

SCC0276 - Aprendizado de Máquin

Aula - *Multi-Layer Perceptron* (MLP)

Profa. Dra. Roseli Aparecida Francelin Romero
SCC - ICMC - USP

2019

Sumário

- 1 **Introdução**
 - Modelo de rede MLP
- 2 **Treinamento de redes MLP**
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 **Funções Utilizadas**
 - Função Softmax

Perceptron multicamadas

- Redes de apenas uma camada representam somente funções linearmente separáveis.
- Redes de múltiplas camadas solucionam essa restrição.
- O desenvolvimento do algoritmo *backpropagation* foi um dos motivos para o ressurgimento da área de redes neurais [Rumelhart *et. al*, 1986].

Sumário

- 1 **Introdução**
 - Modelo de rede MLP
- 2 **Treinamento de redes MLP**
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 **Funções Utilizadas**
 - Função Softmax

Modelo de rede neural com múltiplas camadas.

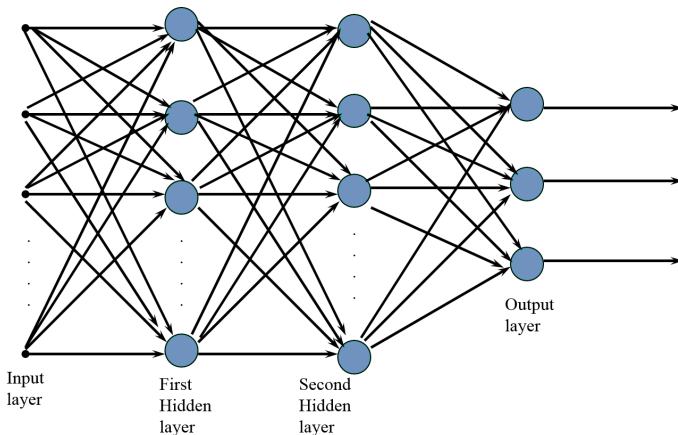


Figura 1: Rede neural *feed-forward* com múltiplas camadas.

Sumário

- 1 Introdução
 - Modelo de rede MLP
- 2 **Treinamento de redes MLP**
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 Funções Utilizadas
 - Função Softmax

- Modelo de rede MLP

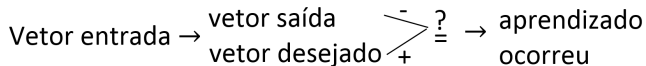
- O algoritmo Backpropagation

- Processo de aprendizado
- Termo Momentum
- Velocidade de aprendizado
- Modos de treinamento
- Generalização

- Função Softmax

Aprendizado da rede

- O esquema de aprendizado da rede pode ser descrito do seguinte modo:



Aprendizado da rede

Pesos (*Gradient Descent Method*)

$$w_{ji}(k+1) = w_{ji}(k) - \eta \frac{\partial E_p(w)}{\partial w_{ji}} \Big|_{w(k)}$$

Onde η é uma constante positiva (velocidade de aprendizado).

- Calculando a derivada parcial do E_p , tem-se:

$$\frac{\partial E_p}{\partial w_{ji}} = \frac{\partial E_p}{\partial y_{pj}} \cdot \frac{\partial y_{pj}}{\partial v_{pj}} \cdot \frac{\partial v_{pj}}{\partial w_{ji}}$$

- Para se calcular $\frac{\partial E_p}{\partial y_{pi}}$, dois casos devem ser considerados:

11 / 45

- **Observação:** os erros são computados no sentido *backward*. O erro foi chamado de *back-propagado* → algoritmo de aprendizado **backpropagation** (BP).

