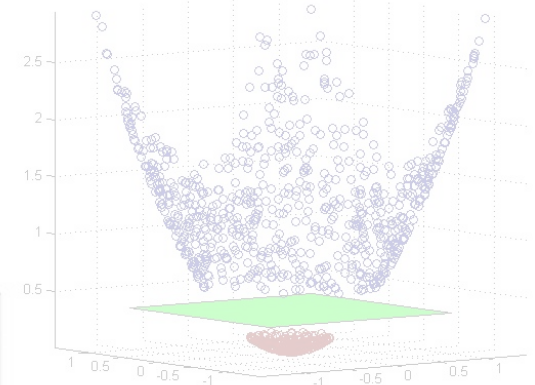
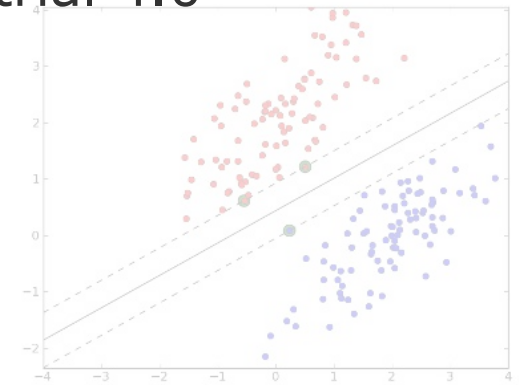
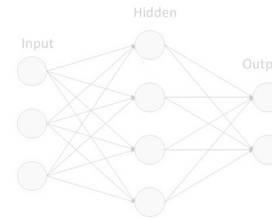
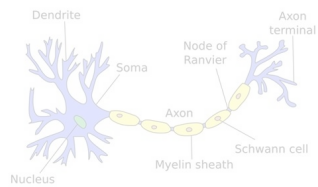
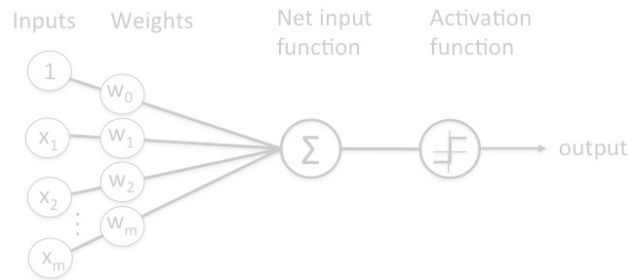
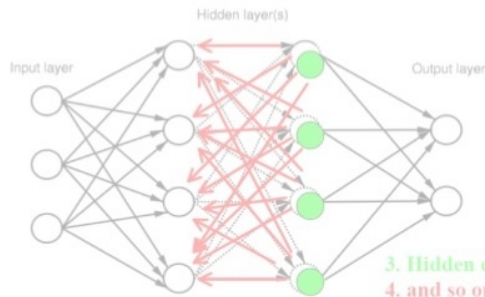


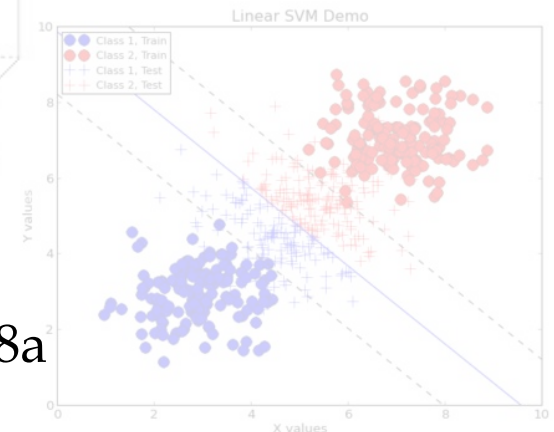
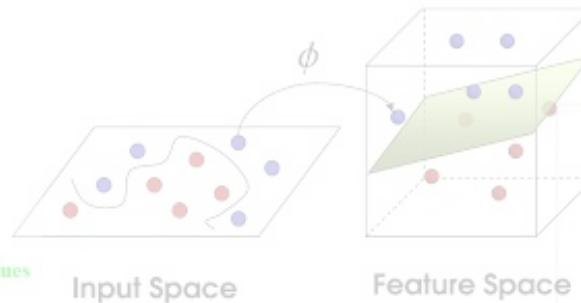
Classificação



Backpropagation Learning



3. Hidden error values
4. and so on ...



Backpropagation Learning

$$E_{out} = d_{out} - out$$

$$E_{total} = \sum_{i=1}^{num(n_{out})} E_{out}^2$$

$$E_{hid} = \sum_{k=1}^{num(n_{hid})} E_{out,k} \cdot w_{out,k}$$

$$dE_{hid} = E_{hid} \cdot (1 - o_{hid}) \cdot o_{hid}$$

David Menotti

www.inf.ufpr.br/menotti/am-18a

Hoje

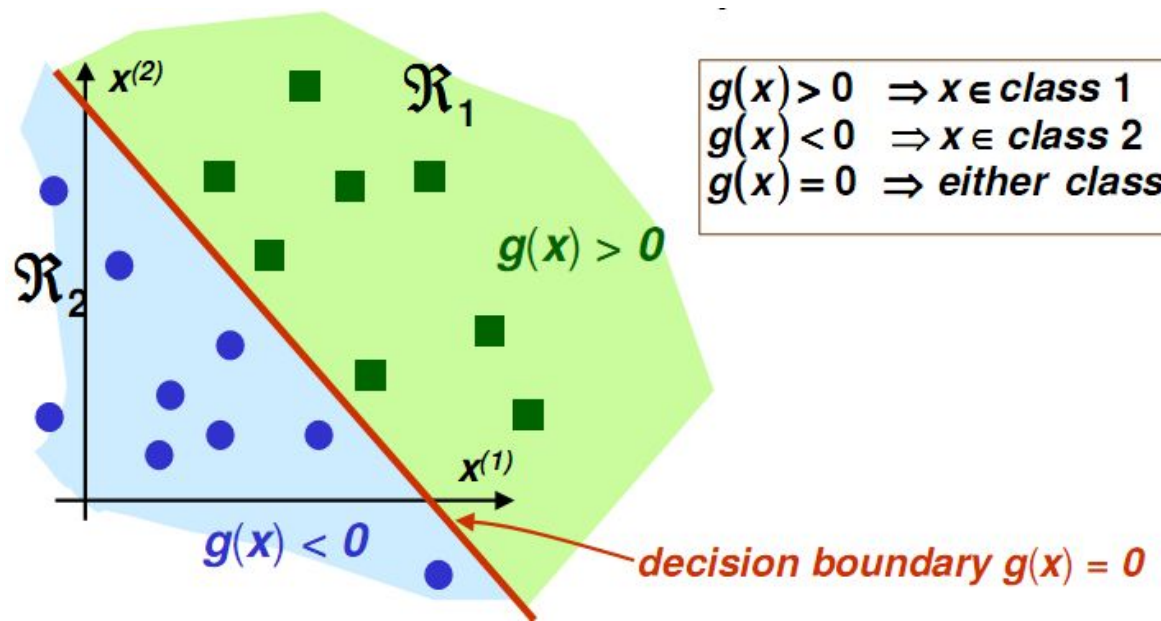
- Funções Discriminantes Lineares
 - Perceptron
 - Support Vector Machines (SVM)
- Multiple Layer Perceptron (MLP)

Funções Discriminantes Lineares

Funções Discriminantes Lineares

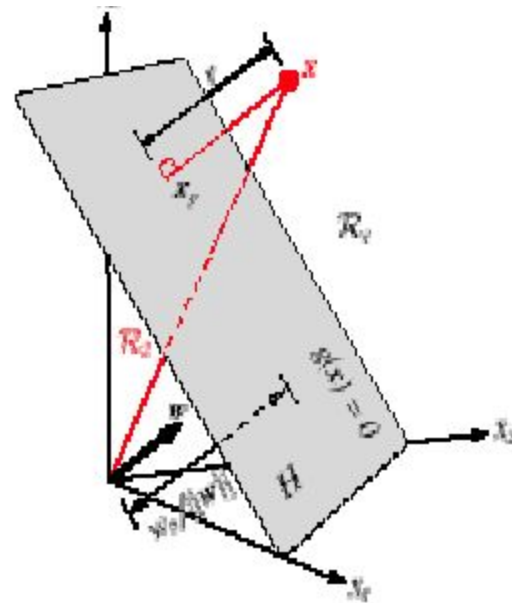
- Em geral, uma função discriminante linear pode ser escrita na forma

- $g(x) = \omega^T x + \omega_0$ é conhecido como o vetor dos pesos e ω_0 representa o *bias*



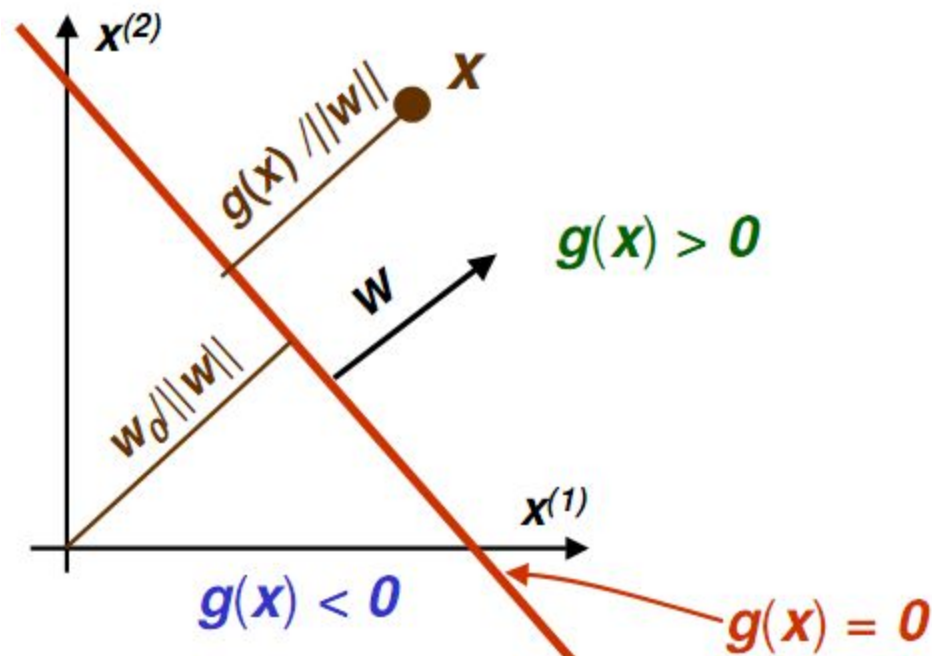
Funções Discriminante Lineares

- $g(x) = \omega^T x + \omega_0$ é um hiperplano
 - Um hiperplano é
 - Um ponto em 1D
 - Uma reta em 2D
 - Um plano em 3D



Funções Discriminante Lineares

- Para duas dimensões,
 - ω determina a orientação do hiperplano
 - ω_0 representa o deslocamento (origem)



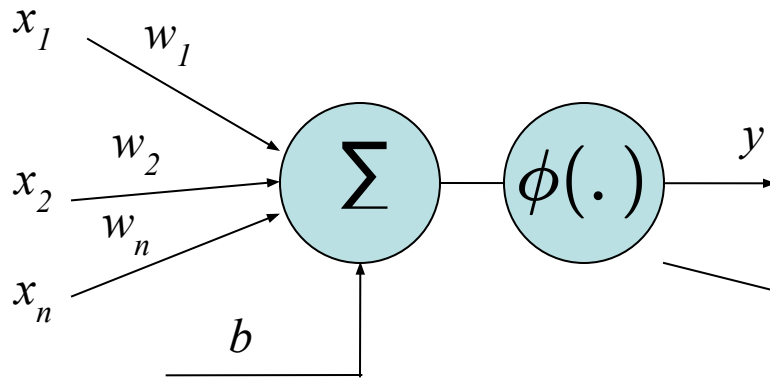
Funções Discriminante Lineares

Perceptron

Perceptron

- Um classificador linear bastante simples, mas bastante importante no desenvolvimento das redes neurais é o **Perceptron**.
 - O perceptron é considerado como sendo a primeira e mais primitiva estrutura de rede neuronal artificial.
 - Concebido por McCulloch and Pits na década de 50.
- Diferentemente do **LDA**, o perceptron não transforma os dados para fazer classificação.
 - Tenta encontrar a melhor fronteira linear que separa os dados.

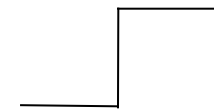
Perceptron



A função de ativação normalmente utilizada no perceptron é a **hardlim** (*threshold*)

$$\phi(z) = \begin{cases} 1 & \text{se } z \geq 0, \\ 0 & \text{se } z < 0 \end{cases}$$

$$y = \phi \left(b + \sum_i \omega_i x_i \right)$$



A **função de ativação** é responsável por determinar a forma e a intensidade de alteração dos valores transmitido de um neurônio a outro.

Perceptron: Algoritmo de Aprendizagem

1. Iniciar os pesos e bias com valores pequenos, geralmente no intervalo [0,3 ; 0,8]
2. Aplicar um padrão de entrada com seu respectivo valor desejado de saída (t_j) e verificar a saída y da rede.
3. Calcular o erro da saída
4. Se $e = 0$, volta ao passo $e = t_j - y$
5. Se $e \neq 0$,
 1. Atualizar pesos
 2. Atualizar o bias
6. Voltar ao passo 2
 - Critério de parada: Todos os padrões classificados corretamente.

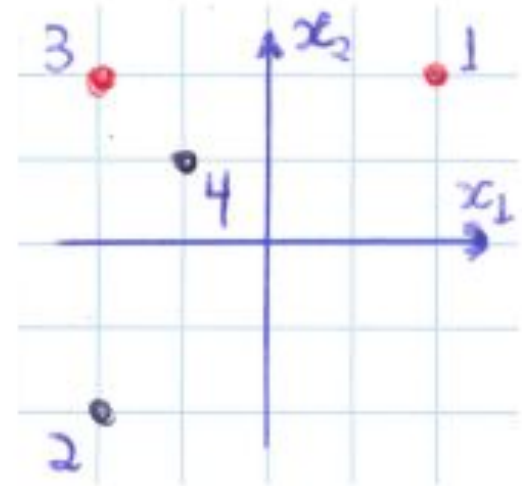
$$\omega_i = \omega_i^{old} + e \times x_i$$

$$b_j = b_j^{old}$$

Perceptron: Exemplo

- Considere o seguinte conjunto de aprendizagem.

	x	t
2	2	0
-2	-2	1
-2	2	0
-1	1	1



Nesse tipo de algoritmo é importante que os dados sejam apresentados ao algoritmo de treinamento de maneira intercalada (shuffle)

Perceptron: Exemplo

Exemplo

- **Nesse** exemplo, vamos inicializar:

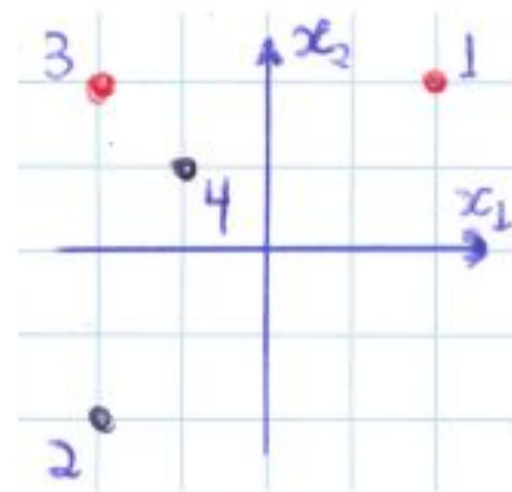
- $\omega = [0, 0]$ e $b = 0$

- E a Função Discriminante?

