

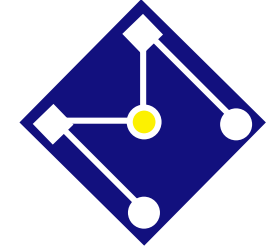


PMR5251 - Avaliação do Comportamento Mecânico de Materiais Utilizando uma Abordagem de ML



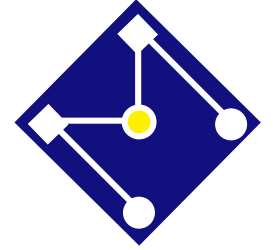
MACHINE LEARNING: PROBLEMAS DE CLASSIFICAÇÃO

Izabel F. Machado
Larissa Driemeier



NOSSAS AULAS

Data	Assunto	Conteúdo principal
01/10	Introdução ao Aprendizado de Máquinas	Teoria conceitual
15/10	Redes Neurais	Teoria e Prática
22/10	Regressão	Teoria e Prática
9/11	Classificação	Teoria e Prática



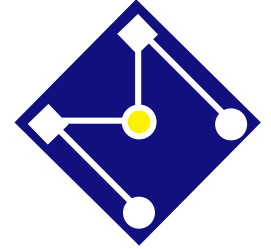
AULA DE HOJE

- Regressão Logística
- Exemplo de aplicação
- Regressão usando redes neurais



REGRESSÃO LOGÍSTICA

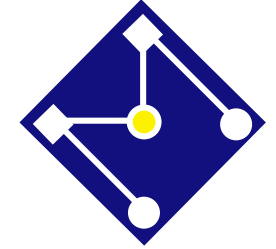
O que é isso?



TÉCNICA DE REGRESSÃO LOGÍSTICA

E se,

- um engenheiro deseja saber se ocorrerá ou não a falha ou não de uma estrutura, em função do estado de tensões, material e geometria?
- Um médico deseja investigar se a probabilidade de ataque cardíaco pode ser predita em função de características sanguíneas, sexo, estilo de vida?
- Uma operadora de telefonia móvel pode querer saber a probabilidade de mudança de plano por parte dos clientes que compõem sua carteira, em função de características como nível de escolaridade, renda, estado civil, número de filhos, tempo de relacionamento com a operadora?
- Um pesquisador quer saber se o passo de seu exoesqueleto é estável, a partir dos movimentos das juntas?
- ...



REGRESSÃO VS CLASSIFICAÇÃO

Um problema de regressão tem um número real como saída.

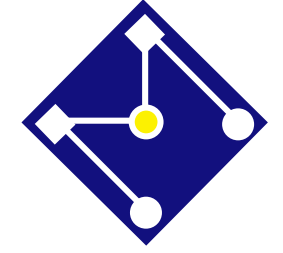
Por exemplo, podemos usar os dados da tabela para estimar o peso de alguém de acordo com sua altura.

Altura (cm)	Peso(kg)
167,1	51,3
181,7	61,9
176,3	69,4
173,3	64,6
172,2	65,5
174,5	55,9
177,3	64,2
177,8	61,9
172,5	51,0
169,6	54,7
168,9	57,8
171,8	51,8
173,5	57,0
170,5	55,5
173,4	52,7

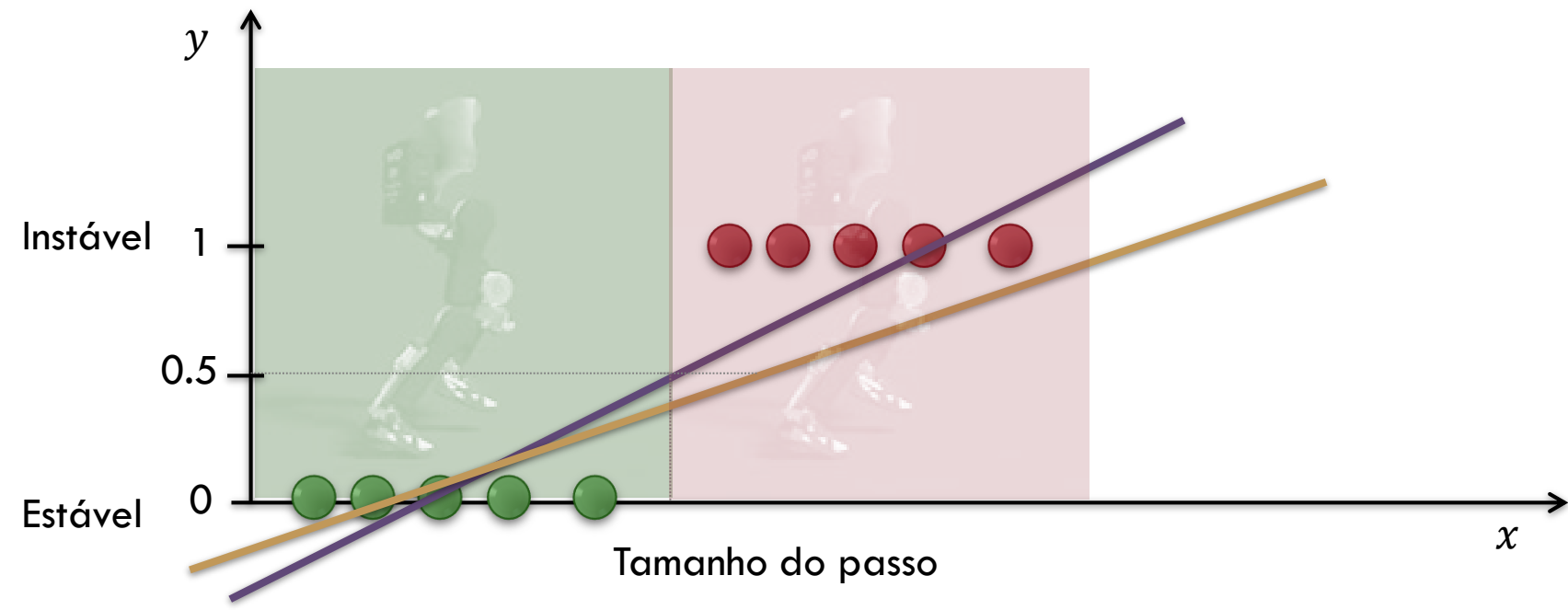
Um problema de classificação tem um valor discreto como saída.