Algoritmo *Backpropagation*

- **Inicialização:** pesos iniciados com valores aleatórios e pequenos ($[-1, +1]$).
- **Treinamento - Repita:**
 - Considere um novo padrão de entrada x_i e seu respectivo vetor de saída t_i desejado do conjunto de treinamento.
 - **Repita:**
 - Apresentar o par (x_i, t_i) . **(modo padrão)**
 - Calcular as saídas dos processadores, começando da primeira camada escondida até a camada de saída.
 - Calcular o erro na camada de saída.
 - Atualizar os pesos de cada processador, começando pela camada de saída, até a camada de entrada.
 - **Até que o erro quadrático médio para esse padrão seja $\leq tol1$.**
- **Até que o erro quadrático médio seja $\leq tol2$ para todos os padrões do conjunto de treinamento.**

Processo de aprendizado

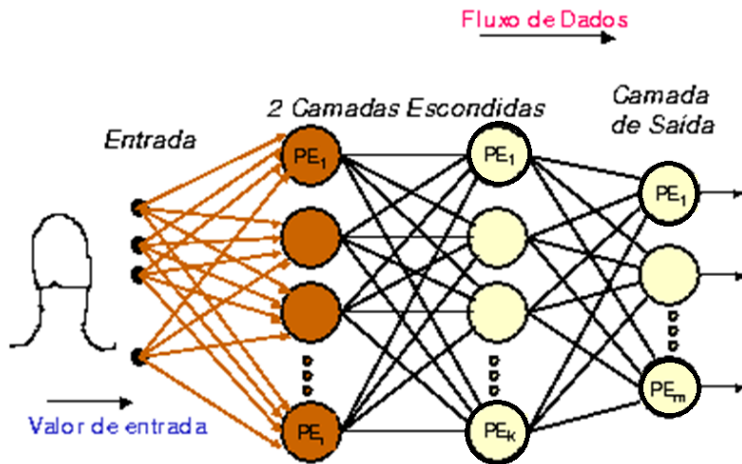


Figura 2: *Feed-forward* (fase 1), primeira camada escondida.

Processo de aprendizado

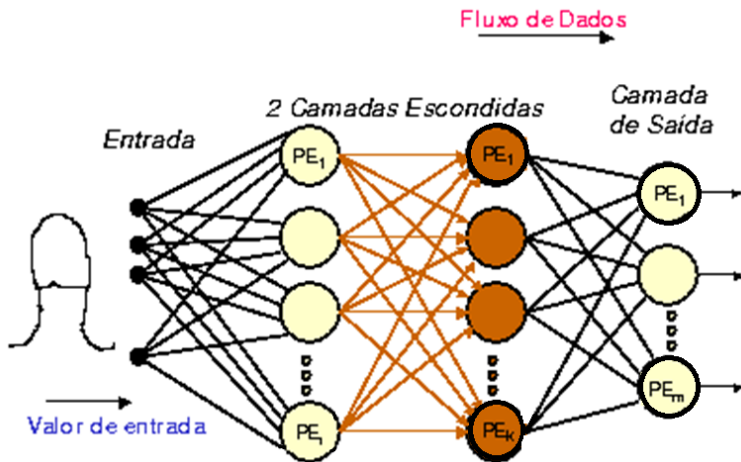


Figura 3: *Feed-forward* (fase 1), segunda camada escondida.

Processo de aprendizado

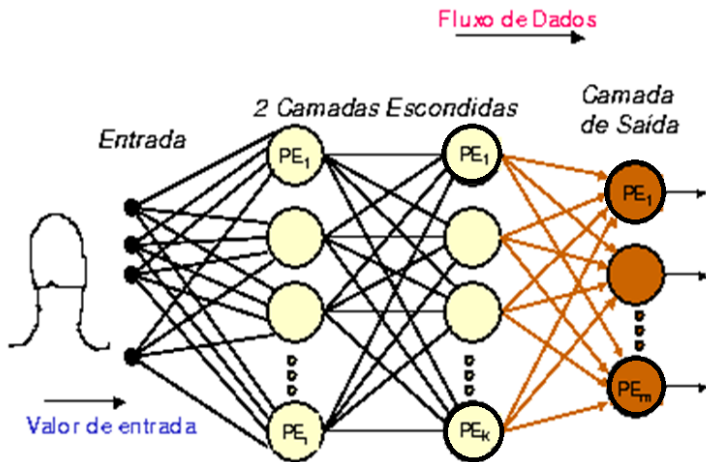


Figura 4: *Feed-forward* (fase 1), camada de saída.