- Fica indefinida!

- $g(x) = \omega^T x + b = 0$

- Primeira amostra definirá um hiperplano



Perceptron: Exemplo

Exemplo

- Apresentando o primeiro padrão (X_1) a rede.

$$y = \text{hardlim}([0, 0][2, 2]^T + 0) = \text{hardlim}(0) = 1$$

- Calcula-se o erro $e = t_i - y = 0 - 1 = -1$
- Como o erro é diferente de 0, atualiza-se os pesos e o bias

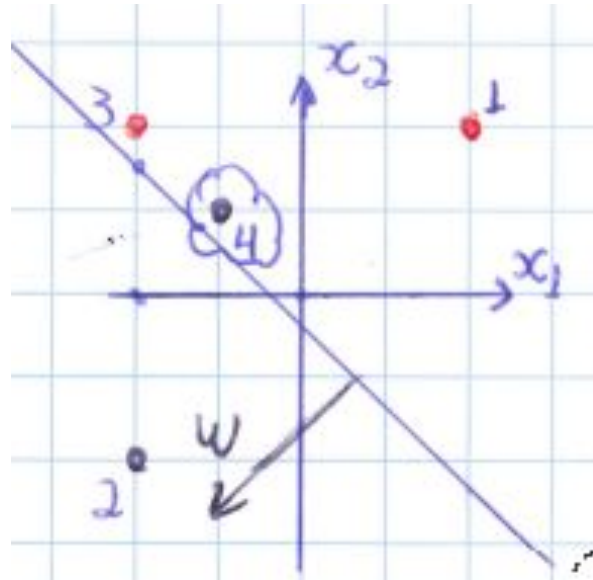
$$\omega = \omega^{old} + e \times x_i = [0, 0] + (-1[2, 2]) = [-2, -2]$$

$$b = b^{old} + e = 0 + (-1) = -1$$

Perceptron

Exemplo

$$-2x_1 - 2x_2 - 1 \geq 0$$



$$\omega = [-2, -2], b = -1$$

Perceptron

Exemplo

- Apresentando o segundo padrão (X_2) a rede:

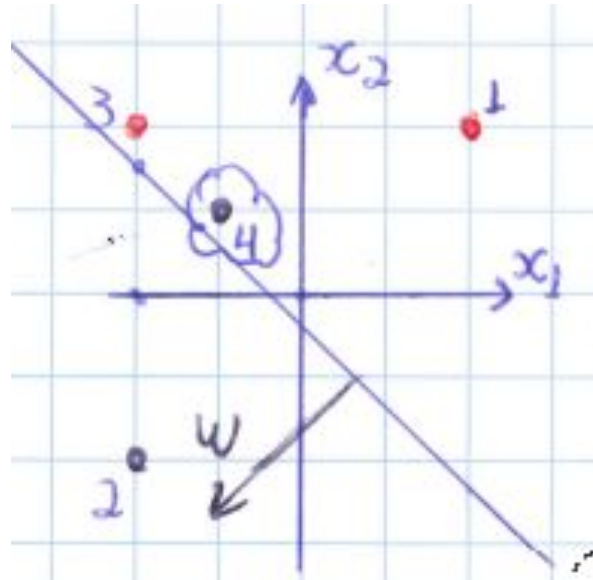
$$y = \text{hardlim}([-2, -2][-2, -2]^T + (-1)) = \text{hardlim}(7) = 1$$

- Calcula-se o erro $e = t_i - y = 1 - 1 = 0$
- Como o erro é 0, os pesos e o bias não precisam ser atualizados.

Perceptron

Exemplo

$$-2x_1 - 2x_2 - 1 \geq 0$$



$$\omega = [-2, -2], b = -1$$

Perceptron

Exemplo

- Apresentando o terceiro padrão (X_3) a rede:

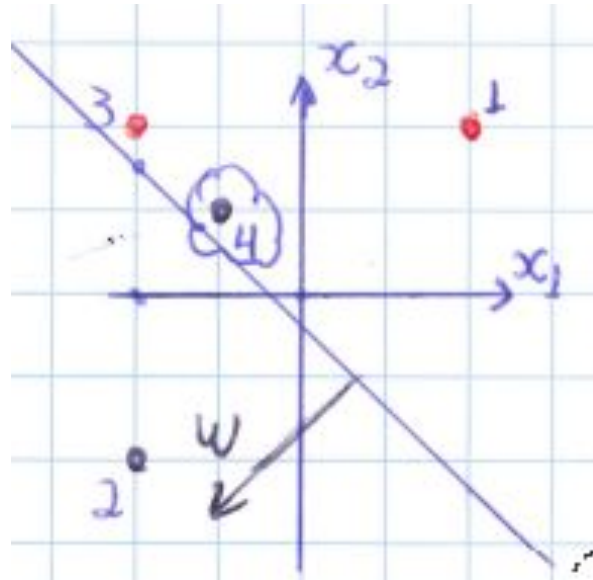
$$y = \text{hardlim} \left([-2, -2] [-2, +2]^T + (-1) \right) = \text{hardlim}(-1) = 0$$

- Calcula-se o erro $e = t_i - y = 0 - 0 = 0$
- Como o erro é 0, os pesos e o *bias* não precisam ser atualizados.

Perceptron

Exemplo

$$-2x_1 - 2x_2 - 1 \geq 0$$



$$\omega = [-2, -2], b = -1$$

Perceptron

Exemplo

- Apresentando o quarto padrão (X_4) a rede:

$$y = \text{hardlim}([-2, -2][-1, +1]^T + (-1)) = \text{hardlim}(-1) = 0$$

- Calcula-se o erro $e = t_i - y = 1 - 0 = 1$
- Como o erro é diferente de 0, atualiza-se os pesos e o bias

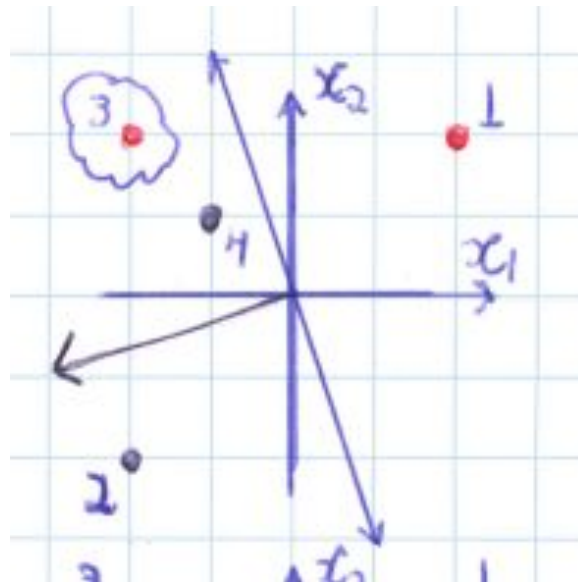
$$\omega = \omega^{old} + e \times x_i = [-2, -2] + (+1[-1, 1]) = [-3, -1]$$

$$b = b^{old} + e = -1 + 1 = 0$$

Perceptron

Exemplo

$$-3x_1 - 1x_2 - 1 \geq 0$$



$$\omega = [-3, -1], b = 0$$

Perceptron

Exemplo

- O processo evolui por várias iterações

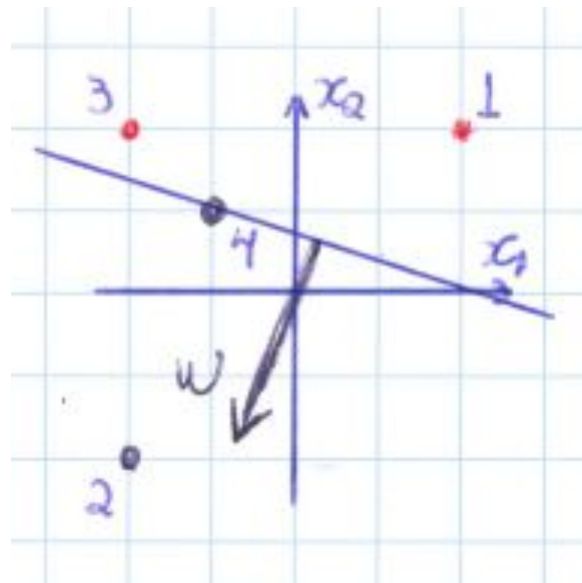
it: 1, i: 0, X=[2.00, 2.00], t= 0.00, y = 1.00, e = -1.00
W=[-2.00, -2.00], B= -1.00
it: 2, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 3, i: 2, X=[-2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 4, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-3.00, -1.00], B= 0.00
it: 5, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 6, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 7, i: 2, X=[-2.00, 2.00], t= 0.00, y = 1.00, e = -1.00
W=[-1.00, -3.00], B= -1.00
it: 8, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-2.00, -2.00], B= 0.00
it: 9, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 10, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 11, i: 2, X=[-2.00, 2.00], t= 0.00, y = 1.00, e = -1.00
W=[0.00, -4.00], B= -1.00
it: 12, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-1.00, -3.00], B= 0.00
it: 13, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 14, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 15, i: 2, X=[-2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 16, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-2.00, -2.00], B= 1.00

it: 17, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 18, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 19, i: 2, X=[-2.00, 2.00], t= 0.00, y = 1.00, e = -1.00
W=[0.00, -4.00], B= 0.00
it: 20, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-1.00, -3.00], B= 1.00
it: 21, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 22, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 23, i: 2, X=[-2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 24, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-2.00, -2.00], B= 2.00
it: 25, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 26, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 27, i: 2, X=[-2.00, 2.00], t= 0.00, y = 1.00, e = -1.00
W=[0.00, -4.00], B= 1.00
it: 28, i: 3, X=[-1.00, 1.00], t= 1.00, y = 0.00, e = 1.00
W=[-1.00, -3.00], B= 2.00
it: 29, i: 0, X=[2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 30, i: 1, X=[-2.00, -2.00], t= 1.00, y = 1.00, e = 0.00
it: 31, i: 2, X=[-2.00, 2.00], t= 0.00, y = 0.00, e = 0.00
it: 32, i: 3, X=[-1.00, 1.00], t= 1.00, y = 1.00, e = 0.00

Perceptron

- O processo acaba quando **todos** os padrões forem classificados corretamente.

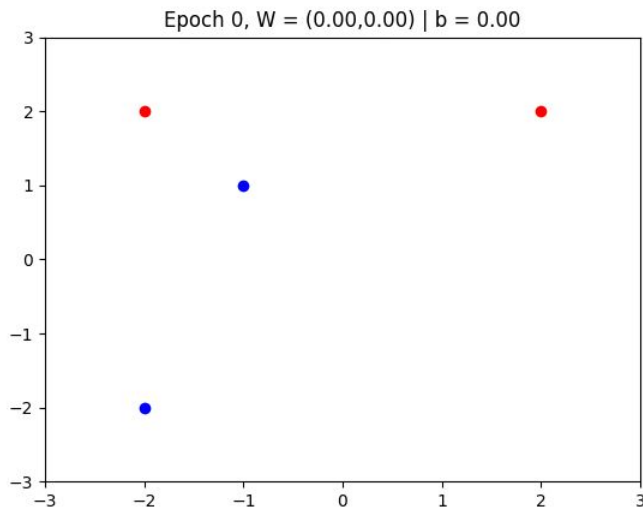
$$-1x_1 - 3x_2 + 2 \geq 0$$



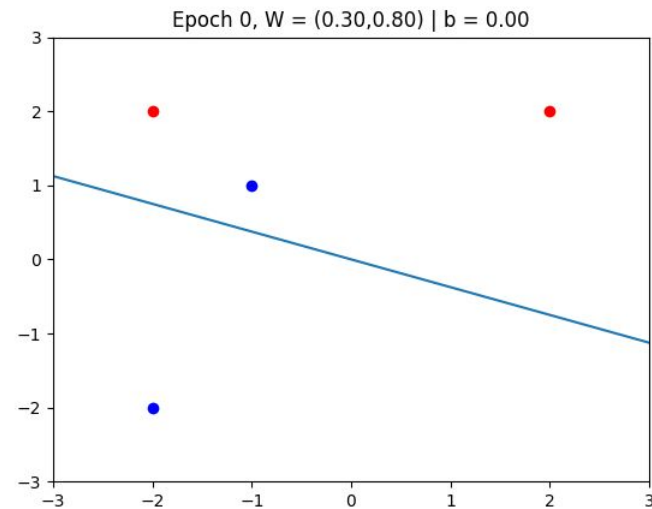
$$\omega = [-1, -3], b = 2$$

Perceptron

Exemplo



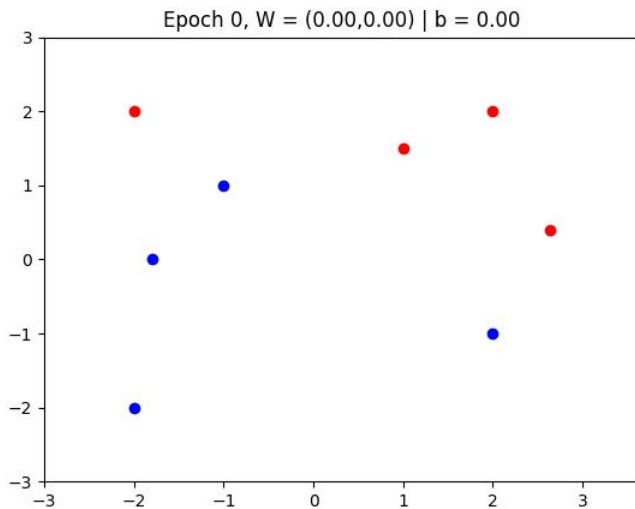
$$0 : \omega = [0.0, 0.0] \quad b = 0.0$$



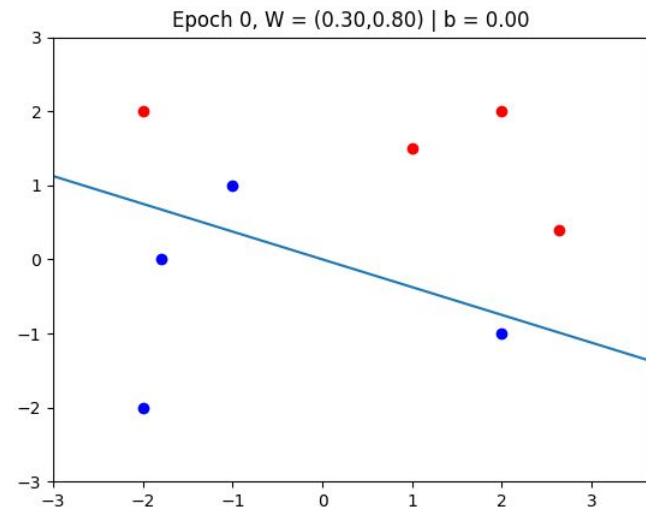
$$0 : \omega = [0.3, 0.8] \quad b = 0.0$$

