://www.ucs.tche.br/revox>, 1999.
- [BIS 95] BISHOP, C. M. **Neural Networks for Pattern Recognition**. Oxford University Press, 1995.
- [CAR 97] CARREIRA-PERPIÑÁN, M. Á. **A Review of Dimension Reduction Techniques**, Technical Report CS-96-09, 1997. <http://www.dcs.shef.ac.uk/~miguel/papers/ps/cs-96-09.ps.gz> (09 de Dezembro de 2000, 01:00).
- [CAR 98] CARVALHO, André, LUDEMIR, Antônio. **Fundamentos de Redes Neurais Artificiais: 11ª Escola de Computação**. Imprinta Gráfica e Editora Ltda, 1998.
- [CHU 92] Churchland, P.S. , Sjenowski, T.J. 1992. **The computacional Brain**. Cambridge, Mass,: MIT Press.
- [DAC 92] Data & Analysis Center for Software. **Artificial Neural Networks Technology**. Disponível por WWW em <http://www.dacs.dtic.mil> (20/09/2000).
- [HAY 94] HAYKIN, Simon. **Neural Networks: A comprehensive Foundation**. New York: Macmillan College Publish Company, 1994.
- [HAR 35] Harrison, R.G. 1935. **On the origin and development of the nervous system studied by methods of experimental embryology**. *Proc.R. Soc. Lon. [Biol.]*, v. **118**, p. 155 – 196.
- [HAR 87] Harrington, A. 1987. **Medicine, Mind, and the Doble Brain: A Study in nineteenth-Century Thought** . Princeton, N.J.: Princeton University Press.
- [KAN 98] Kandell, E.R., Jessel, T.M., Schwartz, J.H. 1998. **Neurociencia y conducta**. Madrid: Prentice Hall, 812 p. Il.
- [KAS 96] KASABOV, Nikola K. **Foundations of Neural Networks, Fuzzy Systems, and Knowledge Engineering**. The MIT Press, 1996.
- [KOV 96] KOVÁCS, Zsolt L. **Redes Neurais Artificiais: Fundamentos e Aplicações**. Segunda Edição, Collegium Cognitio, 1996.
- [LUG 98] LUGER, G. F., Stubblefield W. A., **Artificial Intelligence**. Addison Weslwy, 1998.
- [MINa 00] MINELLO, Luiz Fernando. **A voz Humana**. Disponível por E-mail em dbclfm@unileon.es (16 Dez. 2000).

- [MINb 00] MINELLO, Luiz Fernando. **O Sistema Nervoso**. Disponível por E-mail em dbclfm@unileon.es (17 Dez. 2000).
- [PAT 95] PATTERSON, Dan W. **Artificial Neural Networks: Theory and Applications**. Prentice Hall, 1995.
- [TAFa 96] TAFNER, Malcon A. **Reconhecimento de palavras isoladas usando redes neurais artificiais**. Dissertação de Mestrado. Programa de Pós-Graduação em Engenharia de Produção. Universidade Federal de Santa Catarina. Florianópolis, 1996.
- [TAFb 96] TAFNER, Malcon, XEREZ, Marcos, Rodrigues, Ilson. **Redes Neurais Artificiais: Introdução e Princípios de Neurocomputação**. EKO, 1996.
- [WIN 93] WINSTON, P. H. **Artificial Intelligence**. Third Edition. Addison-Wesley, 1993.
- [ZAD 65] ZADEH, L.A. **Fuzzy sets**. Information and Control, Vol. 8, 1965, pp. 338-353.