Processo de aprendizado

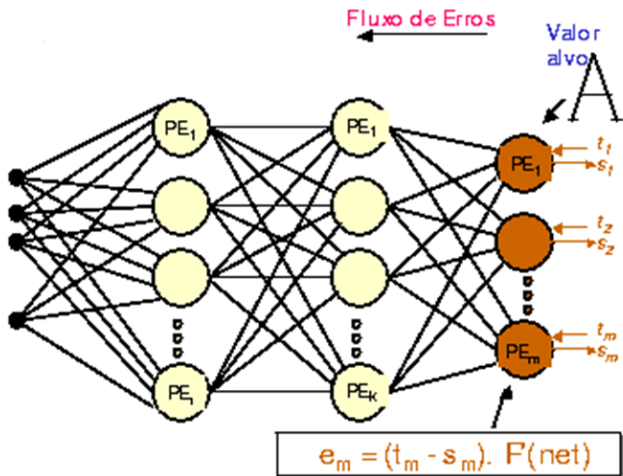


Figura 5: *Feed-backward* (fase 2), cálculo do erro da camada de saída.

Processo de aprendizado

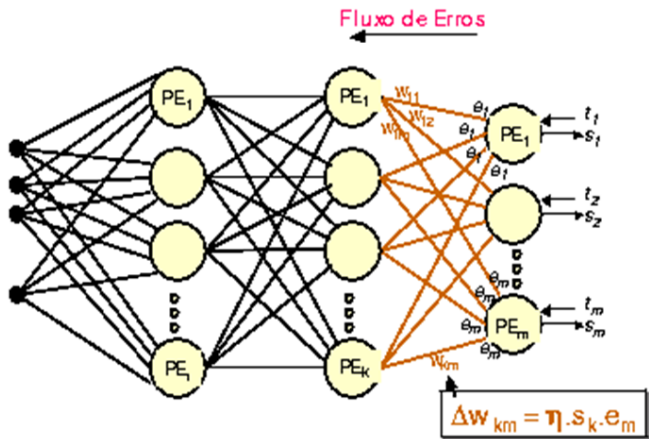


Figura 6: *Feed-backward* (fase 2), atualização dos pesos da camada de saída.

Processo de aprendizado

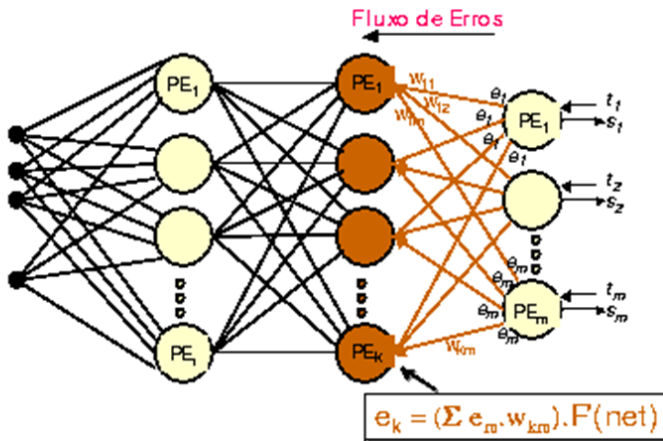


Figura 7: *Feed-backward* (fase 2), cálculo do erro da segunda camada escondida.

Processo de aprendizado

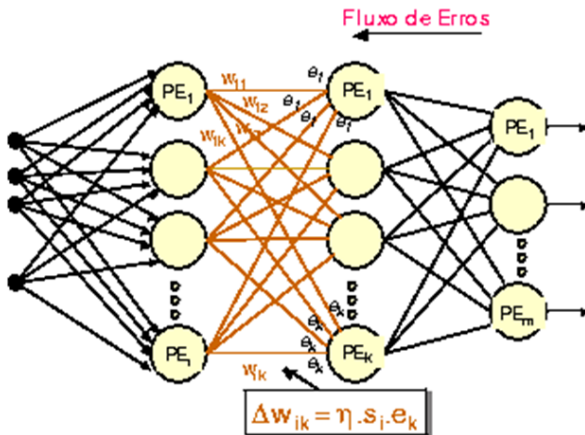


Figura 8: *Feed-backward* (fase 2), atualização dos pesos da segunda camada escondida.