Perceptron

Exemplo - 8 pontos



$$0 : \omega = [0.0, 0.0] \quad b = 0.0$$



$$0 : \omega = [0.3, 0.8] \quad b = 0.0$$

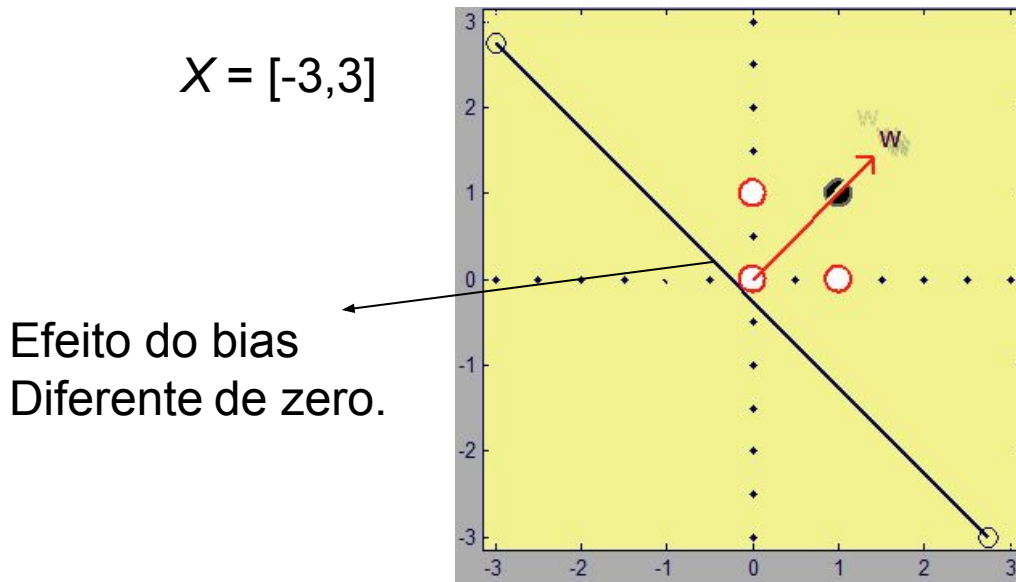
Determinando a fronteira

- No caso bi-dimensional, a fronteira de decisão pode ser facilmente encontrada usando a seguinte equação

$$x_2 = \frac{-\omega_1 \times x_1 - b}{\omega_2} \quad \omega_1 \times x_1 + \omega_2 \times x_2 + b = 0$$

Considere o seguinte exemplo, $\omega = [1, 41; 1, 41], b = 0, 354$

Escolha duas coordenadas x_1 , para então encontrar os x_2 correspondentes



Para $x_1 = -3, x_2 = 2,74$
Para $x_1 = +3, x_2 = -3,25$

Perceptron

Código fonte da simulação

```
#include <stdio.h>

#define TAM 4

double hardlim(double* X, double* W, double B)
{
    double sum = B; int i;
    for(i=0; i<2; i++)
        sum += W[i] * X[i];
    return ( sum >= 0. ) ? 1. : 0.;
}
```

```
int main(void){
    int c,n,i,it;
    double y,e;
    double X[][2] = {{2.,2.},{-2.,-2.},{-2.,2.},{-1.,1.}};
    double T[] = {0.,1.,0.,1.};

    double W[] = {0.,0.}, B = 0.;

    it=0;n = 0;c=0;
    while (c<4) {
        it++;

        y = hardlim(X[n],W,B);
        e = T[n] - y;

        printf("it:%3d, i:%2d, X=[%6.2lf,%6.2lf], "
            , it,n,X[n][0],X[n][1]);
        printf("t=%6.2lf, y = %6.2lf, e = %6.2lf\n"
            ,T[n],y,e);

        if ( e != 0. ) {
            for(i=0; i<2; i++) {
                W[i] = W[i] + e*X[n][i];
            }
            B = B + e;
            printf("W=[%6.2lf,%6.2lf],",W[0],W[1]);
            printf("B=%6.2lf\n",B);
            c=0;
        }
        else
            c++;
        n = (n==TAM-1) ? 0 : n+1;
    }
}
```

Perceptron

Código fonte da visualização

```
#!/usr/bin/python # -*- encoding: iso-8859-1 -*-
# Aprendizagem de Máquina / Rayson Bartoski Laroca dos Santos
import sys,math,time,pylab
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_svmlight_file, make_classification

i = 0
def predict(row, weights):
    activation = weights[0]
    for i in range(len(row)):
        activation += weights[i + 1] * row[i]
    # hardlim
    return 1.0 if activation >= 0.0 else 0.0

def update_weights(weights, error, X):
    global i

    w1 = weights[1] + error * X[i][0]
    w2 = weights[2] + error * X[i][1]
    bias = weights[0] + error

    return bias,w1,w2

def plot_line(X,y,weights,epoch,time):
    for x,label in zip(X,y):
        if label == 0:
            plt.scatter(x[0],x[1], color='red')
        else:
            plt.scatter(x[0],x[1], color='blue')

# equação da reta - x * w1 + y * w2 + bias = 0
xx = []
yy = []
x_min = min(X.transpose()[0])
x_max = max(X.transpose()[0])
y_min = min(X.transpose()[1])
y_max = max(X.transpose()[1])

# point 1
x = x_max + 1
if weights[2] == 0:
    y = 0
else:
    y = (x * weights[1] + weights[0]) / - weights[2]
xx.append(x)
yy.append(y)

# point 2
x = x_min - 1

if weights[2] == 0:
    y = 0
else:
    y = (x * weights[1] + weights[0]) / - weights[2]
xx.append(x)
yy.append(y)

# plot perceptron line
fig = pylab.gcf()
fig.canvas.set_window_title('Epoch ' + str(epoch))

if weights[1] <= 0 and weights[2] <= 0:
    plt.plot(xx,yy)
plt.title('Epoch {0:1d}, W = ({1:.2f},{2:.2f}) | b = {3:.2f}'
        .format(epoch,weights[1],weights[2],weights[0]))
plt.xlim(x_min-1,x_max+1)
plt.ylim(y_min-1,y_max+1)
fig.savefig('./perceptron'+str(epoch))

# updates every x seconds
plt.pause(x seconds)
plt.clf()

def main(data):
    global i

    t = 0.50 # tempo entre um plot e outro.
    t_final = 2 # tempo que o plot fica no final

    if data == None:
        print 'No dataset given | Usage: perceptron.py <data>'
        # print 'Loading default data...'
        # time.sleep(2)
        # X = np.array([[2,2],[-2,-2],[-2,2],[-1,1],[1,1.5],[2,-1],[2.64,6]
        # y = np.float64([0,1,0,1,0,1,0,1])

        print 'Loading random data...'
        time.sleep(2)
        X, y = make_classification(n_samples=30, n_features=2,
                                n_redundant=0, n_clusters_per_class=1)

    else:
        print "Loading data..."
        X, y = load_svmlight_file(data)
        X = X.toarray()

    weights = [0, 0.3, 0.8]
    #weights = [0, 0, 0]

    print 'W = ({0:.2f},{1:.2f}) | b = {2:.2f}'.format(weights[1],weights[

epoch = 0
correct = 0

plot_line(X,y,weights,epoch, t)
while correct < len(X) and epoch < 25:
    if i == len(X):
        i = 0
        prediction = predict(X[i], weights)
        gt = y[i]
        print(" Pattern %d -> gt=%d, Predicted=%d" % (i, gt, prediction))
        if gt == prediction:
            correct += 1
        else:
            error = gt - prediction

            weights = update_weights(weights, error, X)

    print '\nUpdate Weights:'
    print ' Weights = ({0:.2f},{1:.2f}) | bias = {2:.2f}'.format(weights[1],wei

    # i = 0
    correct = 0

    if i == 0:
        epoch += 1
        plot_line(X,y,weights,epoch, t)

        i += 1

    print '\nResults:'
    print ' Epochs:', epoch
    print ' Weights = ({0:.2f},{1:.2f}) | bias = {2:.2f}'.format(weights[1],wei
    plot_line(X,y,weights,epoch,t_final)

if __name__ == "__main__":
    time1 = time.time()

    if len(sys.argv) == 1:
        main(None)
    elif len(sys.argv) != 2:
        sys.exit("Usage: python perceptron.py data.dat")
    else:
        main(sys.argv[1])

    time2 = time.time()
    seconds = (time2-time1)
    m, s = divmod(seconds, 60)
    h, m = divmod(m, 60)
```

Perceptron

Usando um Perceptron em python

Load the required elements

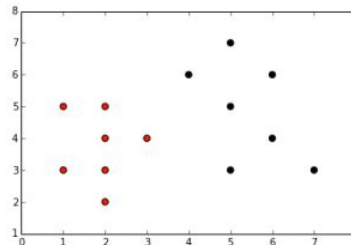
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 #import sklearn.linear_model.perceptron as p
4 from sklearn.linear_model import perceptron
5
6 # Needed to show the plots inline
7 %matplotlib inline
```

Set the data

```
1 # Data
2 d = np.array([
3 [2, 1, 2, 5, 7, 2, 3, 6, 1, 2, 5, 4, 6, 5],
4 [2, 3, 3, 3, 3, 3, 4, 4, 4, 4, 5, 5, 5, 6, 6, 7]
5 ])
6
7 # Labels
8 t = np.array([0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1])
```

Plot the Data (see what it looks like)

```
1 colormap = np.array(['r', 'k'])
2 plt.scatter(d[0], d[1], c=colormap[t], s=40)
```



Example Data

Perceptron

Usando um Perceptron em python

- Treinando e classificando

```
Python
# rotate the data 180 degrees
d90 = np.rot90(d)
d90 = np.rot90(d90)
d90 = np.rot90(d90)

# Create the model
net = perceptron.Perceptron(n_iter=100, verbose=0, random_state=None, fit_intercept=True, eta0=0.002)
net.fit(d90,t)

# Print the results
print "Prediction " + str(net.predict(d90))
print "Actual      " + str(t)
print "Accuracy    " + str(net.score(d90, t)*100) + "%"
```



Levando os dados para formato linha

As we can see, the model has *solved* the problem.

```
1 Prediction [0 0 0 1 1 0 0 1 0 0 1 1 1]
2 Actual    [0 0 0 1 1 0 0 1 0 0 1 1 1]
3 Accuracy  100.0%
```

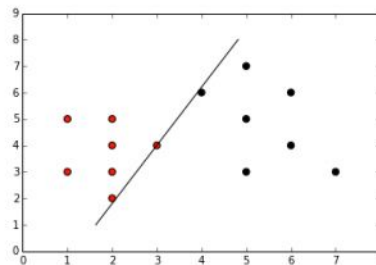
Perceptron

Usando um Perceptron em python

- Visualizando a função criada

Plot the Decision Boundary