Por exemplo, “gosta de abacaxi na pizza” e “não gosta de abacaxi na pizza” são opções discretas. Não há meio termo.

Idade	Gosta de abacaxi na pizza
42	1
65	1
50	1
76	1
96	1
50	1
91	0
58	1
25	1
23	1
75	1
46	0
87	0
96	0
45	0
32	1
63	0
21	1
26	1
93	0
68	1
96	0

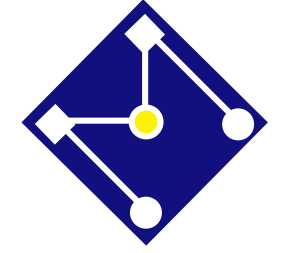


ESTABILIDADE DE MARCHA

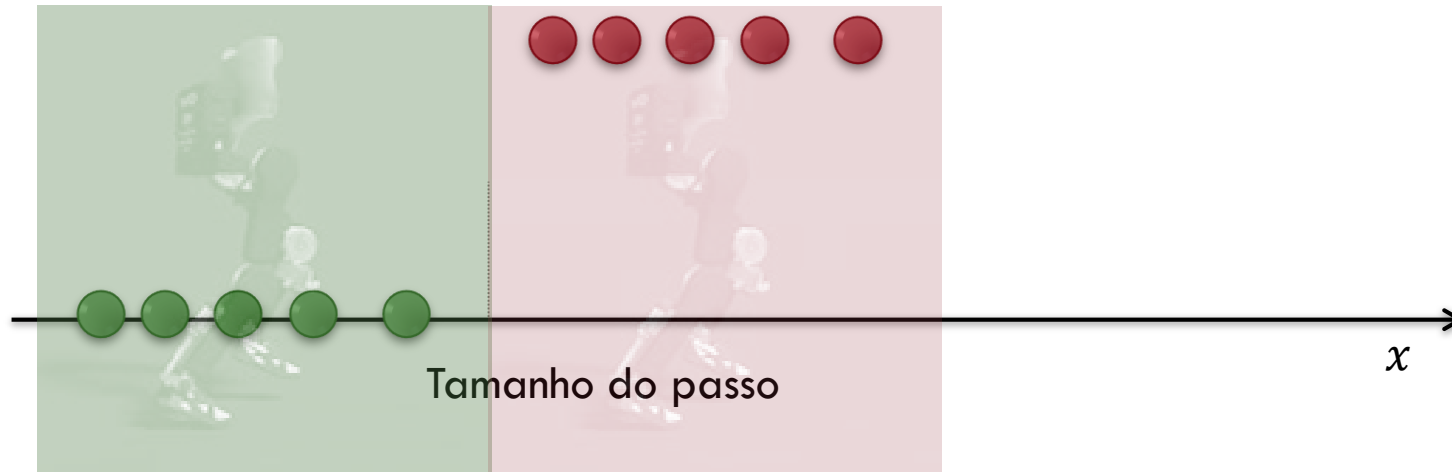


O resultado, y_i , assume o valor 1 (em nosso caso, isso representa uma marcha instável) com probabilidade p_i e o valor 0 com probabilidade $1 - p_i$.

É a probabilidade p_i que modelamos em relação às variáveis independentes.



ESTABILIDADE DE MARCHA



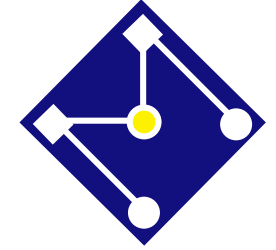
$$h(x) = \hat{y} = \omega^T x$$

$$(p_i) = \omega_0 + \omega_1 x_1^{(i)} + \omega_2 x_2^{(i)} + \dots + \omega_n x_n^{(i)}$$

transformação

Podemos formular o problema da seguinte forma: “o passo é instável?” ou, melhor ainda, “qual a **probabilidade do passo ser instável?**”. Teoricamente, **passos instáveis** deveriam ter uma probabilidade de **1.0** (de serem instáveis), ao passo que **passos estáveis** deveriam ter uma probabilidade de **0.0** (de serem instáveis).

Assim, **passos instáveis** pertencem à **classe positiva (SIM, 1**, eles são instáveis), e **passos estáveis pertencem à classe negativa (NÃO, 0**, eles não são instáveis).



RAZÃO DE CHANCES (ODDS RATIO)

A razão de chances descreve a relação entre a probabilidade de sucesso e de falha. Se um evento ocorre com probabilidade p , a razão de chances deste evento é

$$OR = \frac{p}{1 - p}$$

para 1.

	Não compra	Compra	Total
Mulher	106	159	265
Homem	125	121	246

$$OR_M = \frac{159/265}{106/265} = 1.5$$

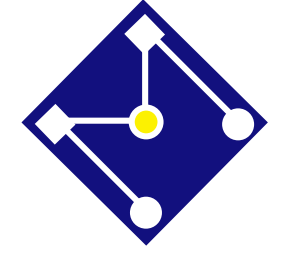
$$OR_H = \frac{121/265}{125/265} = 0.968$$

Quanto maior OR , melhor é a chance de sucesso.