Processo de aprendizado

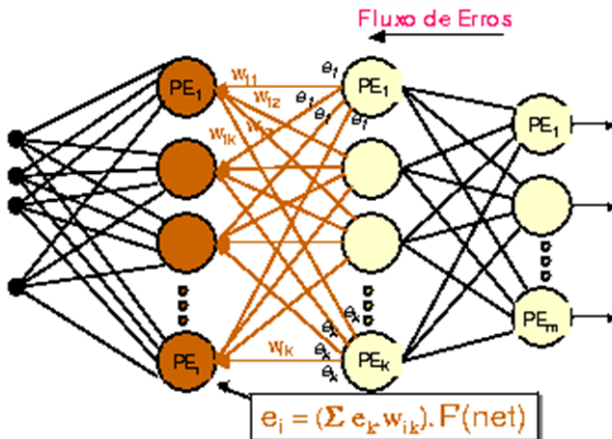


Figura 9: *Feed-backward* (fase 2), cálculo do erro da primeira camada escondida.

Processo de aprendizado

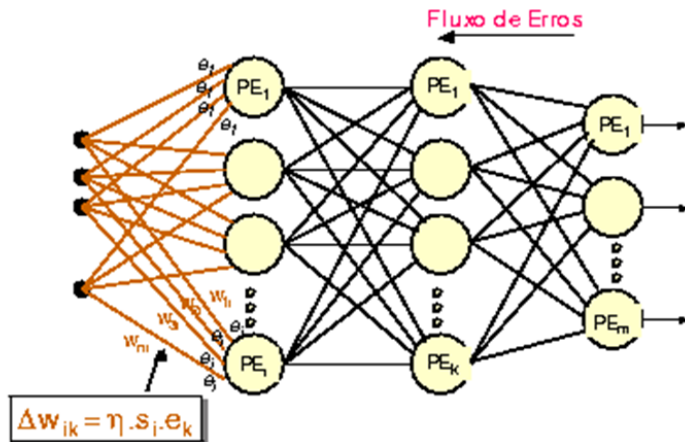


Figura 10: *Feed-backward* (fase 2), atualização dos pesos da primeira camada escondida.

Processo de aprendizado

- Este procedimento de aprendizado é repetido diversas vezes, até que, **para todos processadores de camada de saída e para todos padrões de treinamento**, o erro seja menor do que o especificado.

Sumário

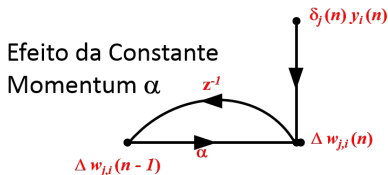
- 1 Introdução
 - Modelo de rede MLP
- 2 Treinamento de redes MLP
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - **Termo Momentum**
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 Funções Utilizadas
 - Função Softmax

Efeito da constante α

- É um método simples de aumentar a velocidade do aprendizado e evitar o perigo de instabilidade, como mostrado por Rumelhart *et al.*, 1986.

$$\Delta w_{ji}(n) = \eta \delta_j(n) y_i(n) + \alpha \Delta w_{ji}(n-1) \quad (1)$$

- Onde α é geralmente um número positivo chamado **constante momentum**.