```
1 # Plot the original data
2 plt.scatter(d[0], d[1], c=colormap[t], s=40)
3
4 # Output the values
5 print "Coefficient 0 " + str(net.coef_[0,0])
6 print "Coefficient 1 " + str(net.coef_[0,1])
7 print "Bias " + str(net.intercept_)
8
9 # Calc the hyperplane (decision boundary)
10 ymin, ymax = plt.ylim()
11 w = net.coef_[0]
12 a = -w[0] / w[1]
13 xx = np.linspace(ymin, ymax)
14 yy = a * xx - (net.intercept_[0]) / w[1]
15
16 # Plot the line
17 plt.plot(yy,xx, 'k-')
```



Decision Boundary – Values on either side of the line are in different groups

Funções Discriminante Lineares

***Support Vector Machines
(SVM)***

SVM

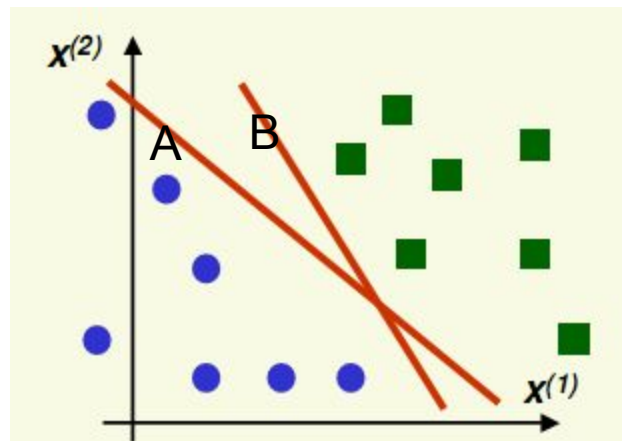


- Proposto Vladimir Vapnik / Histórico
 - Introdução (1963)
 - Kernel Trick (1992)
 - Soft Margin (1995)
- Um dos mais importantes acontecimentos na área de **reconhecimento de padrões** nos últimos 25 anos.
- Tem sido largamente utilizado com sucesso para resolver diferentes problemas.

- Vapnik & Lerner (1963). Pattern Recognition using Generalized Portrait Method, ARC, 24.
- Boser, Guyon & Vapnik (1992). "A training algorithm for optimal margin classifiers" – COLT '92.
- Cortes, C. & Vapnik (1995) V. N. Support-vector Networks. Machine Learning 20 (3), 1995

SVM - Introdução

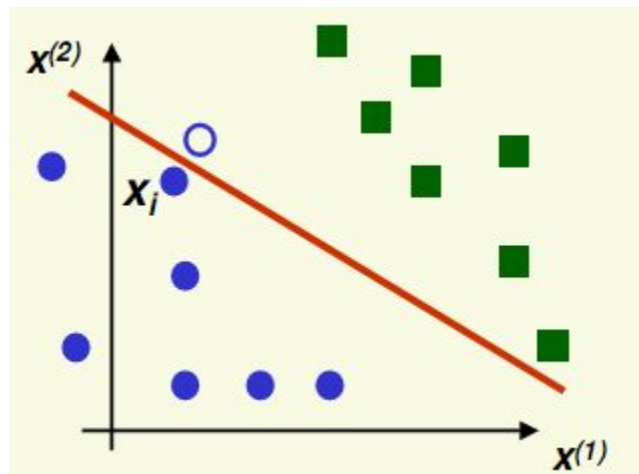
- Como vimos anteriormente, o perceptron é capaz de construir uma fronteira se os dados forem linearmente separáveis.



Mas qual a fronteira que deve ser escolhida??

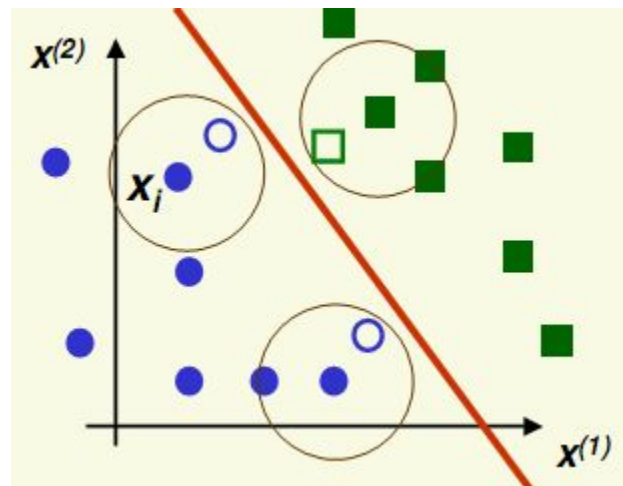
SVM - Introdução

- Suponha que a fronteira escolhida é a A.
- Como ela está bem próxima da classe azul, seu **poder de generalização é baixo**
 - Note que um novo elemento (dados não usados no treinamento), bem próximo de um azul será classificado erroneamente.



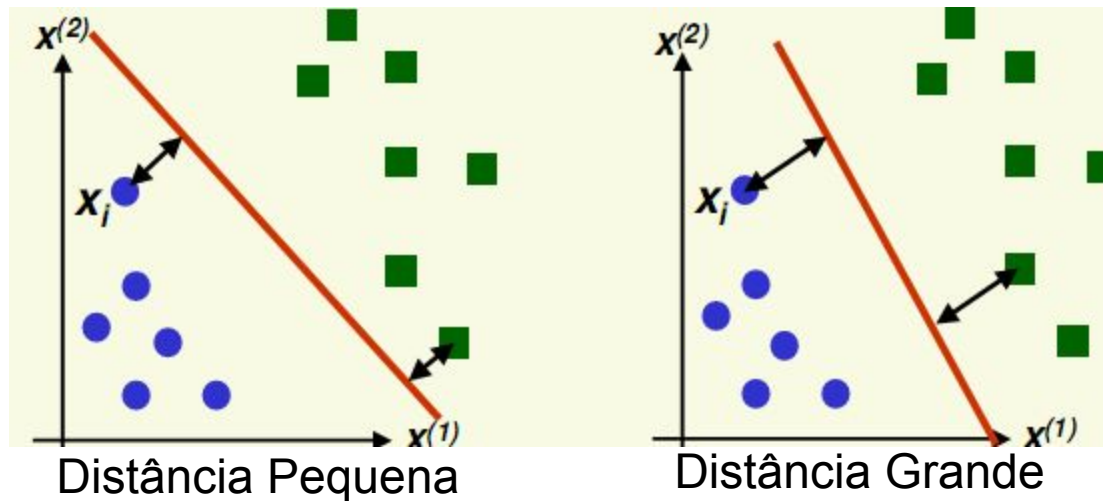
SVM - Introdução

- Escolhendo a fronteira B, podemos notar que o poder de generalização é bem melhor.
- Novos dados são corretamente classificados, pois temos uma fronteira mais distante dos dados de treinamento.



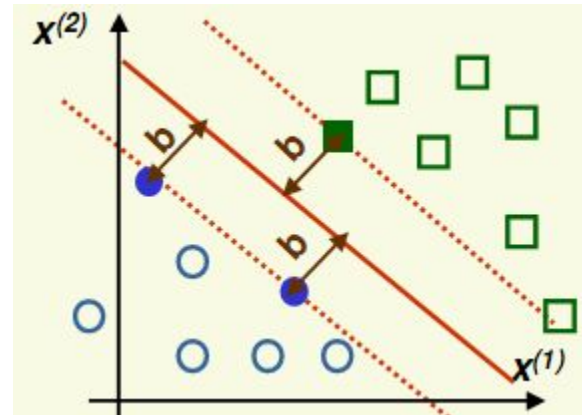
Maximização da Margem

- O conceito por trás do SVM é a **maximização da margem**, ou seja, maximizar a distância da margem dos dados de treinamento



Hiperplano ótimo: Distância da margem para o exemplo da classe positiva é igual a distância da margem para o exemplo da classe negativa.

Vetores de Suporte



- São os exemplos da base de treinamento mais próximos do hiperplano.
 - O hiperplano é definido unicamente pelos vetores de suporte, os quais são encontrados durante o treinamento.
 - Minimização de uma função quadrática
 - Alto custo computacional.

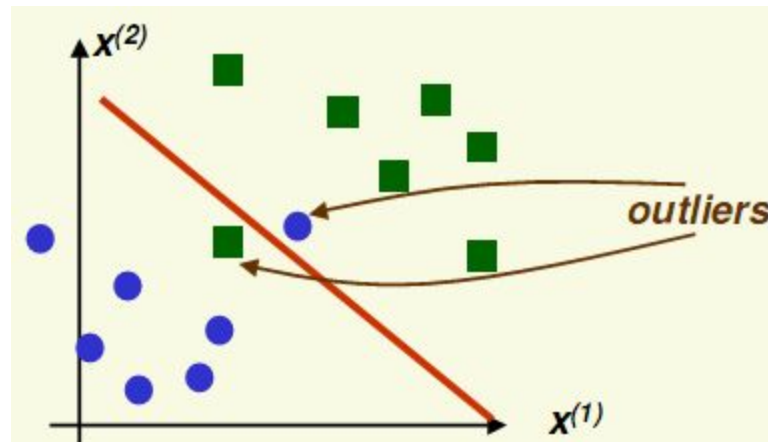
SVM: Decisão

$$f(x) = \sum_i \alpha_i y_i K(x, x_i) + b$$

- A função de decisão pode ser descrita pela fórmula acima, na qual,
 - K é a função de **kernel**,
 - α e b são os parâmetros encontrados durante o treinamento,
 - x_i e y_i são:
 - os vetores de características e
 - o *label* da classe,respectivamente.

Soft Margin

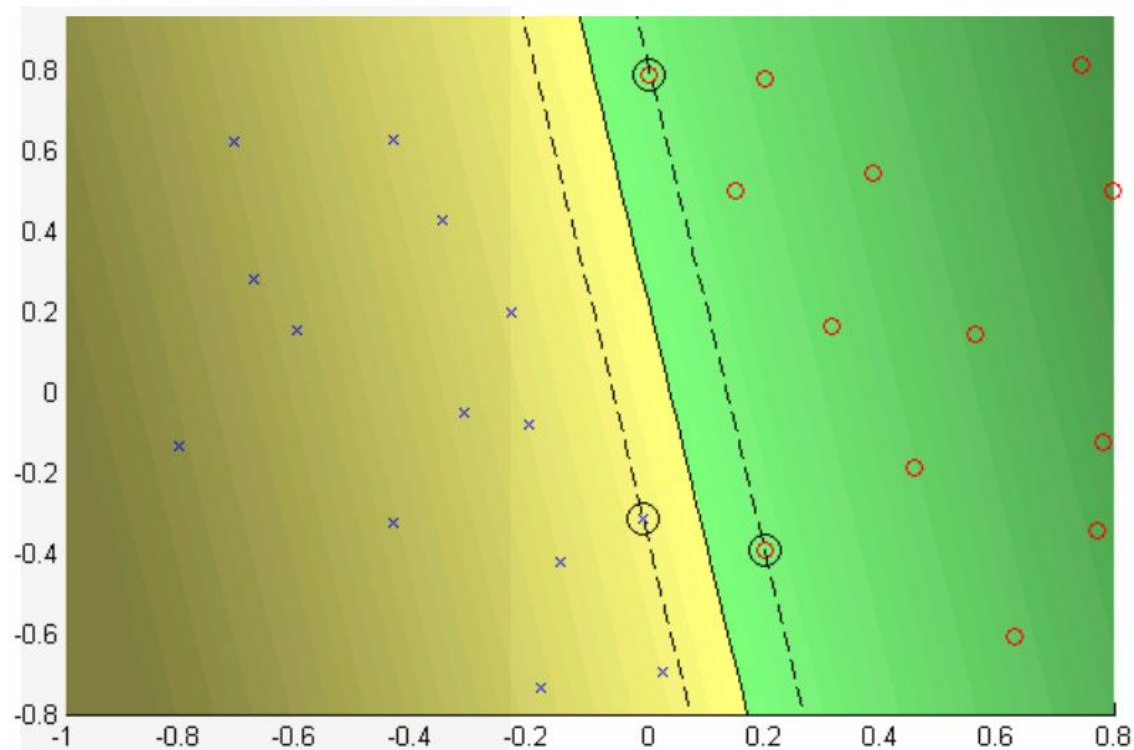
- Mesmo para dados que não podem ser separados linearmente, o SVM ainda pode ser apropriado.
- Isso é possível através do uso das “variáveis de folga” (**parâmetro C**).



Para um bom desempenho, os dados devem ser “quase” linearmente separáveis

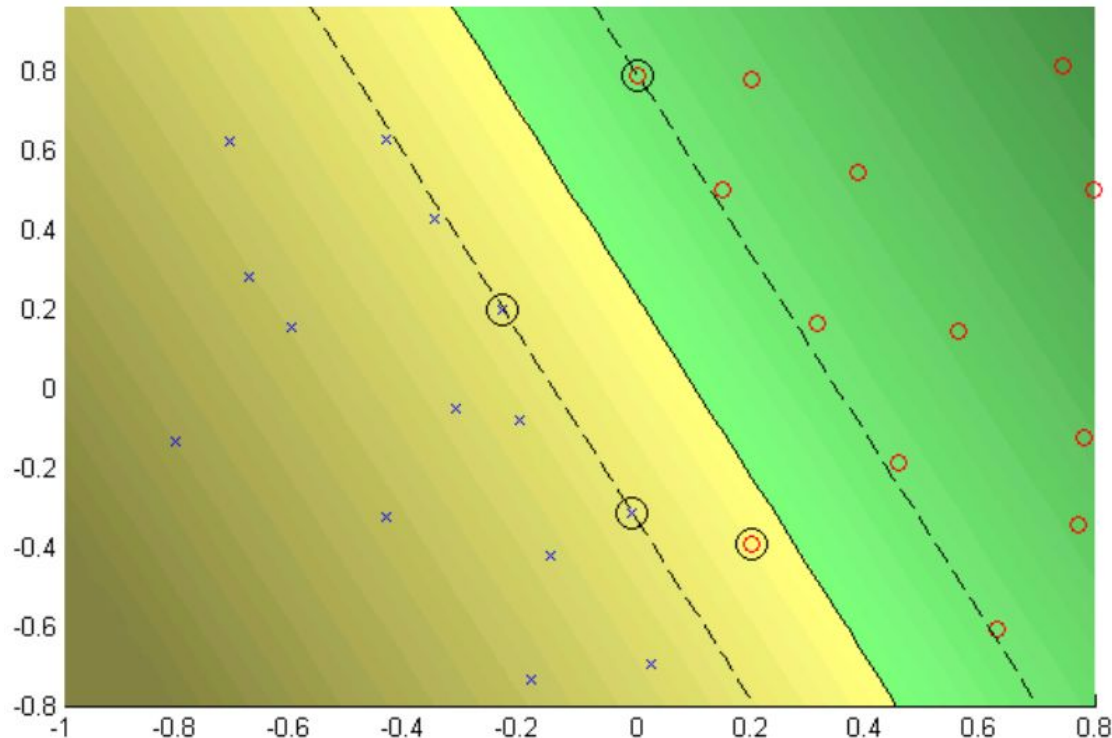
Soft Margin

- Quanto maior o número da variável de folga (C), mais *outliers* serão descartados.
- $C \rightarrow \infty$ **hard margin**



Soft Margin

- Quanto maior o número da variável de folga (C), mais *outliers* serão descartados.
- **C** = 10 **soft margin**



Mapeamento não Linear

- A grande maioria dos problemas reais não são linearmente separáveis.
- A pergunta então é:
 - “Como resolver problemas que não são linearmente separáveis com um classificador linear?”

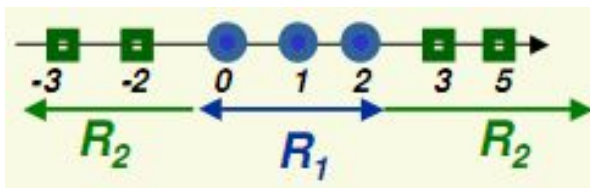
Projetar os dados em um espaço onde os dados são linearmente separáveis.



Mapeamento não Linear

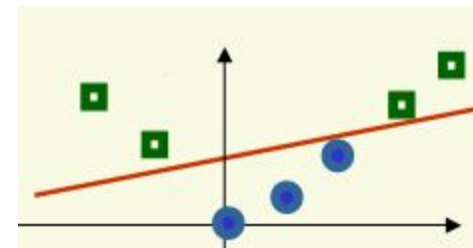
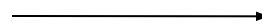
- Projetar os dados em outra dimensão usando uma função de kernel (*kernel trick*).
- Encontrar um hiperplano que separe os dados nesse espaço.

Em qual dimensão esses dados seriam linearmente separáveis?



1D

$$\phi(x, x^2)$$



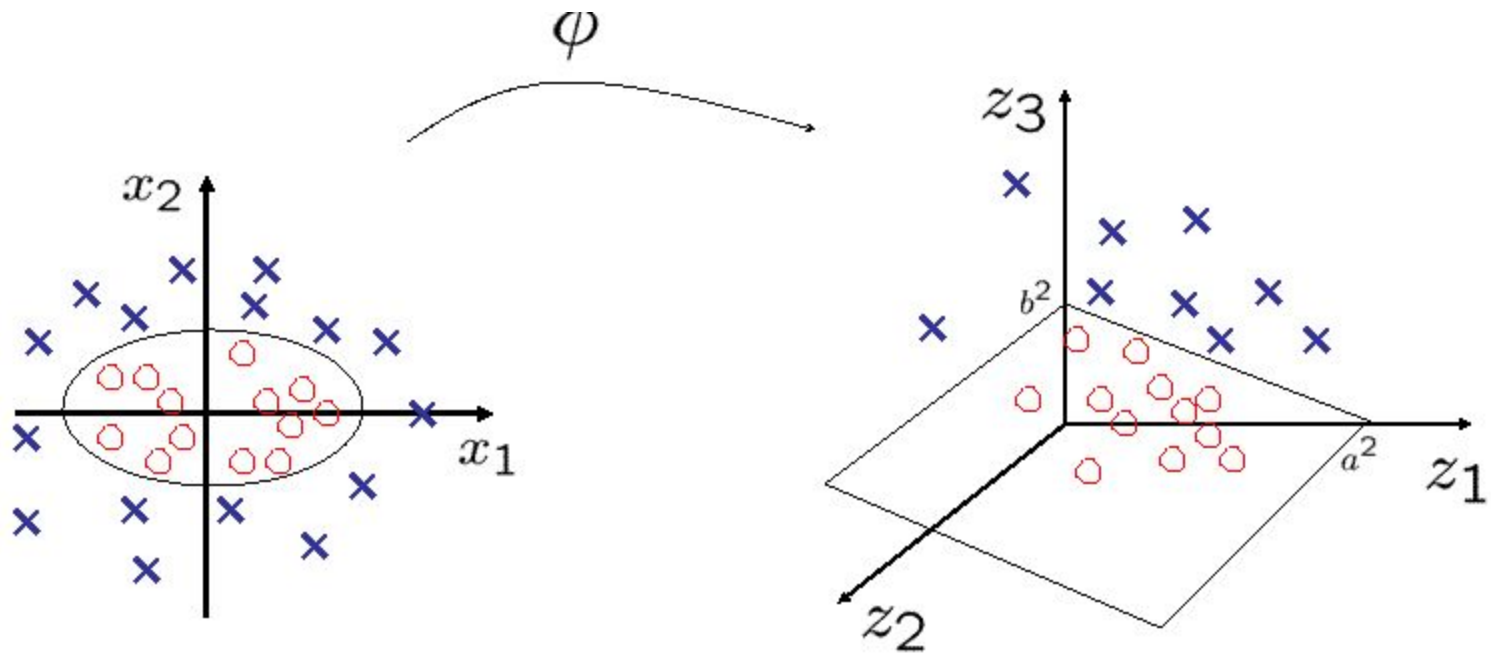
2D

Kernel Trick

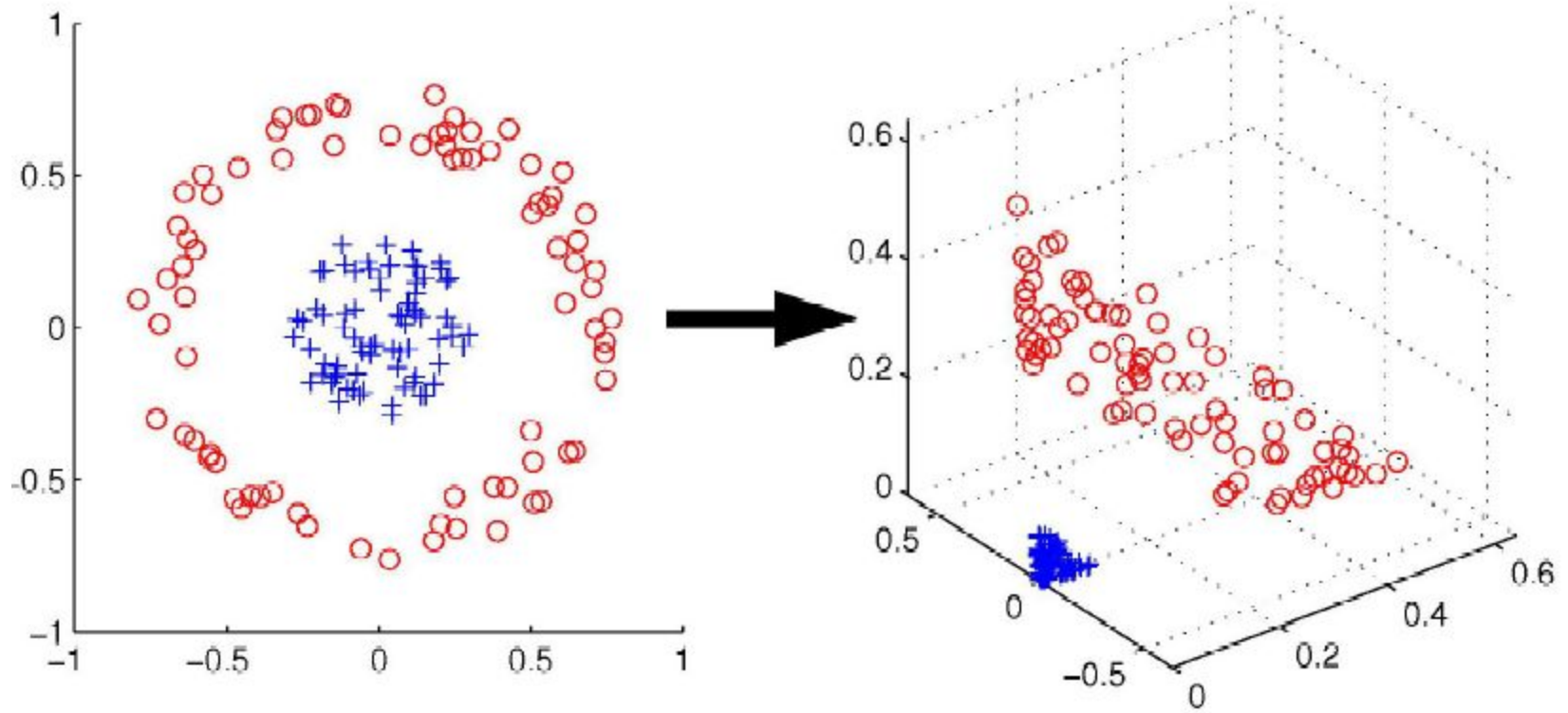
- A função que projeta o espaço de entrada no espaço de características é conhecida como *Kernel*