Esse valor pode variar de 0 a ∞ .


$$OR = \omega_0 + \omega_1 x_1^{(i)} + \omega_2 x_2^{(i)} + \dots + \omega_n x_n^{(i)}$$





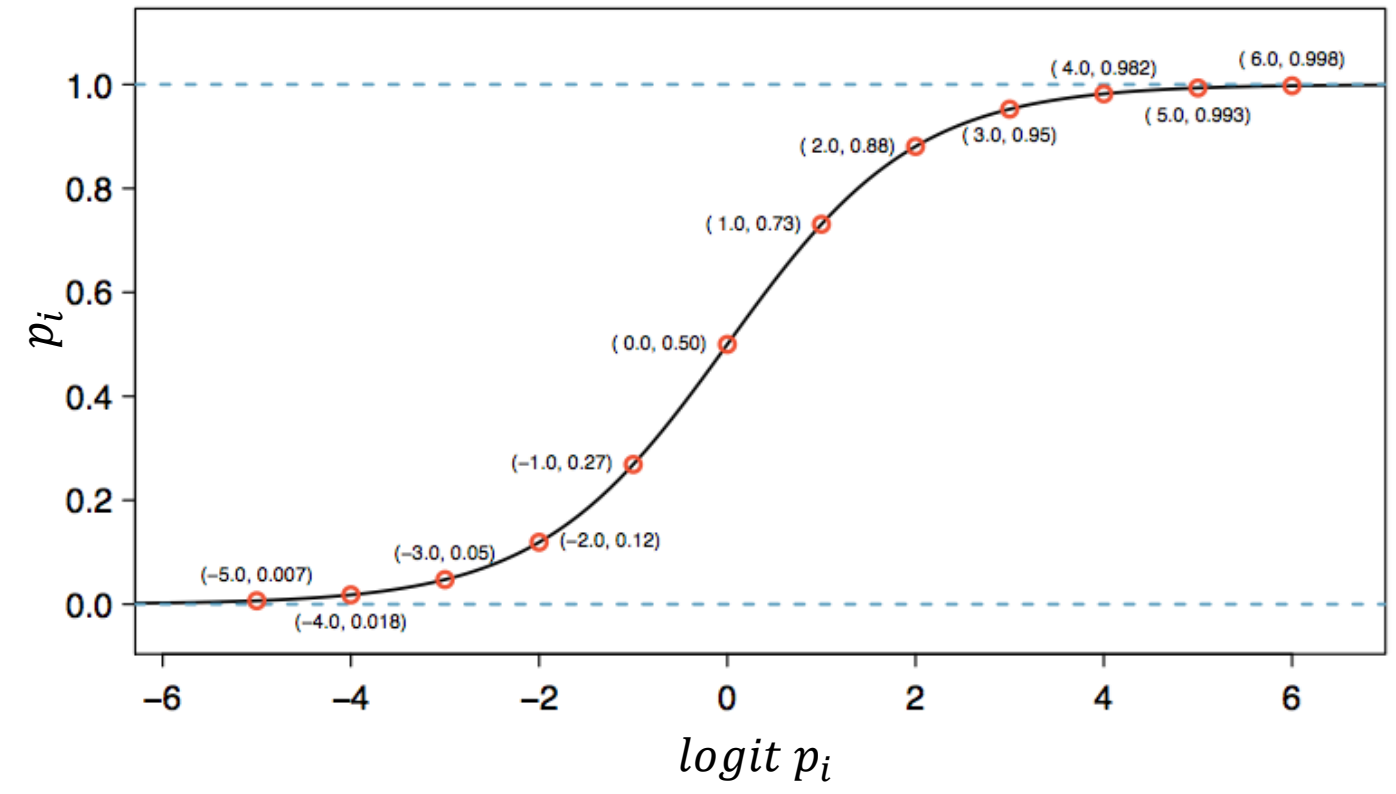
O PROBLEMA DE CLASSIFICAÇÃO

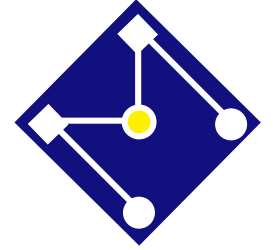
Uma transformação comum para p_i é a transformação de *logit*, que pode ser escrita como

$$\text{logit } OR_i = \ln\left(\frac{p_i}{1-p_i}\right) = \boldsymbol{\omega}^T \boldsymbol{x}$$


Resolvendo para p_i tal que,

$$p_i = g(\boldsymbol{\omega}^T \boldsymbol{x})$$





O PROBLEMA DE CLASSIFICAÇÃO

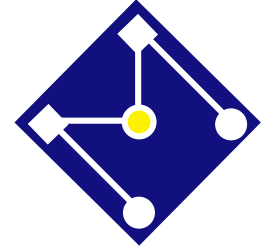
De forma completa,

$$p_i = p(y = k | \mathbf{x}^{(i)}; \boldsymbol{\omega})$$

descreve a probabilidade para do dado i pertencer à classe $k = 1, 2, \dots, K$, dado que conhecemos a entrada \mathbf{x} , parametrizada por $\boldsymbol{\omega}$.

$$p(y = k | \mathbf{x}^{(i)}; \boldsymbol{\omega}) = g(\boldsymbol{\omega}^T \mathbf{x}^{(i)})$$

Como $g(\boldsymbol{\omega}^T \mathbf{x})$ é um modelo de probabilidade, $0 \leq g(\boldsymbol{\omega}^T \mathbf{x}) \leq 1$ para qualquer \mathbf{x} .