A equação (α) é chamada
REGRA DELTA
GENERALIZADA. Se $\alpha = 0$
 \Rightarrow REGRA DELTA

Sumário

- 1 **Introdução**
 - Modelo de rede MLP
- 2 **Treinamento de redes MLP**
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - **Velocidade de aprendizado**
 - Modos de treinamento
 - Generalização
- 3 **Funções Utilizadas**
 - Função Softmax

Velocidade de aprendizado

- O algoritmo BP fornece uma aproximação para a trajetória no espaço dos pesos.
- Quanto menor o valor de η , menores as mudanças nos pesos e mais suave será a trajetória.
 - **Aprendizado lento.**
- Se η é muito grande, o aprendizado torna-se rápido, porém a rede pode tornar-se **instável**.

Sumário

- 1 **Introdução**
 - Modelo de rede MLP
- 2 **Treinamento de redes MLP**
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - **Modos de treinamento**
 - Generalização
- 3 **Funções Utilizadas**
 - Função Softmax

Modos de treinamento

- Aprendizado BP resulta de muitas apresentações de um conjunto de treinamento de exemplos.
- Uma apresentação **completa** do conjunto de treinamento corresponde a 1 ciclo (*epoch*).
- O processo de aprendizado é repetido ciclo após ciclo, até que os pesos sinápticos e níveis *threshold* se **estabilizem**.
- Tomar os pesos em uma forma aleatória → pesquisa no espaço dos pesos **estocástica**.

Modo padrão

- **(1) Modo padrão:**

- Atualização nos pesos é feita após a apresentação de cada exemplo de treinamento.
- Um ciclo consistindo de N exemplos de treinamento, arranjados na ordem:

$$\{[\mathbf{x}_1, \mathbf{d}_1], [\mathbf{x}_2, \mathbf{d}_2], \dots, [\mathbf{x}_N, \mathbf{d}_N]\}$$

- $[\mathbf{x}_1, \mathbf{d}_1] \rightarrow$ cálculos *forward/backward* e atualização dos pesos.
- $[\mathbf{x}_2, \mathbf{d}_2] \rightarrow$ cálculos *forward/backward* e atualização dos pesos.
- \vdots
- $[\mathbf{x}_N, \mathbf{d}_N] \rightarrow$ cálculos *forward/backward* e atualização dos pesos.

Modo padrão

- Dessa forma, a variação média nas mudanças dos pesos é:

$$\begin{aligned}\Delta \hat{w}_{ji} &= \frac{1}{N} \sum_{n=1}^N \Delta w_{ji}(n) \\ &= -\frac{\eta}{N} \sum_{n=1}^N \frac{\partial E(n)}{\partial w_{ji}(n)} \Rightarrow\end{aligned}$$

$$\boxed{\Delta \hat{w}_{ji} = -\frac{\eta}{N} \sum_{n=1}^N e_j(n) \frac{\partial e_j(n)}{\partial w_{ji}(n)}} \quad (2)$$

Modo *batch*

- **(2) Modo batch:**

- Atualização dos pesos é feita depois da apresentação de todos os exemplos de treinamento que constituem um ciclo.
- Para um ciclo particular, função custo com o erro quadrático médio:

$$\mathcal{E}_{av} = \frac{1}{2N} \sum_{n=1}^N \sum_{j \in C} e_j^2(n) \quad (3)$$

- Onde C denota o conjunto de índices correspondentes aos neurônios da camada de saída e e_j é o sinal do erro do neurônio j correspondente ao exemplo de treinamento w .

$$\Delta w_{ji} = -\eta \frac{\partial \mathcal{E}_{av}}{\partial w_{ji}} \implies$$

$$\Delta w_{ji} = \frac{\eta}{N} \sum_{n=1}^N e_j(n) \frac{\partial e_j(n)}{\partial w_{ji}}$$

(4)

Sumário

- 1 Introdução
 - Modelo de rede MLP
- 2 Treinamento de redes MLP
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 Funções Utilizadas
 - Função Softmax

Generalização

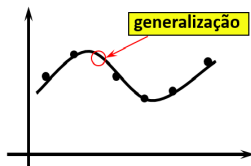
Aprendizado BP

conjunto treinamento + algoritmo BP \Rightarrow pesos sinápticos

GENERALIZAR

"GENERALIZAÇÃO": termo da psicologia.