- Baseado no teorema de Cover
 - Dados no espaço de entrada são transformados (transf. não linear) para o espaço de características, onde são linearmente separáveis.
- O vetor $\phi(x_i)$ representa a “imagem” induzida no espaço de características pelo vetor de entrada

Exemplo

$$\begin{aligned} \phi : \quad \mathbb{R}^2 &\longrightarrow \mathbb{R}^3 \\ (x_1, x_2) &\longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \end{aligned}$$



Exemplo



Kernel Trick

- Permite construir um hiperplano no espaço de característica sem ter que considerar o próprio espaço de características de forma explícita.
- Toda vez que um produto interno entre vetores deve ser calculado, utiliza-se o kernel.
- Uma função de kernel deve satisfazer o teorema de Mercer (1909) para ser válida.
 - Função definida positiva simétrica

Exemplos de Kernel

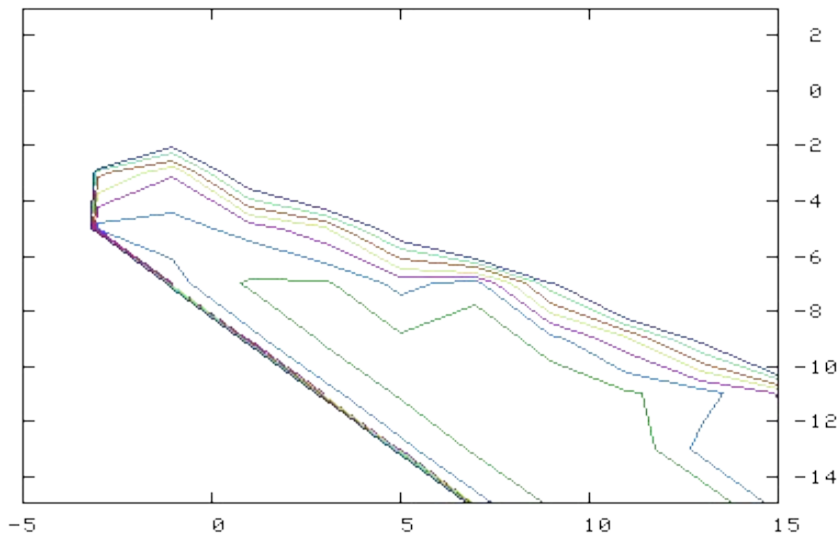
Kernel	Inner Product Kernel
Linear	$K(x, y) = (x \cdot y)$
Gaussian	$K(x, y) = \exp\left(-\frac{\ x - x_i\ ^2}{2\sigma^2}\right)$
Polynomial	$K(x, y) = (x \cdot y)^P$
Tangent Hyperbolic	$K(x, y) = \tanh(x \cdot y - \Theta)$

- *Kernels* complexos
 - Mais parâmetros

Exemples de Kernel

- *Grid Search & RBF (Gaussian)*
 - *Gamma & C (Margin)*

heart scale

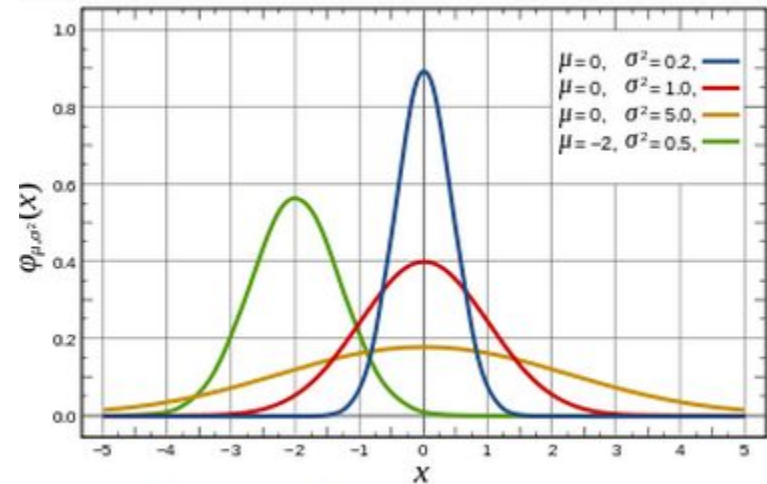


84
83.5
83
82.5
82
81.5
81

$\lg(\gamma)$

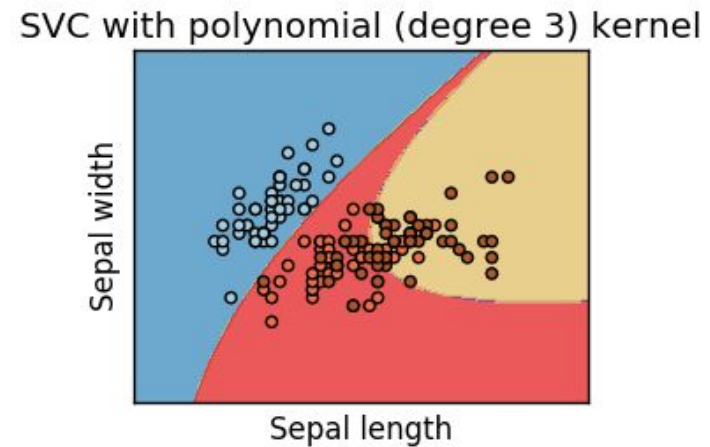
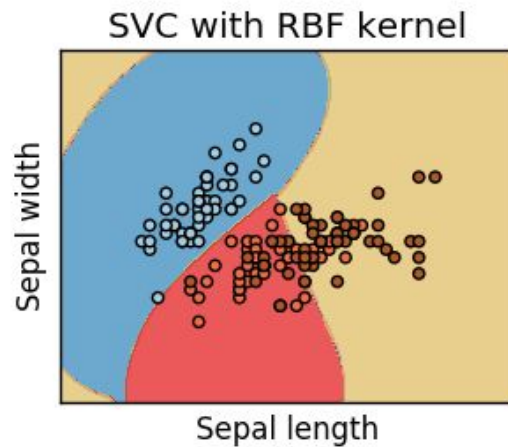
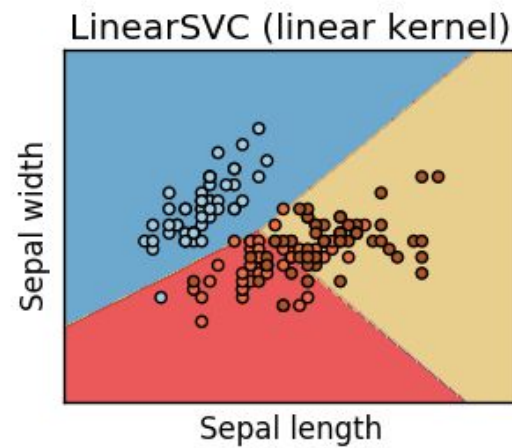
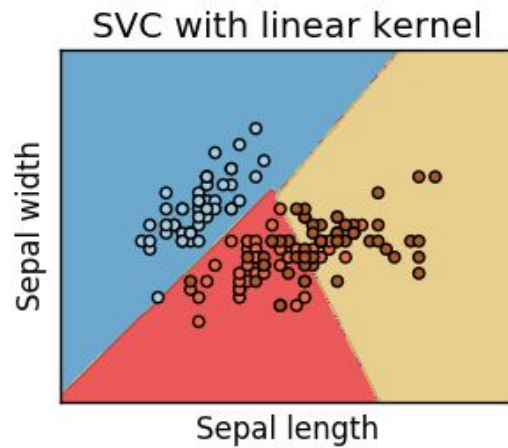
$\lg(C)$

Probability density function



The red curve is the standard normal distribution

Exemplos de Kernel

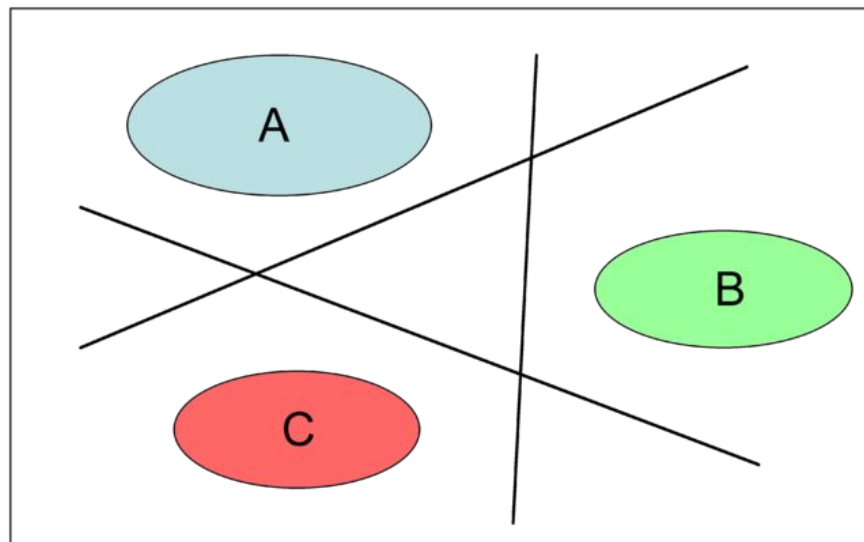


Tomada de Decisão

- SVM são classificadores binários, ou seja, separam duas classes.
- Entretanto, a grande maioria dos problemas reais possuem mais que duas classes.
- Como utilizar os SVMs nesses casos?
 - Um-contra-todos
 - Pairwise

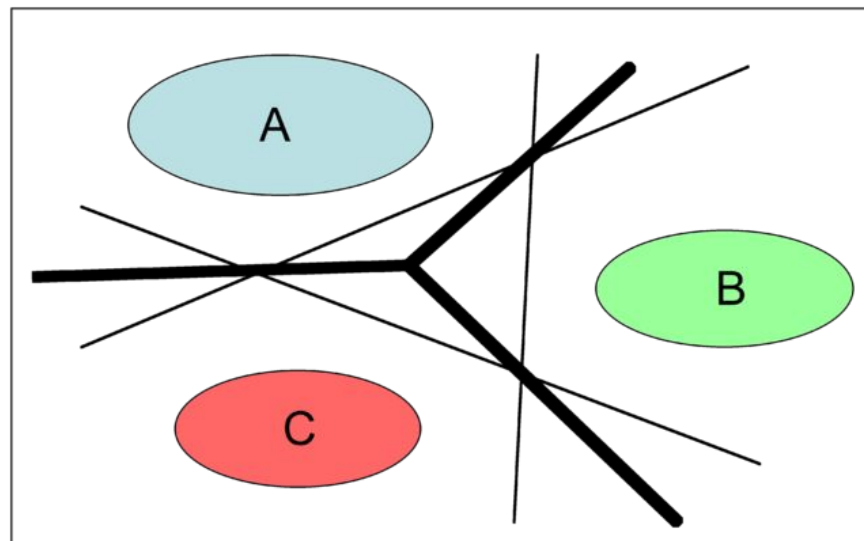
Um-Contra-Todos

- Aqui, o número de classificadores é igual a q .
- Treina-se um classificador c_i para a primeira classe, usando-se como contra exemplos as outras classes, e assim por diante.
- Para se obter a decisão final pode-se utilizar uma estratégia de **votos**.



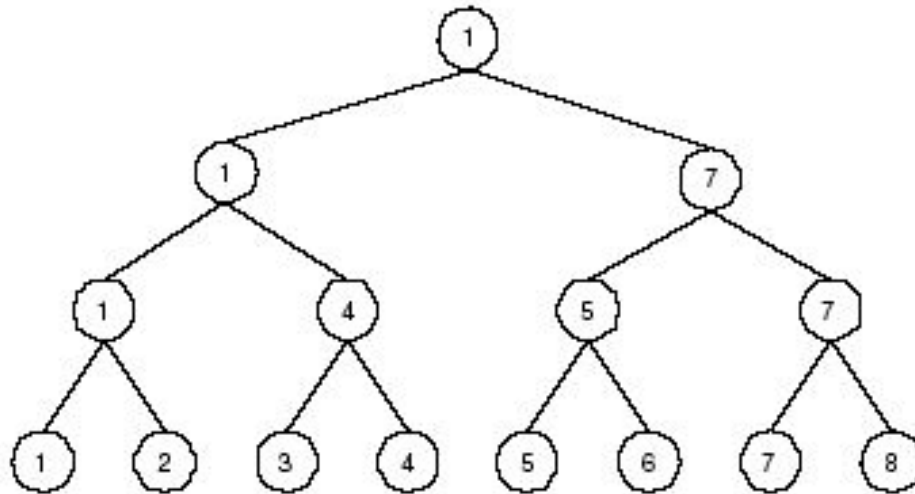
Um-Contra-Todos

- Aqui, o número de classificadores é igual a q .
- Treina-se um classificador c_i para a primeira classe, usando-se como contra exemplos as outras classes, e assim por diante.
- Para se obter a decisão final pode-se utilizar uma estratégia de **votos**.



Pairwise

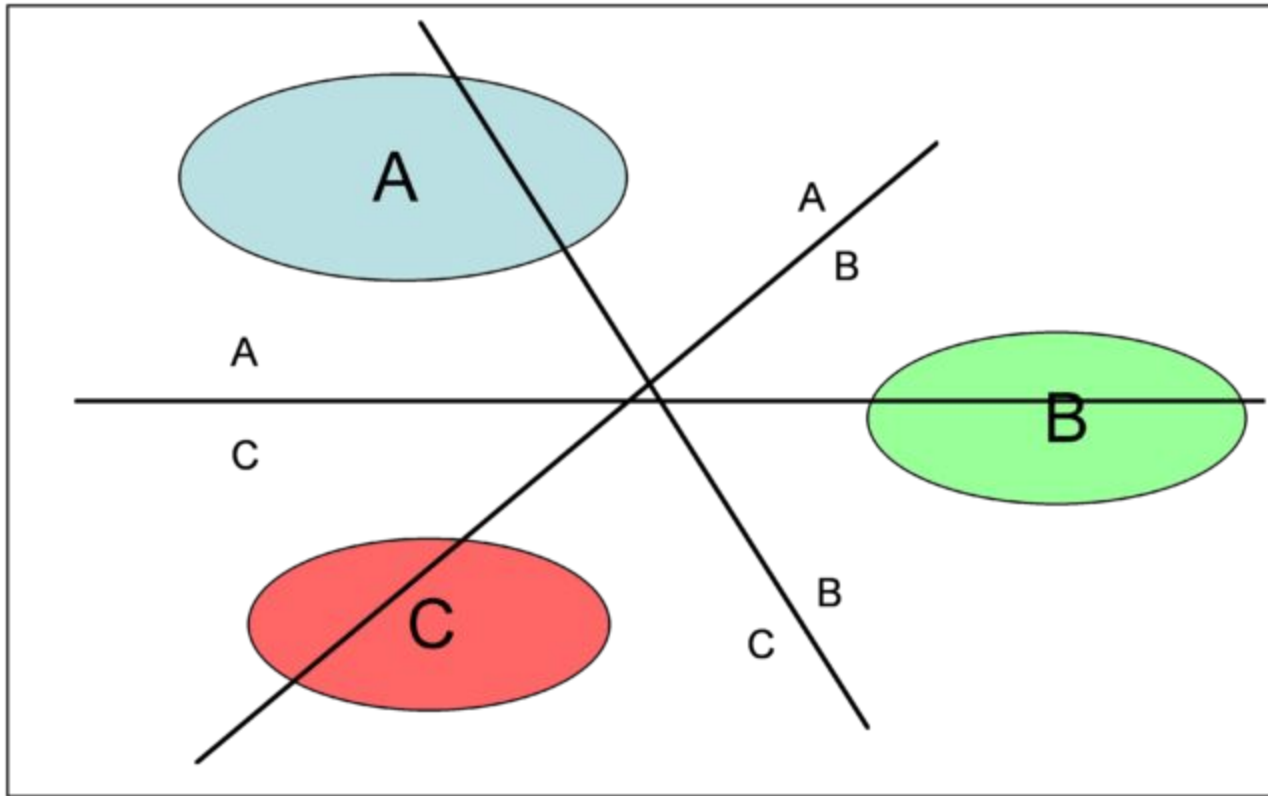
- Consiste em treinar classificadores *pairwise* e arranjá-los em uma árvore



A competição se dá nos níveis inferiores, e o ganhador chegará ao nó principal da árvore.

Número de classificadores para q classes = $q(q-1)/2$.

Pairwise

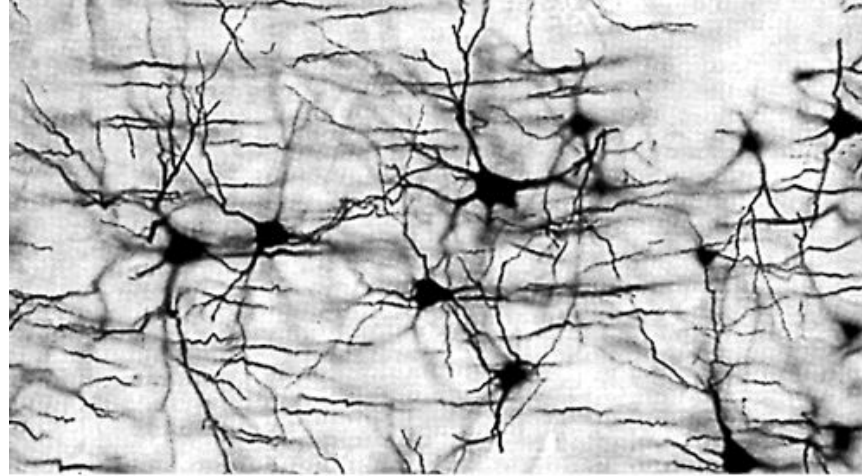


Exercício

- Utilizar a ferramenta LibSVM para realizar classificação usando SVM.

Multiple Layer Perceptron

Redes Neurais



- Cérebro humano.
 - Mais fascinante processador existente.
 - Aprox. 10 bilhões de neurônios conectados através de sinapses.
 - Sinapses transmitem estímulos e o resultado pode ser estendido por todo o corpo humano.

Redes Neuroniais Artificiais: Um breve histórico

- 1943 – McCulloch e Pitts.
 - Sugeriram a construção de uma máquina baseada ou inspirada no cérebro humano.
- 1949 – Donald Hebb.
 - Propõe uma lei de aprendizagem específica para as sinápses dos neurônios.
- 1957/1958 - Frank Rosenblatt.
 - Estudos aprofundados — pai da neuro-computação.
 - Perceptron.
 - Criador do primeiro neuro-computador a obter sucesso (Mark I).

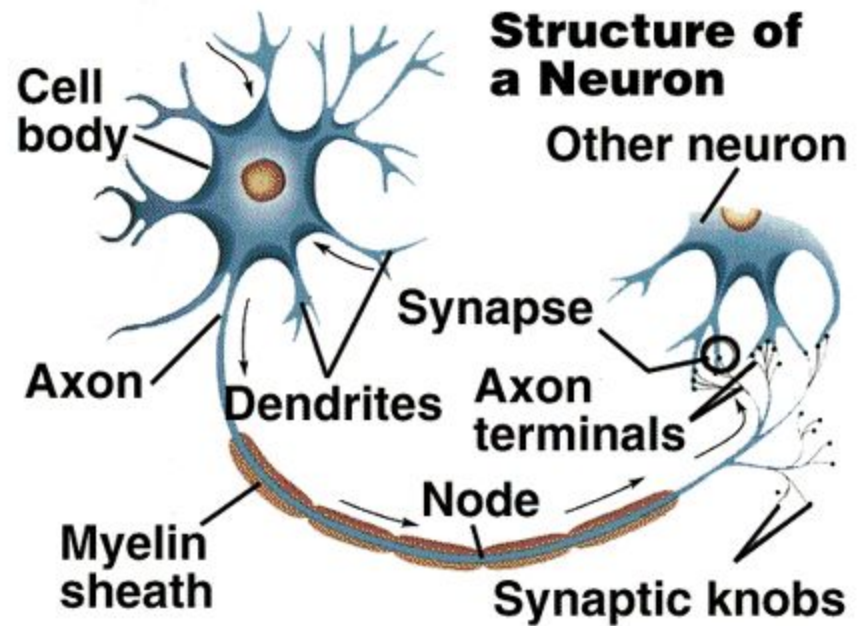
Mark I - Perceptron



Um breve histórico (cont)

- 1958-1967.
 - Várias pesquisas mal sucedidas.
- 1967-1982.
 - Pesquisas silenciosas.
- 1986 – David Rumelhart.
 - Livro “*parallel distributed processing*”.
 - Algoritmo eficaz de aprendizagem.
- 1987.
 - Primeira conferência IEEE Int. Conf. On Neural Nets.

N



- Dendritos:
 - Receber os estímulos transmitidos por outros neuronios.
- Corpo (somma).
 - Coletar e combinar informações vindas de outros neurônios.
- Axônio.
 - Transmite estímulos para outras células.

Redes Neurais Artificiais

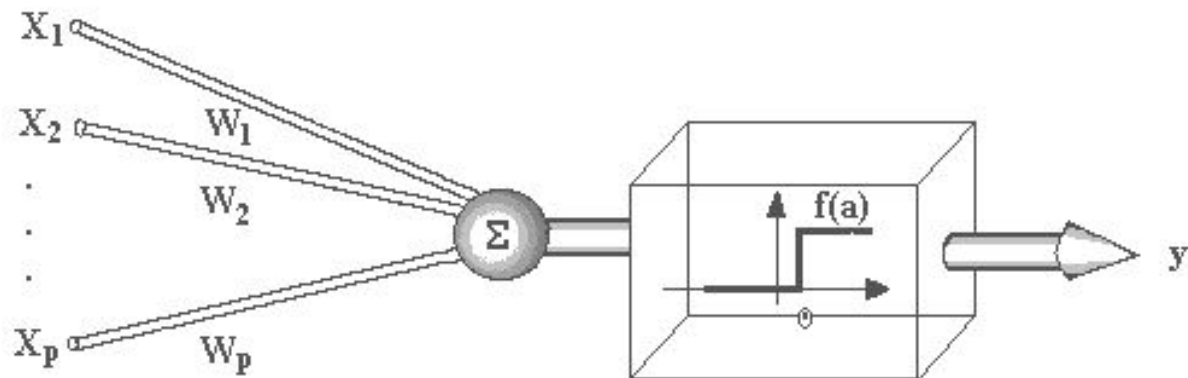
- Técnica computacional que apresenta um modelo inspirado na estrutura do neurônio.
- Simula o cérebro humano:
 - Aprendendo, errando e fazendo descobertas.