CLASSIFICAÇÃO BINÁRIA VS MULTICLASSE

Para problemas de classificação binária ($k = 1, 2$, e $y = 0, 1$), aprenderemos um modelo $g(x)$ para o qual

$$p(y = 1|x; \omega) \text{ é modelado por } g(\omega^T x)$$

$$p(y = 0|x; \omega) \text{ é modelado por } 1 - g(\omega^T x)$$

Para problemas multiclasse ($k = 1, 2, \dots, K$), o classificador retorna uma função com valor vetorial $g(x)$, em que

$$p(y = 1|x; \omega) \text{ é modelado por } g_1(\omega^T x)$$

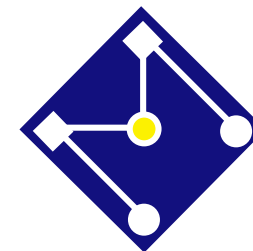
$$p(y = 2|x; \omega) \text{ é modelado por } g_2(\omega^T x)$$

$$\vdots$$

$$p(y = K|x; \omega) \text{ é modelado por } g_M(\omega^T x)$$

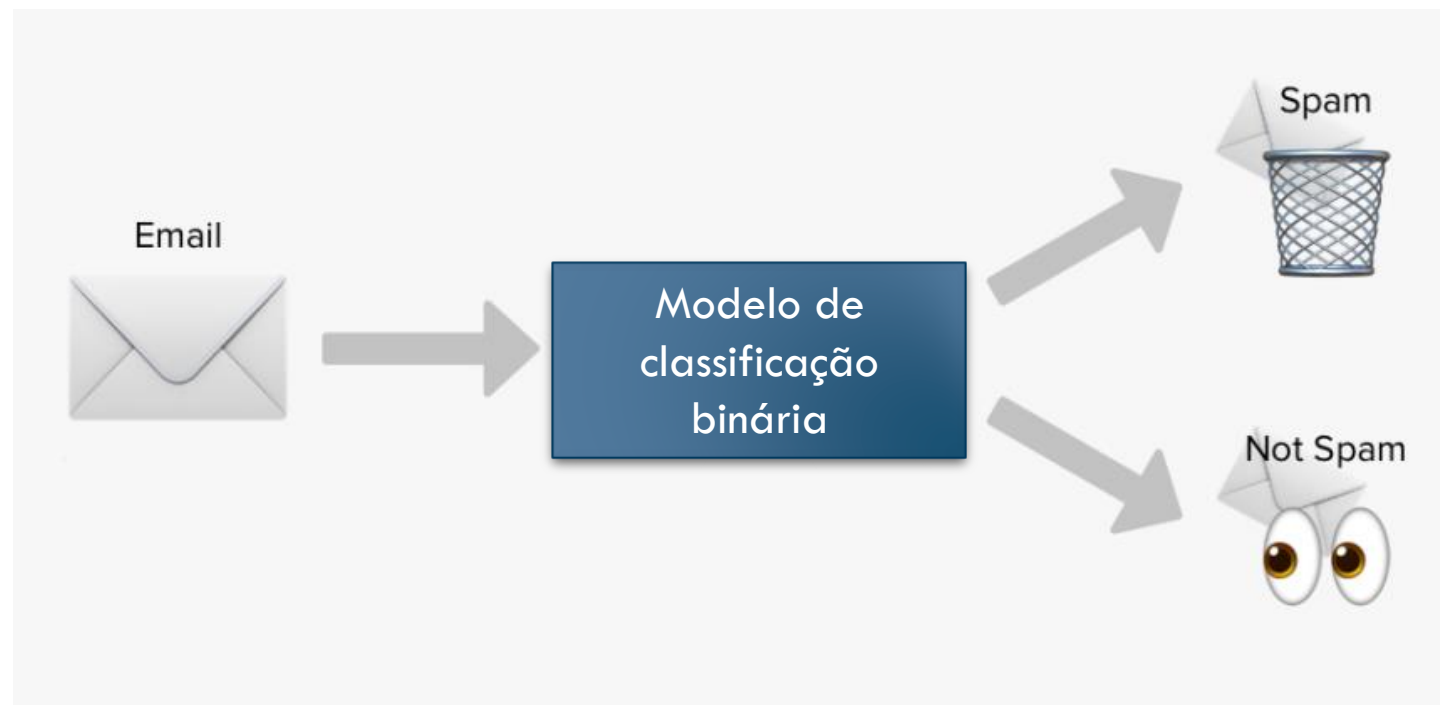
$$g(x) = \begin{bmatrix} g_1(\omega^T x) \\ g_2(\omega^T x) \\ \vdots \\ g_K(\omega^T x) \end{bmatrix}$$

cada elemento $g_k(\omega^T x)$ de $g(x)$ corresponde à probabilidade condicional da classe $p(y = k|x; \omega)$

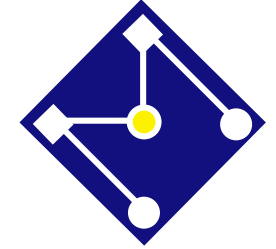


CLASSIFICAÇÃO BINÁRIA

Por enquanto, vamos nos concentrar no problema de classificação binária no qual y pode assumir apenas dois valores, 0 e 1 (a maior parte do que aprendermos também será generalizada para o caso de várias classes).



https://www.pngitem.com/middle/ibiboIT_spam-not-spam-machine-learning-hd-png-download/



RESUMO: REGRESSÃO LOGÍSTICA BINÁRIA

$$y^{(i)} = 0,1$$

$$\ln\left(\frac{p_i}{1-p_i}\right) = \omega_0 + \omega_1 x_1^{(i)} + \omega_2 x_2^{(i)} + \dots + \omega_n x_n^{(i)}$$

$$\ln\left(\frac{p_i}{1-p_i}\right) = \omega^T x$$

$$p_i = g(\omega^T x) = \frac{1}{1 + e^{-\omega^T x}}$$

Temos uma hipótese de como mc
nosso problema.