Processo de aprendizado pode ser visto como um Método de Aproximação de Funções



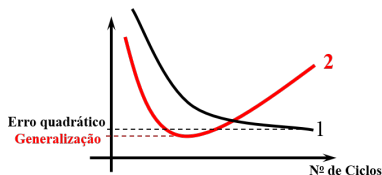
generalização : efeito de uma boa aproximação não linear dos dados de entrada, tamanho e eficiência do conjunto treinamento, arquitetura da rede, complexidade física do problema

- **Problema:** determinar o melhor número de nós na camada intermediária.
- Estatisticamente, esse problema é equivalente a determinar o tamanho do conjunto de parâmetros usado para modelar o conjunto de dados. Existe um limite no tamanho da rede.
- Esse limite deve ser tomado lembrando que é melhor treinar a rede para **produzir a melhor generalização** do que treinar a rede para representar perfeitamente um conjunto de dados.
- Isso pode ser feito usando **validação cruzada**.

- Conjunto de dados:
 - Treinamento (75%)
 - Teste (25%)
- Conjunto de treinamento:
 - Um subconjunto para validação do modelo.
 - Um subconjunto para treinamento.
- Validar o modelo em um conjunto diferente do usado para estimá-lo.

- Usa-se o subconjunto de validação para avaliar o desempenho de diferentes candidatos do modelo (diferentes topologias) e, então, escolhe-se uma delas.
- O modelo escolhido é treinado sobre o conjunto de treinamento inteiro e a capacidade de generalização é medida no conjunto de teste.
- A validação cruzada pode ser usada para decidir quando o treinamento de uma rede deve ser encerrado.

Tamanho do conjunto de treinamento



Curva 1: poucos parâmetros (*underfitting*)

Curva 2: muitos parâmetros (*overfitting*)

- Em ambos os casos:
 - ① O desempenho do erro na generalização exibe um mínimo.
 - ② O mínimo no caso *overfitting* é menor e mais definido.
- Pode-se obter boa generalização se a rede é projetada com muitos neurônios, contanto que o treinamento seja cessado após um número de ciclos correspondente ao mínimo da curva do **erro** obtida na **validação cruzada**.

Sumário

- 1 Introdução
 - Modelo de rede MLP
- 2 Treinamento de redes MLP
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 Funções Utilizadas
 - Função Softmax

Sumário

- 1 **Introdução**
 - Modelo de rede MLP
- 2 **Treinamento de redes MLP**
 - O algoritmo Backpropagation
 - Processo de aprendizado
 - Termo Momentum
 - Velocidade de aprendizado
 - Modos de treinamento
 - Generalização
- 3 **Funções Utilizadas**
 - Função Softmax

Problema

- Em redes neurais *feedforward*, o aprendizado pode tornar-se lento ou pouco efetivo, devido a alguns fatores
- A otimização da função *erro quadrático médio* provoca saturação dos neurônios de saída, tornando o aprendizado mais lento.
- Quando uma entrada x_j é próxima de zero, o peso correspondente w_j se ajusta muito devagar.
- Otimizar a função *cross-entropy*, em vez do erro quadrático médio.
- Usar a função de ativação *softmax* nos neurônios de saída.

Função cross-entropy

- Pode ser a função a ser otimizada, alternativamente ao erro quadrático médio.

$$\text{cross-entropy} = - \sum y' \log y$$

Função Softmax

- A função *softmax* pode ser utilizada como função de ativação na última camada de uma rede neural.
- Seja v_i a combinação linear das entradas e pesos do neurônio de saída i , define-se:

$$\text{softmax}(v_i) = \frac{\exp(v_i)}{\sum_j \exp(v_j)}$$

- Aplicando a função *softmax* na camada de saída, obtém-se uma **distribuição de probabilidade** válida, em que a soma das saídas dos neurônios de saída é igual a 1.
 - A saída de cada neurônio representa a probabilidade de a classe correspondente ser a classificação correta.