- Estrutura de processamento de informação distribuída paralelamente na forma de um grafo direcionado.

Cérebro X Computador

Parâmetro	Cérebro	Computador
Material	Orgânico	Metal e Plástico
Velocidade	Milisegundos	Nanosegundos
Tipo de Processamento	Paralelo	Sequencial
Armazenamento	Adaptativo	Estático
Controle de Processos	Distribuído	Centralizado
Número de Elementos Processados	10^{11} a 10^{14}	10^5 a 10^6
Ligações entre Elementos Processados	10.000	< 10

Neurônio artificial

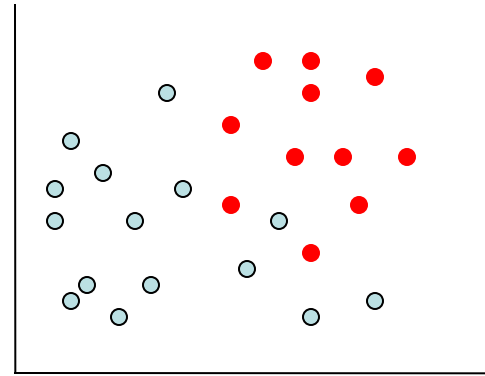
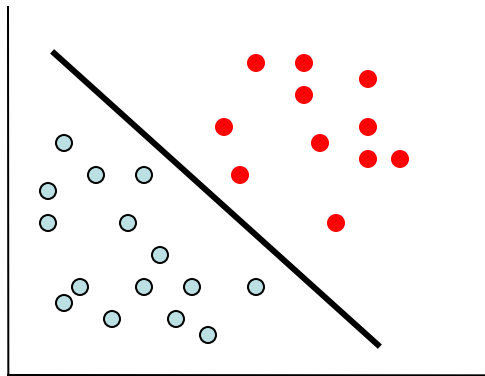
- Sinais são apresentados à entrada.
- Cada sinal é multiplicado por um peso.
- Soma ponderada produz um nível de ativação.
- Se esse nível excede um limite (*threshold*) a unidade produz uma saída.



[PERCEPTRON]

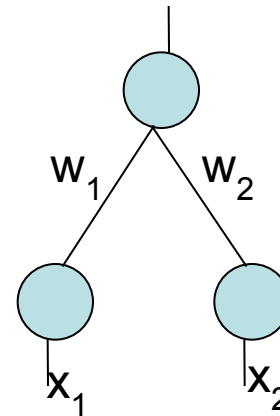
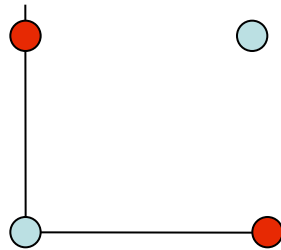
Perceptron

- Resolve problemas linearmente separáveis somente.
- Problemas reais, entretanto, na maioria das vezes são mais complexos.



Exemplo

- Considere por exemplo, o problema XOR



Como visto anteriormente com SVMs, podemos afirmar que em altas dimensões os problemas são linearmente separáveis.

Sendo assim, vamos mudar o problema de \mathbb{R}^2 para \mathbb{R}^3 adicionando uma terceira característica.

Exemplo

X1	X2	X3	Output
1	1	1	0
1	0	0	1
0	1	0	1
0	0	0	0

- Nesse caso, a característica adicionada (X_3) é a operação AND entre X_1 e X_2
- O fato de adicionarmos essa característica faz com que o problema torne-se linearmente separável.

Adicionando uma camada

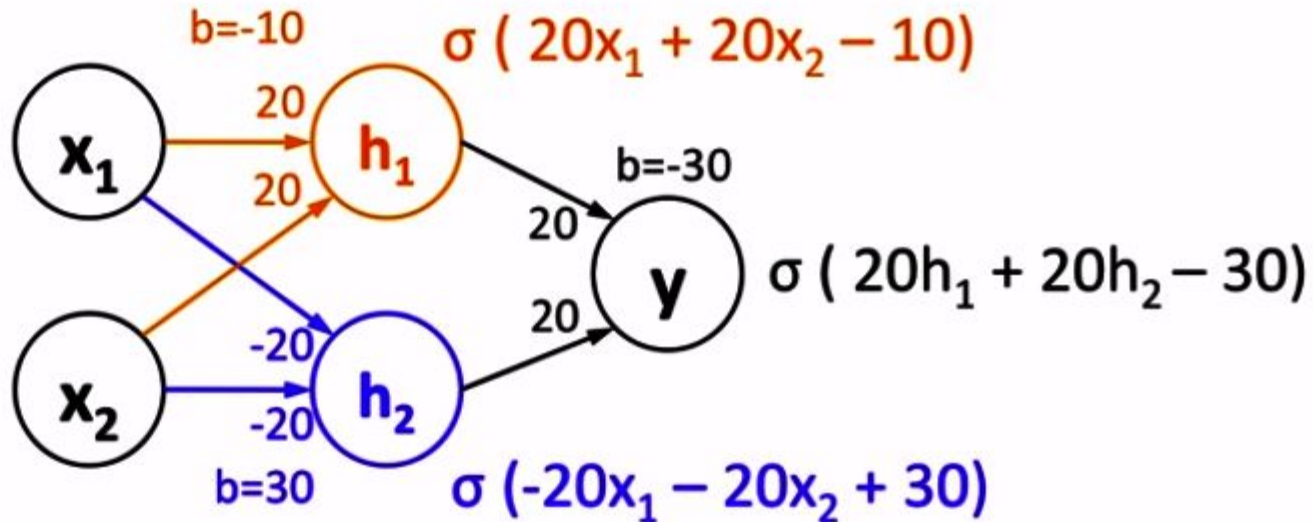
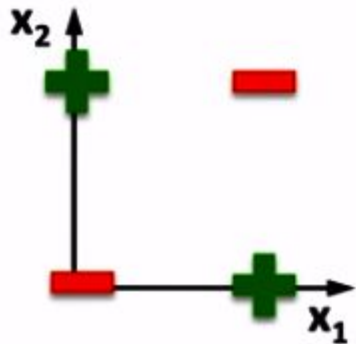
- Outra maneira de resolver esse problema consiste em adicionar uma camada extra (**escondida**) entre as camadas de entrada e saída.
- Dado uma quantidade suficiente de neurônios na camada escondida, é possível resolver qualquer tipo de problemas
 - Claro que isso depende das características de entrada.

Camada Escondida

- Essa camada pode ser vista como um extrator de características, ou seja, a grosso modo, o neurônio escondido seria uma característica a mais
 - O que torna o problema do XOR linearmente separável.

Uma outra rede: XOR

Linear classifiers cannot solve this



$$\sigma(20 \cdot 0 + 20 \cdot 0 - 10) \approx 0$$

$$\sigma(20 \cdot 1 + 20 \cdot 1 - 10) \approx 1$$

$$\sigma(20 \cdot 0 + 20 \cdot 1 - 10) \approx 1$$

$$\sigma(20 \cdot 1 + 20 \cdot 0 - 10) \approx 1$$

$$\sigma(-20 \cdot 0 - 20 \cdot 0 + 30) \approx 1$$

$$\sigma(-20 \cdot 1 - 20 \cdot 1 + 30) \approx 0$$

$$\sigma(-20 \cdot 0 - 20 \cdot 1 + 30) \approx 1$$

$$\sigma(-20 \cdot 1 - 20 \cdot 0 + 30) \approx 1$$

$$\sigma(20 \cdot 0 + 20 \cdot 1 - 30) \approx 0$$

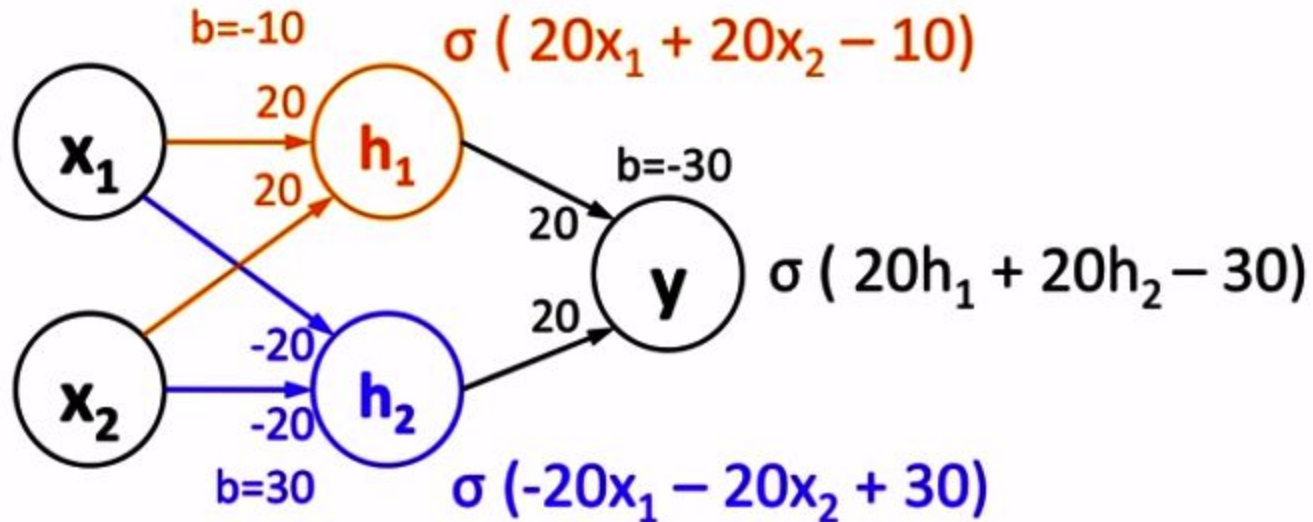
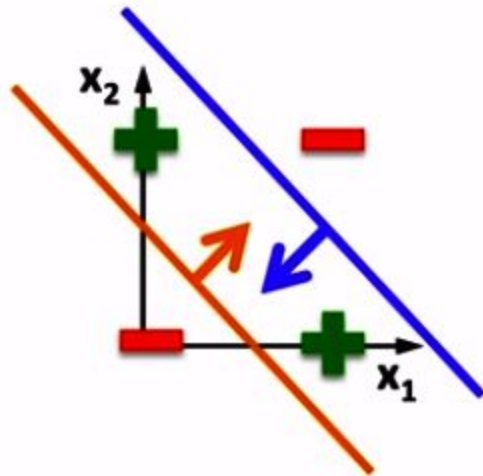
$$\sigma(20 \cdot 1 + 20 \cdot 0 - 30) \approx 0$$

$$\sigma(20 \cdot 1 + 20 \cdot 1 - 30) \approx 1$$

$$\sigma(20 \cdot 1 + 20 \cdot 1 - 30) \approx 1$$

Uma outra rede: XOR

Linear classifiers cannot solve this

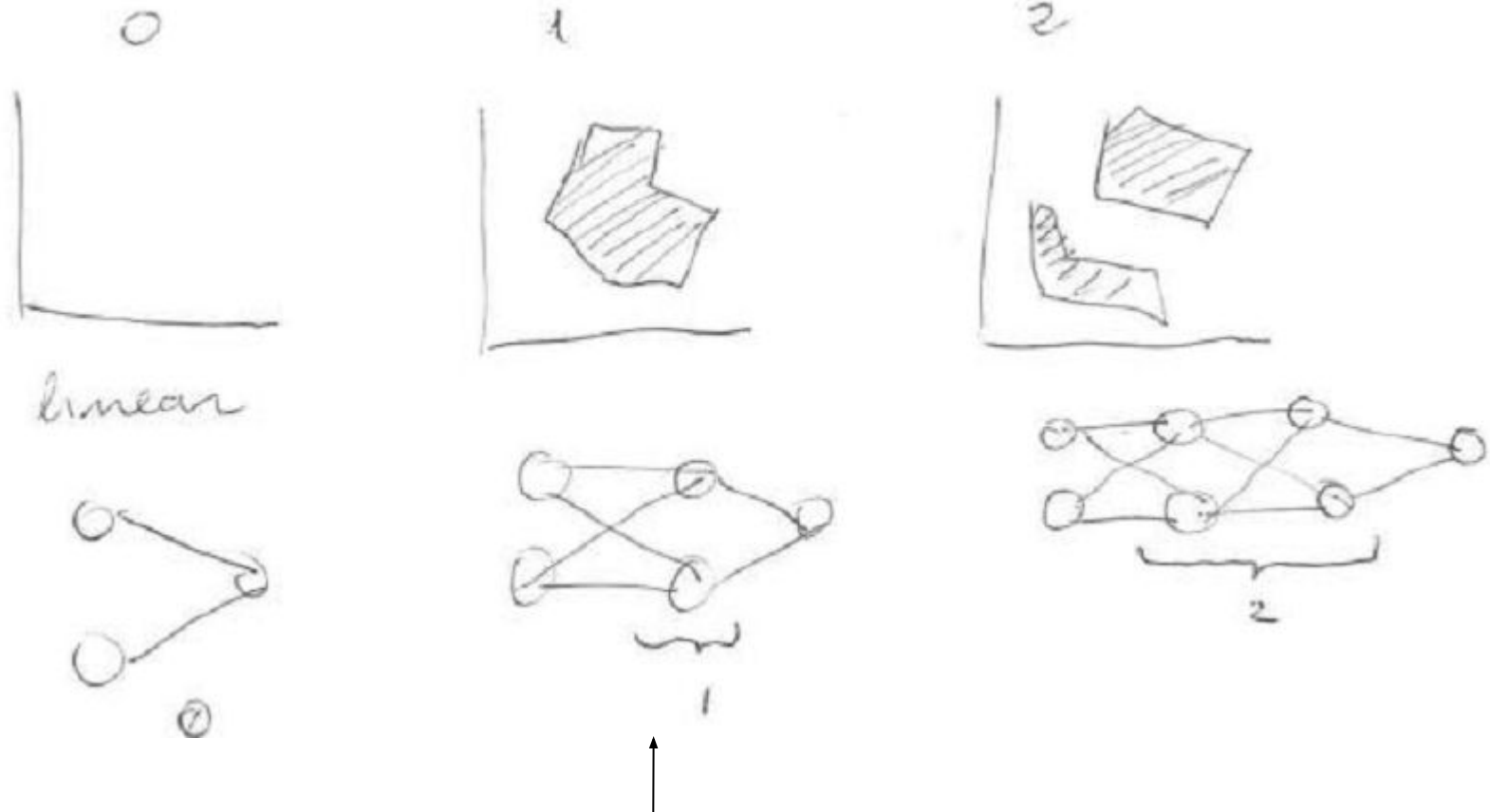


$$\begin{aligned} \sigma(20 \cdot 0 + 20 \cdot 0 - 10) &\approx 0 \\ \sigma(20 \cdot 1 + 20 \cdot 1 - 10) &\approx 1 \\ \sigma(20 \cdot 0 + 20 \cdot 1 - 10) &\approx 1 \\ \sigma(20 \cdot 1 + 20 \cdot 0 - 10) &\approx 1 \end{aligned}$$

$$\begin{aligned} \sigma(-20 \cdot 0 - 20 \cdot 0 + 30) &\approx 1 \\ \sigma(-20 \cdot 1 - 20 \cdot 1 + 30) &\approx 0 \\ \sigma(-20 \cdot 0 - 20 \cdot 1 + 30) &\approx 1 \\ \sigma(-20 \cdot 1 - 20 \cdot 0 + 30) &\approx 1 \end{aligned}$$

$$\begin{aligned} \sigma(20 \cdot 0 + 20 \cdot 1 - 30) &\approx 0 \\ \sigma(20 \cdot 1 + 20 \cdot 0 - 30) &\approx 0 \\ \sigma(20 \cdot 1 + 20 \cdot 1 - 30) &\approx 1 \\ \sigma(20 \cdot 1 + 20 \cdot 1 - 30) &\approx 1 \end{aligned}$$

Camadas X Fronteiras



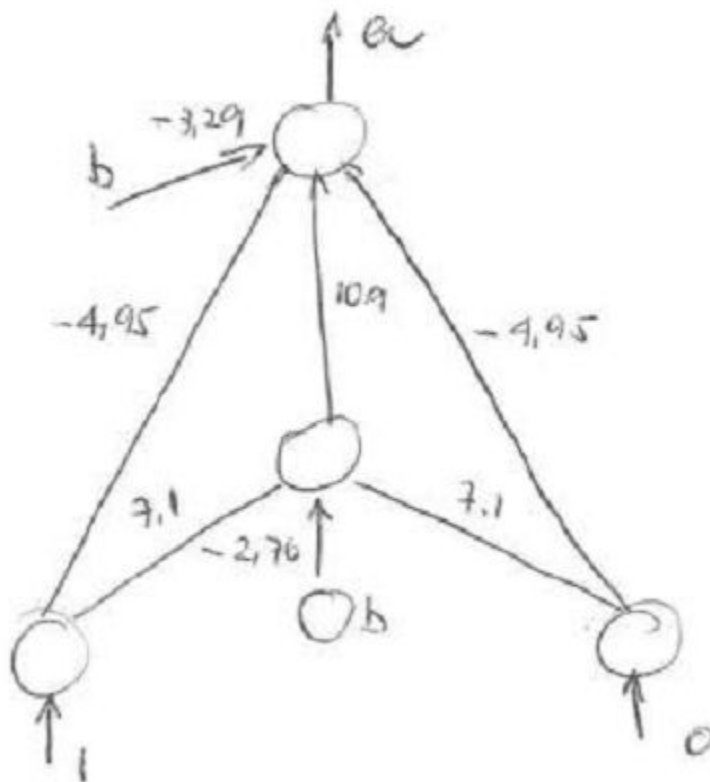
Geralmente uma camada escondida resolve a grande maioria dos problemas.

O problema de atribuição de créditos

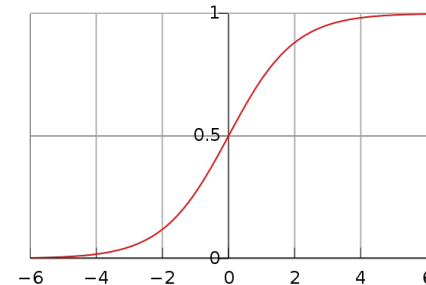
- Quando temos uma camada escondida, surge o problema de atribuição de créditos aos neurônios desta camada
 - Não existem “*targets*” como na camada de saída.
 - Período entre 1958 e 1982 foi dedicado a resolver esse problema
 - Solução foi o algoritmo conhecido como Backpropagation.
 - David E. Rumelhart et al (1986)

MLP para o problema XOR

- Considere o problema XOR e a seguinte rede já treinada.



A saída é calculada de maneira similar ao perceptron, mas a MLP geralmente usa **sigmoid** como função de ativação.