$$h(x) = g(\omega^T x) = \frac{1}{1 + e^{-\omega^T x}}$$

$$e^{\ln\left(\frac{p_i}{1-p_i}\right)} = e^{\omega^T x}$$

$$\frac{p_i}{1-p_i} = e^{\omega^T x}$$

$$p_i = (1-p_i)e^{\omega^T x}$$

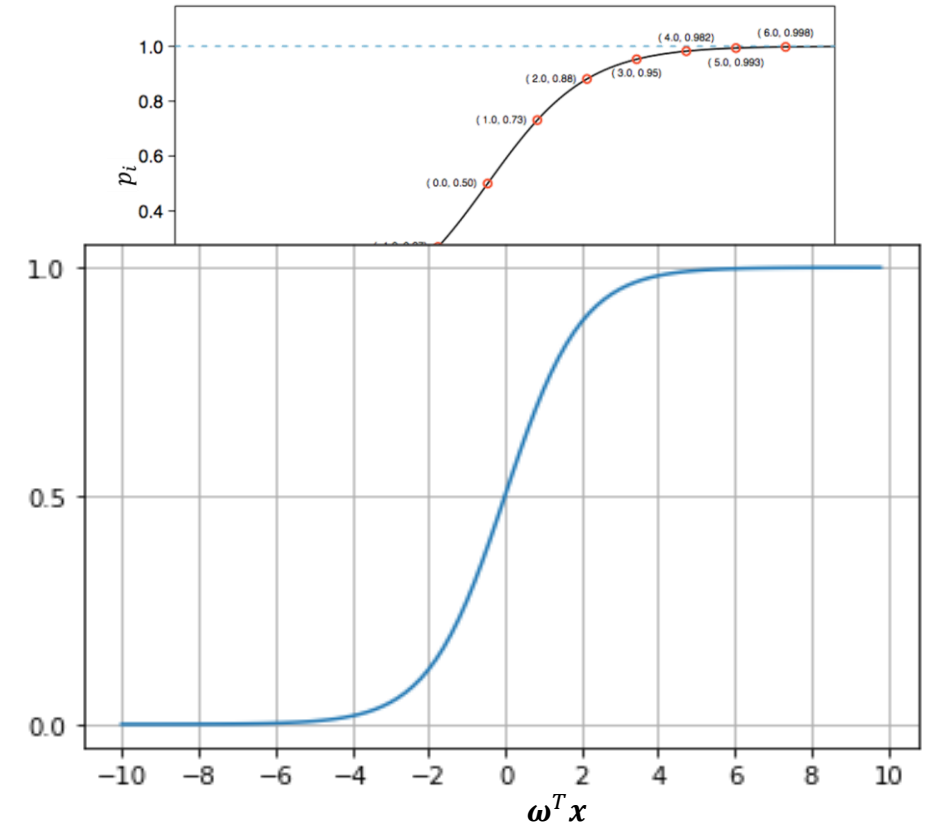
$$p_i + p_i e^{\omega^T x} = e^{\omega^T x}$$

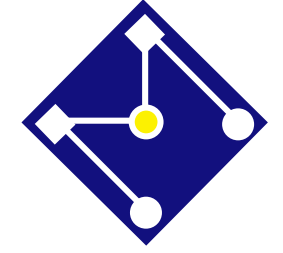
$$p_i(1 + e^{\omega^T x}) = e^{\omega^T x}$$

$$p_i = \frac{e^{\omega^T x}}{1 + e^{\omega^T x}} \frac{e^{-\omega^T x}}{e^{-\omega^T x}}$$

$$p_i = \frac{1}{1 + e^{-\omega^T x}}$$

ω ???



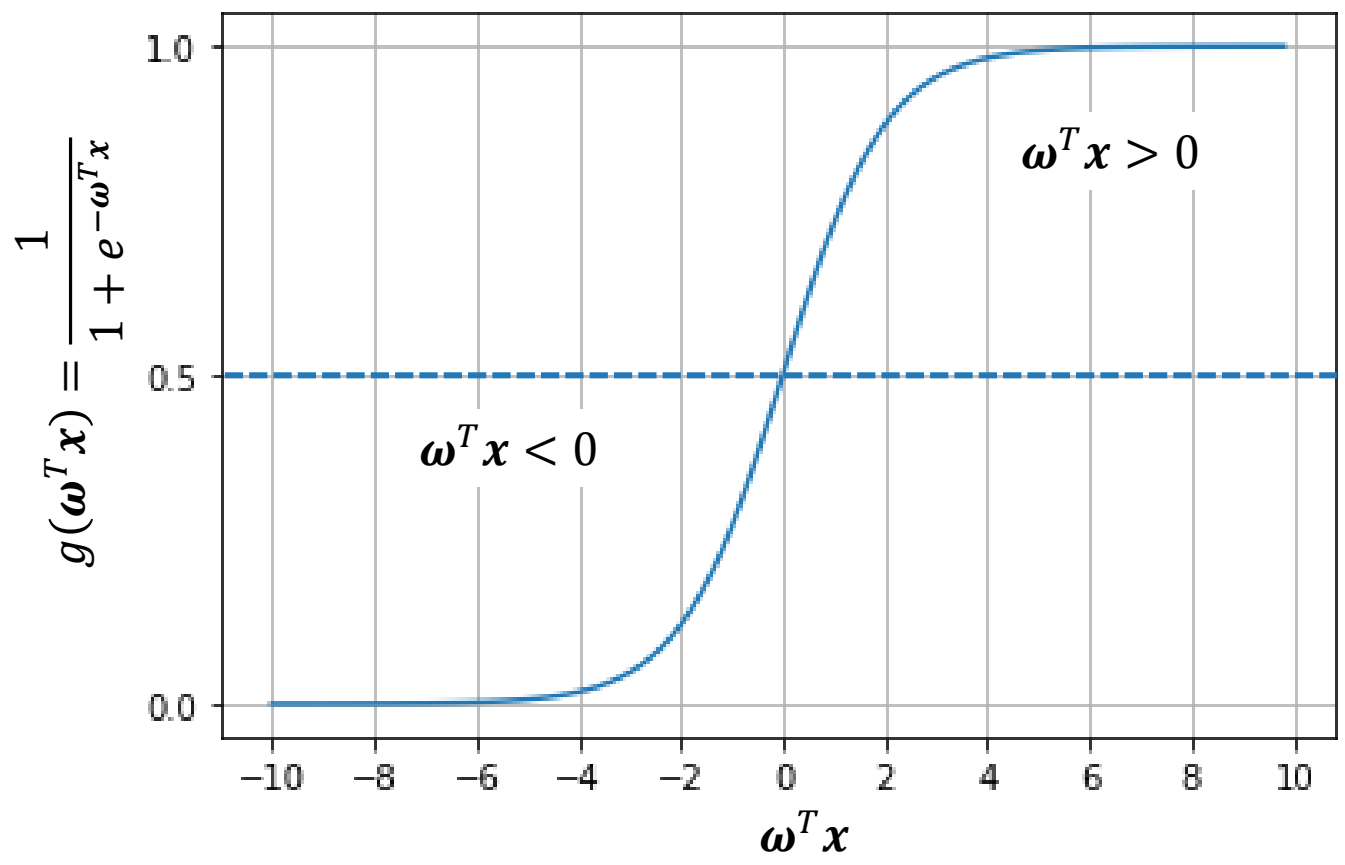


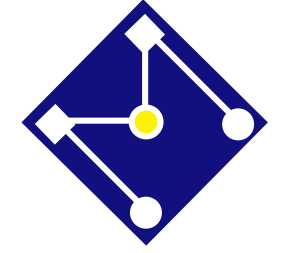
HIPÓTESE h

$$h_{\omega}(\mathbf{x}) = \frac{1}{1 + e^{-\omega^T \mathbf{x}}} = g(\omega^T \mathbf{x})$$

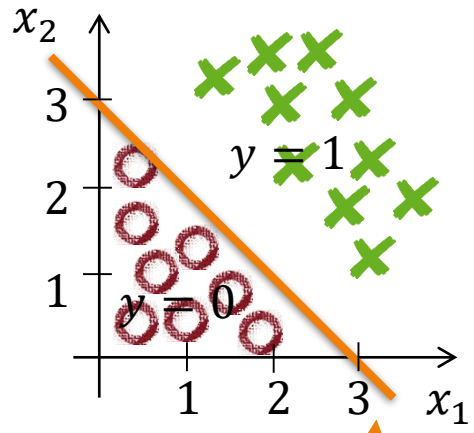
$$h_{\omega}(\mathbf{x}) = \begin{cases} > 0.5 & \omega^T \mathbf{x} > 0 \\ < 0.5 & \omega^T \mathbf{x} < 0 \end{cases}$$

Se a soma ponderada de entradas for maior que zero, a classe prevista é 1, caso contrário, será 0. Portanto, **o limite de decisão que separa as duas classes pode ser encontrado configurando a soma ponderada das entradas como $\omega^T \mathbf{x} = 0$.**





CONTORNO DE DECISÃO



Contorno de decisão
 É uma propriedade da hipótese e dos parâmetros definidos.

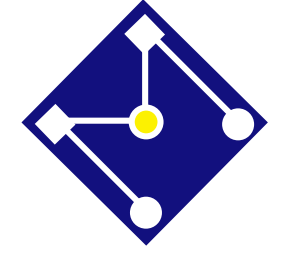
$$\boldsymbol{\omega} = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$h_{\boldsymbol{\omega}}(\mathbf{x}) = g(\boldsymbol{\omega}^T \mathbf{x}) = g(\omega_0 + \omega_1 x_1 + \omega_2 x_2)$$

Preveremos $y = 1$ se:

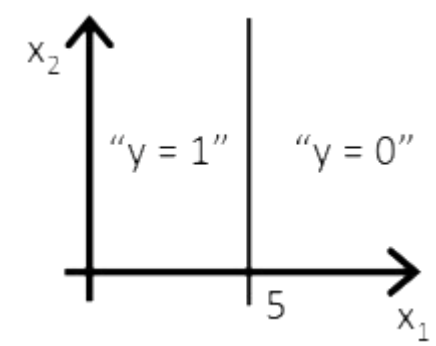
$$-3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 = 3 \quad h_{\boldsymbol{\omega}}(\mathbf{x}) = 0.5$$

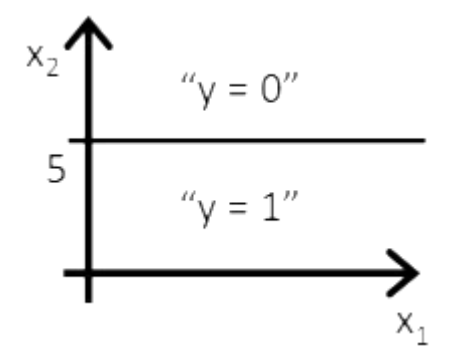


QUESTÃO PARA VOCÊ PENSAR

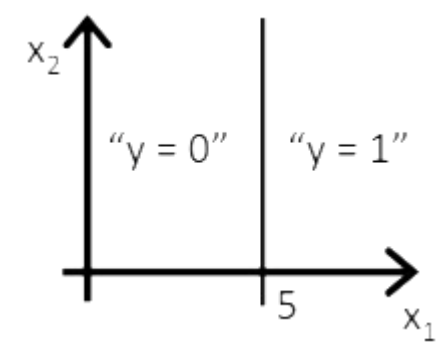
Considere a regressão logística com 2 características, $h_{\omega}(x) = g(\omega^T x) = g(\omega_0 +$



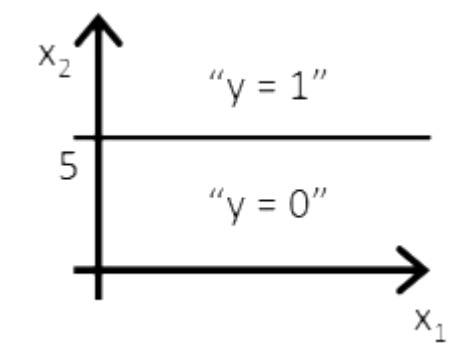
(a)