$$v = \frac{1}{1 + e^{-s}}$$

Para grandes quantidades de dados o resultado da sigmoid pode ser interpretado como estimativa de probabilidade.

MLP para o problema XOR

- Para calcular a saída da rede, primeiramente devemos encontrar o valor do neurônio escondido.

$$1 * 7.1 + 1 * -2.76 + 0 * 7.1 = 4.34$$

- Passando o valor 4.34 na função de ativação,
 - temos 0.987.

- O valor de saída é então calculado

$$1 * -4.95 + 0 * -4.95 + 0.987 * 10.9 + 1 * -3.29 = 2.52$$

- Passando 2.52 na função de ativação,
 - temos a saída igual a 0.91

MLP para o problema XOR

- Após classificarmos todos os padrões de entrada teríamos

1	0	0.91	←
0	0	0.08	
0	1	1.00	←
1	1	0.10	

Podemos usar como regra de decisão um limiar igual a 0.9, por exemplo.

BackProp

- Seja o_j o valor de ativação para o neurônio j .
- Seja f uma função de ativação.

$$net_j = \sum_{i=1}^n \omega_{ij} o_i$$

- Seja w_{ij} o peso entre os neurônios i e j .
- Seja net_j a entrada para o neurônio j , a qual é calculada por

$$o_j = f(net_j)$$

onde n é o número de unidades ligadas ao neurônio j

BackProp

- O treinamento acontece da seguinte maneira:
 1. Inicializar os valores dos pesos e neurônios aleatoriamente.
 2. Apresentar um padrão a camada de entrada da rede
 3. Encontrar os valores para as camadas escondidas e a camada de saída.
 4. Encontrar o erro da camada de saída.
 5. Ajustar os pesos através da retropropagação dos erros (*Backpropagation*)
 1. Diminuir o erro a cada iteração
 6. Encontrar o erro na camada escondida
 7. Ajustar os pesos.

Critério de parada é o erro desejado

BackProp

- O valor corrente de ativação de um neurônio k é o_k e o target é t_k
- Após realizar os passos 1, 2, e 3, o próximo passo consiste em calcular o erro, o qual pode ser realizado através da seguinte equação

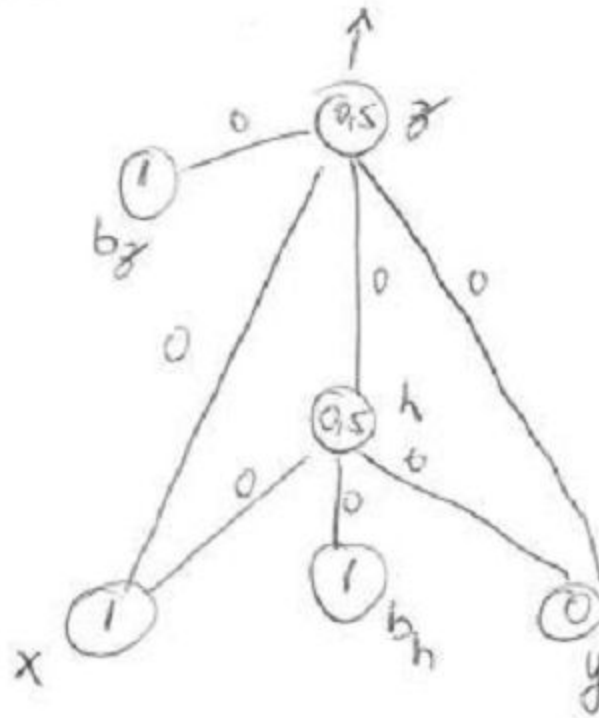
$$\delta_k = (t_k - o_k) o_k (1 - o_k)$$

- Observe que:

$$f(f - 1) = f', \text{ se } f(x) = \frac{1}{1 - e^{-x}}$$

BackProp: Exemplo

- Considere uma rede inicializada da seguinte forma



Nesse caso $\delta_z = (1 - 0.5) \times 0.5 \times (1 - 0.5) = 0.125$

BackProp: Exemplo

- A fórmula para atualizar os pesos W entre o neurônio i e j é

$$\omega_{ij} = \omega_{ij} + \eta \delta_j o_i$$

- onde **eta** é uma constante pequena e positiva chamada de “taxa de aprendizagem” (*learning rate*)
- Usando uma taxa de 0.1, temos

BackProp: Exemplo

$$W_{xz} = 0 + 0.1 \times 0.125 \times 1 = 0.0125$$

$$W_{yz} = 0 + 0.1 \times 0.125 \times 0 = 0$$

$$W_{hz} = 0 + 0.1 \times 0.125 \times 0.5 = 0.00625$$

$$W_{bz} = 0 + 0.1 \times 0.125 \times 1 = 0.0125$$

A fórmula para calcular o erro dos neurônios da camada escondida é

$$\delta_i = o_i(1 - o_i) \sum_k \delta_k w_{ik}$$

Como temos somente um neurônio no nosso exemplo

$$\delta_h = o_h(1 - o_h) \delta_z w_{hz}$$

BackProp: Exemplo

Nesse caso teríamos

$$\delta_h = 0.5(1 - 0.5)0.125 \times 0.00625 = 0.000195313$$

Para atualizar os pesos, usamos a mesma fórmula

$$W_{hx} = 0 + 0.1 \times 0.0001953 \times 1 = 0.00001953$$

$$W_{hy} = 0 + 0.1 \times 0.0001953 \times 0 = 0$$

$$W_{hbh} = 0 + 0.1 \times 0.0001953 \times 1 = 0.00001953$$

Com os novos pesos calculados, a saída da rede seria 0.507031

BackProp: Exemplo

- Após aplicarmos todos os padrões, teríamos o seguinte

1	0	0.4998
0	0	0.4998
0	1	0.4998
1	1	0.4997

- Usando uma taxa de aprendizagem = **0,1** ,
 - o algoritmo levará **20.000** iterações para convergir.
- Uma solução para melhorar isso seria
 - aumentar a **learning rate**.
- Se usarmos *learning rate* = **2**,
 - o algoritmo converge em **480** iterações.

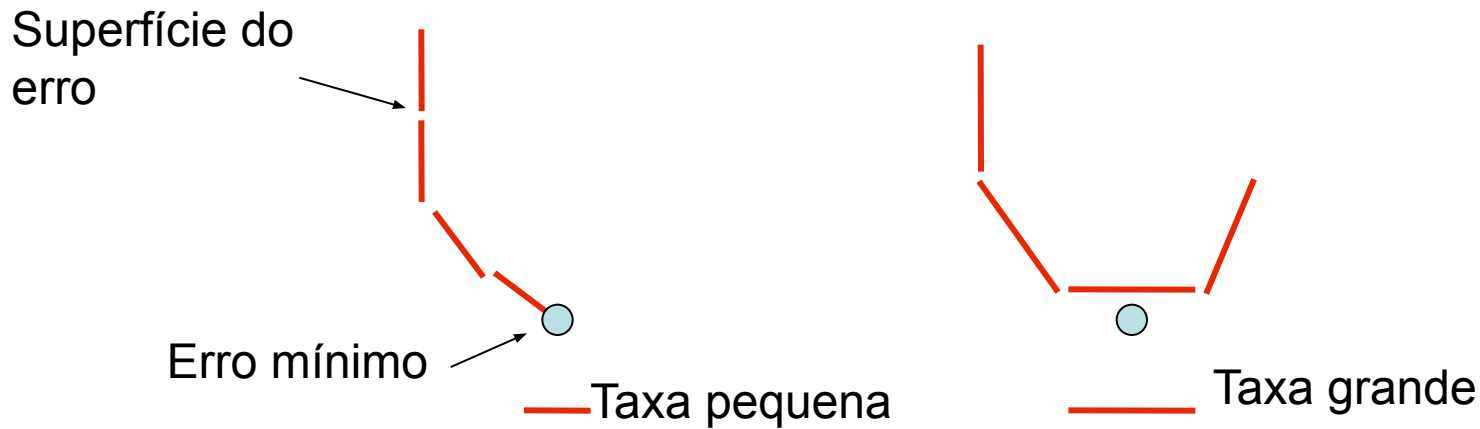
Aspectos Práticos

- Alguns aspectos práticos devem ser considerados na utilização de redes neuronais **MLP**.
 - Taxa de aprendizagem
 - Momentum
 - Online vs batch
 - Shuffle
 - Normalização
 - Inicialização dos pesos
 - Generalização

Y. LeCun et al,
Efficient Backprop, 1998

Taxa de Aprendizagem

- Taxas muito pequenas tornam o processo bastante lento.
- Taxas muito grandes tornam o processo rápido.
 - Podem não trazer os resultados ideais.

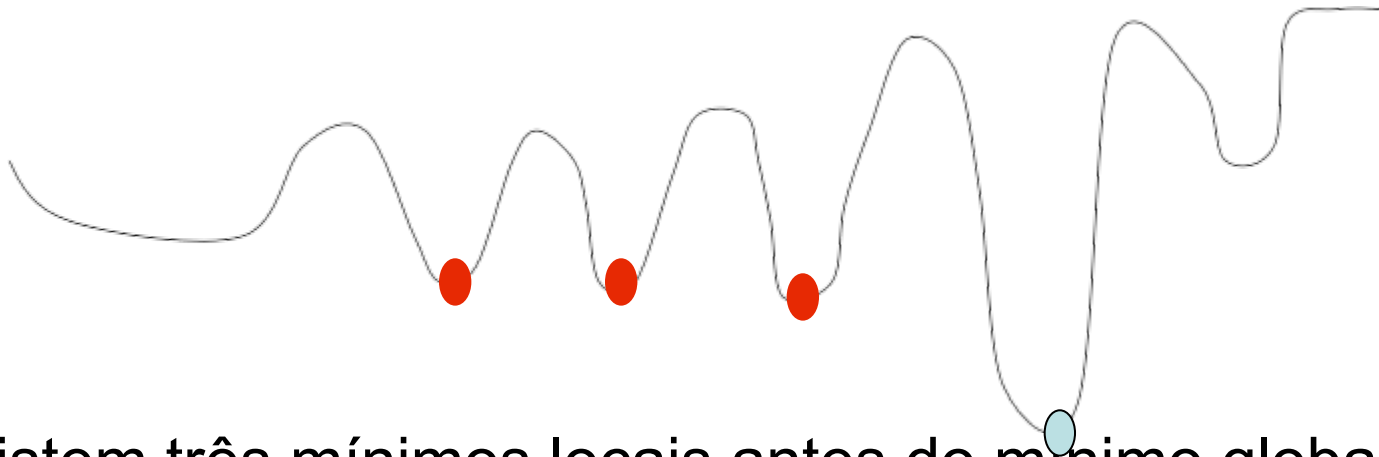


Taxa de Aprendizagem

- O ideal é começar com uma taxa grande e reduzir durante as iterações.
- Permite a exploração global no início (*exploration*) a local (*exploitation*) quando o algoritmo estiver próximo do ótimo global.
- Geralmente valores entre 0.05 e 0.75 fornecem bons resultados.

Momentum

- É uma estratégia usada para evitar mínimos locais.
- Considere a seguinte superfície



- Existem três mínimos locais antes do mínimo global.

Momentum

- Considere, por exemplo, que uma bola seja solta no início da superfície.
- Se ela tiver “*momentum*” suficiente, ela vai conseguir passar os três mínimos locais e alcançar o mínimo global.

$$\omega_{ij} = \alpha\omega_{ij} + \eta\delta_j o_k$$

- Normalmente $0 \leq \alpha \leq 0.9$

On-line vs Batch

- A diferença está no momento em que os pesos são atualizados.
- Na versão *on-line*, os pesos são atualizados a cada padrão apresentado a rede.
- Na versão *batch*, todos os pesos são somados durante uma iteração/época (todos os padrões) e só então os pesos são atualizados.
- *Versão batch*
 - Interessante para aplicações que possam ser paralelizadas.
- *Versão on-line*
 - Geralmente converge mais rapidamente.

Misturando Exemplos (*Shuffle*)

- Redes neuronais aprendem melhor quando diferentes exemplos de diferentes classes são apresentados a rede.
- Uma prática muito comum consiste em apresentar um exemplo de cada classe a rede
 - Isso garante que os pesos serão atualizados levando-se em consideração todas as classes.

Misturando Exemplos (*Shuffle*)

- Se apresentarmos à rede todos os exemplos de uma classe, e assim por diante, os pesos finais tenderão para a última classe
 - A rede vai “esquecer” o que ela aprendeu antes.

Normalização

- A normalização é interessante quando existem características em diversas unidades dentro do vetor de características.
- Nesses casos, valores muito altos podem saturar a função de ativação.
- **Soma 1**: Uma maneira bastante simples de normalizar os dados consiste em somar todas as características e dividir pela soma

$$x'_i = \frac{x_i}{\sum_i x_i}$$

- Garante que a “energia” inserida na rede seja 1

Normalização

- **Z-score**: Para redes neuronais MLP, geralmente é interessante ter as características com média próxima de zero

$$x'_i = \frac{x_i - \mu(X_i)}{\sigma(X_i)}$$

- Melhora o tempo de convergência durante a aprendizagem.

Normalização

- **Max-min**: Pode-se também normalizar cada característica para ficar entre 0 e 1

$$x'_i = \frac{\max(X_i) - x_i}{\max(X_i) - \min(X_i)}$$

- Pode “sobrecarregar” a rede

Normalização

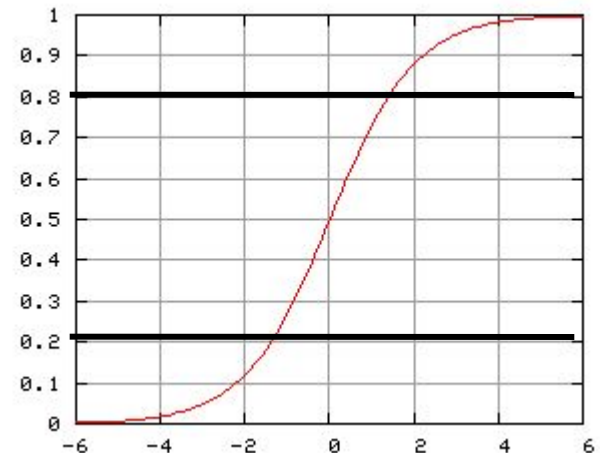
- As características devem ser não correlacionadas se possível
 - Quando temos poucas características podemos verificar isso facilmente.
 - Com várias características, o problema se torna muito mais complexo.
 - Métodos de seleção de características para escolher características descorrelacionadas

Inicialização dos Pesos

- Assumindo que os dados foram normalizados e uma *sigmoid* está sendo usada como função de ativação.
 - Os **pesos** (espaço de separação) devem ser inicializados aleatoriamente mas de modo que fiquem na região linear da sigmoid

Com pesos muito altos ou muito baixo a sigmoid deve saturar

- Gradientes pequenos
- Aprendizagem muito lenta.

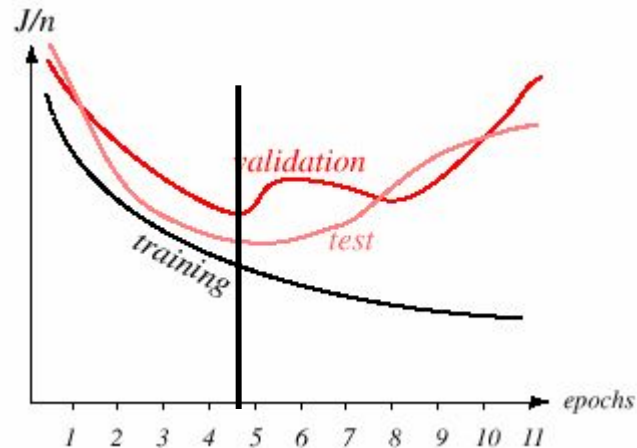


Generalização

- Um aspecto bastante importante quando treinamos um classificador é garantir que o mesmo tenha um bom poder de generalização.
 - Evitar *overfitting*.
- A maneira clássica de se garantir uma boa generalização consiste em reservar uma parte da base para **validar** a generalização.

Generalização

- A cada iteração, devemos monitorar o desempenho na base de validação.
- Não é raro observar o seguinte desempenho

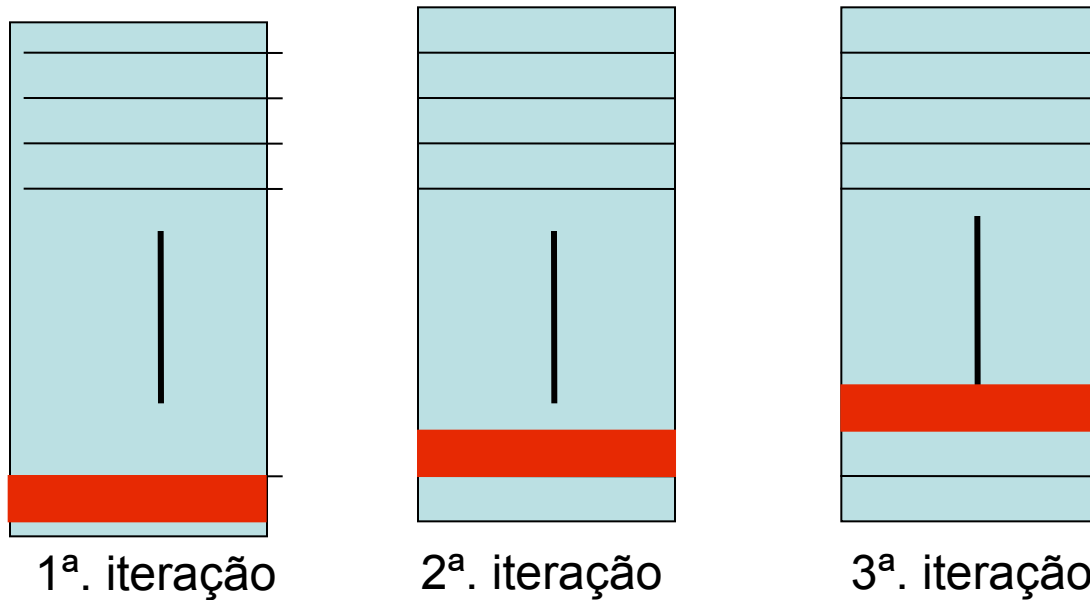


Generalização

- Uma outra estratégia é a validação cruzada.
 - Interessante quando a base não é muito grande
 - Separar alguns exemplos para validação pode prejudicar o treinamento.
- Consiste em dividir a base de aprendizagem em n partições iguais e usar $n-1$ partições para aprendizagem e uma partição para validação.

Generalização

- A cada iteração de aprendizado a partição usada para validação é trocada.



Tamanho da Rede

- Geralmente uma camada escondida é suficiente.
- Em poucos casos você vai precisar adicionar uma segunda camada escondida.
- Não existe uma fórmula matemática para se encontrar o número de neurônios.
 - Empírico
- Dica prática
 - Comece com uma rede menor, pois a aprendizagem vai ser mais rápida.

Referências

- Luiz E. S Oliviera, **Classificadores Lineares**, DInf / UFPR, 2014.
- Victor Lavrenko, **Introductory Applied Machine Learning**, University of Edinburgo, 2014.