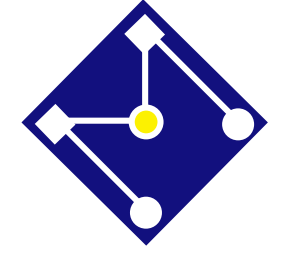
(b)



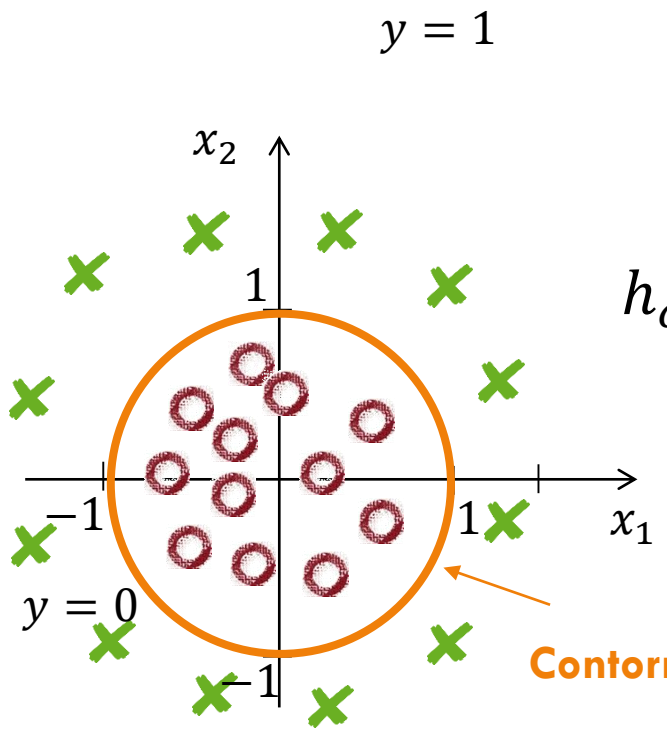
(c)



(d)



CONTORNO NÃO LINEAR



$$\omega = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \omega_3 \\ \omega_4 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$h_\omega(\mathbf{x}) = g(\omega^T \mathbf{x}) = g(\omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_1^2 + \omega_4 x_2^2)$$

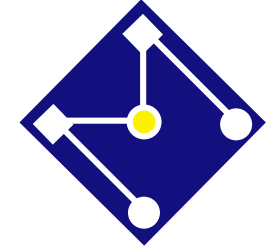
Preveremos $y = 1$ se:

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

Contorno de decisão

$$x_1^2 + x_2^2 = 1$$

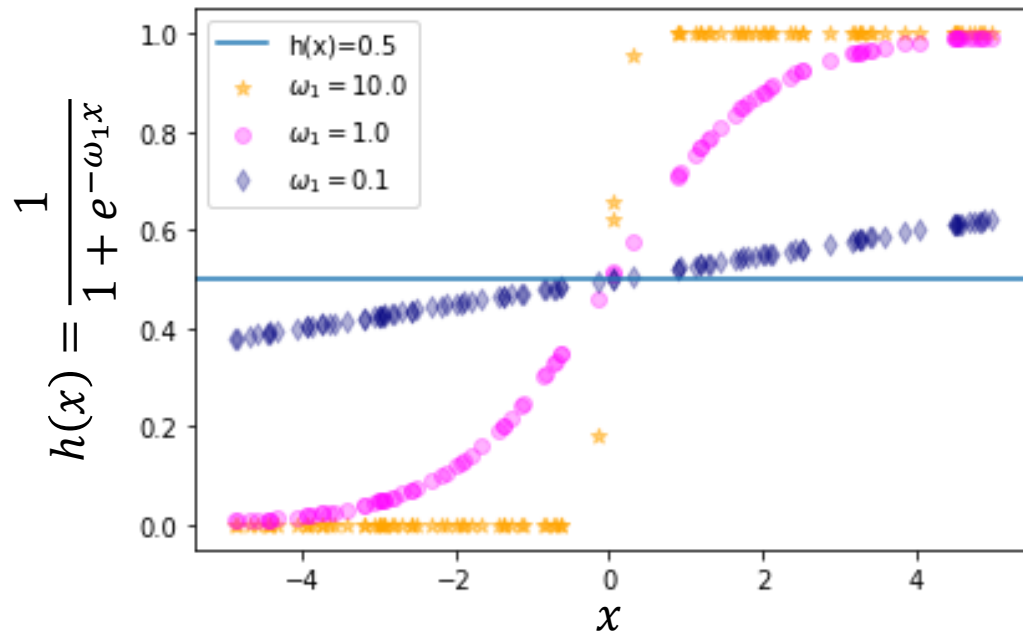
$$h_\omega(\mathbf{x}) = 0.5$$



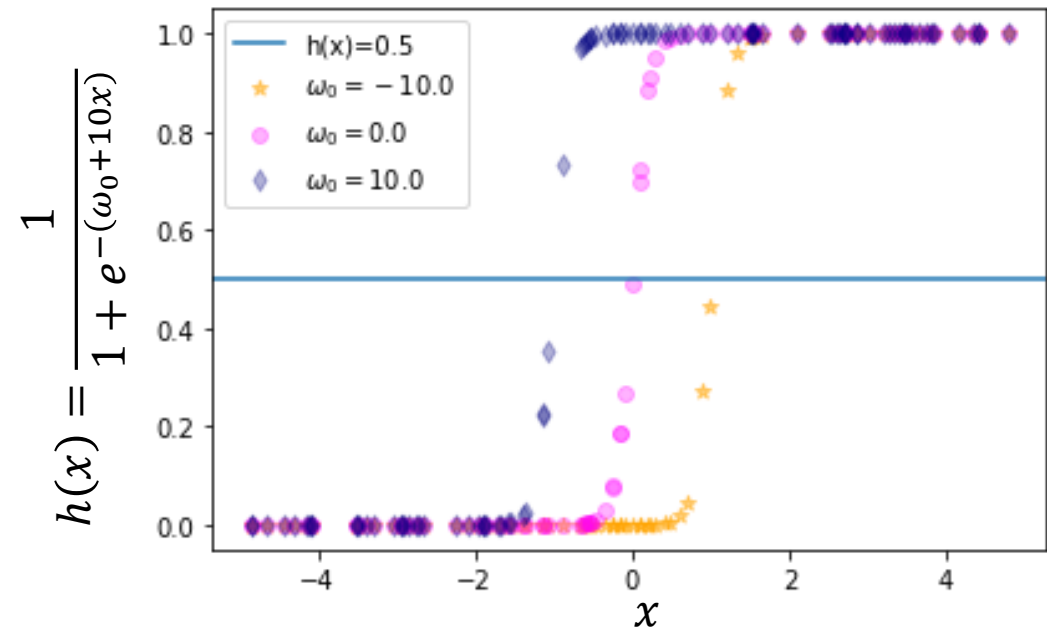
ANALISANDO O CASO LINEAR:

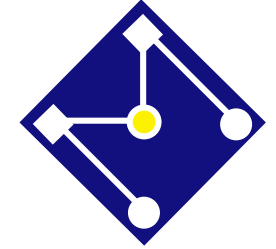
$$\omega^T \mathbf{x} = \omega_0 + \omega_1 x$$

$\omega_0 = 0$



$\omega_1 = 10$





EXEMPLO

Suponha que a probabilidade de um cliente adquirir uma assinatura de uma revista por mala direta é,

$$\text{prob}(\text{evento}) = \frac{1}{1 + e^{-(-1.143 + 0.452x_1 + 0.029x_2 - 0.242x_3)}}$$

x_1 é sexo (1 para feminino e 0 para masculino), x_2 é idade e x_3 é estado civil (1 para solteiro e 0 para casado).

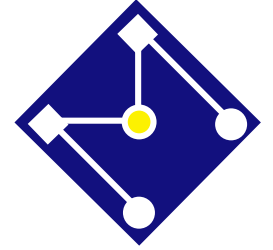
Uma pessoa do sexo feminino, com 40 anos de idade e casada, irá adquirir a assinatura da revista?

$$\text{prob}(\text{evento}) = \frac{1}{1 + e^{-(-1.143 + 0.452 \times 1 + 0.029 \times 40 - 0.242 \times 0)}} = 0.61$$

Sim, irá adquirir a revista.

Exemplo extraído de:

L. P. Favero; P. Belfiore; F. L. da Silva; B. L. Chan. *Análise de dados: modelagem multivariada para tomada de decisões*, Ed. Campus



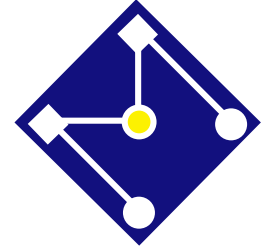
NOSSO PROBLEMA, ENTÃO É...

Conjunto de m dados treinamento: $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$
onde

$$\mathbf{x}^{(i)} \in \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}, x_0^{(i)} = 1, y \in \{0,1\}$$

Como escolho ω ????

$$h_{\omega}(\mathbf{x}) = g(\omega^T \mathbf{x}) = \frac{1}{1 + e^{-\omega^T \mathbf{x}}}$$

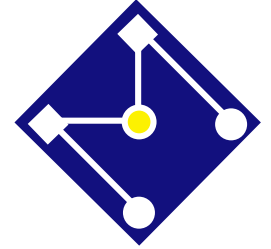


ENTENDENDO O CUSTO MATEMATICAMENTE...

Para obter a função custo, a interpretaremos estatisticamente com o método de máxima verossimilhança. Maximizar a função de verossimilhança equivale a encontrar o valor de ω que torna a observação de \mathbf{y} a **mais provável** possível,

$$\hat{\omega} = \arg \max_{\omega} p(\mathbf{y}|\mathbf{x}; \omega)$$

onde $p(\mathbf{y}|\mathbf{x}; \omega)$ é a densidade de probabilidade de todas as saídas observadas \mathbf{y} nos dados de treinamento, dadas todas as entradas \mathbf{x} e parâmetros ω . Isso determina matematicamente o que significa **mais provável**.



UM POUCO DA NOMENCLATURA A SER USADA

A probabilidade, de acordo com a definição de distribuição de Bernoulli,

$$p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\omega}) = [h_{\omega}(\mathbf{x}^{(i)})]^{y^{(i)}} [1 - h_{\omega}(\mathbf{x}^{(i)})]^{1-y^{(i)}}$$

Onde $h_{\omega}(\mathbf{x}^{(i)}) = \frac{1}{1+e^{-\boldsymbol{\omega}^T \mathbf{x}^{(i)}}}$

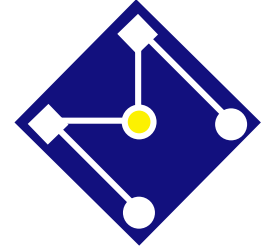
Ou seja,

se a resposta da classificação binária é $y^{(i)} = 1$

$$p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\omega}) = h_{\omega}(\mathbf{x}^{(i)})$$

se a resposta da classificação binária é $y^{(i)} = 0$

$$p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\omega}) = 1 - h_{\omega}(\mathbf{x}^{(i)})$$



CONT...

As m observações são independentes e, portanto,

$$p(\mathbf{y}|\mathbf{x}; \boldsymbol{\omega}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\omega}) = \prod_{i=1}^m h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)})^{y^{(i)}} [1 - h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)})]^{1-y^{(i)}}$$

Por razões numéricas, geralmente é melhor considerar o logaritmo de $p(\mathbf{y}|\mathbf{x}; \boldsymbol{\omega})$

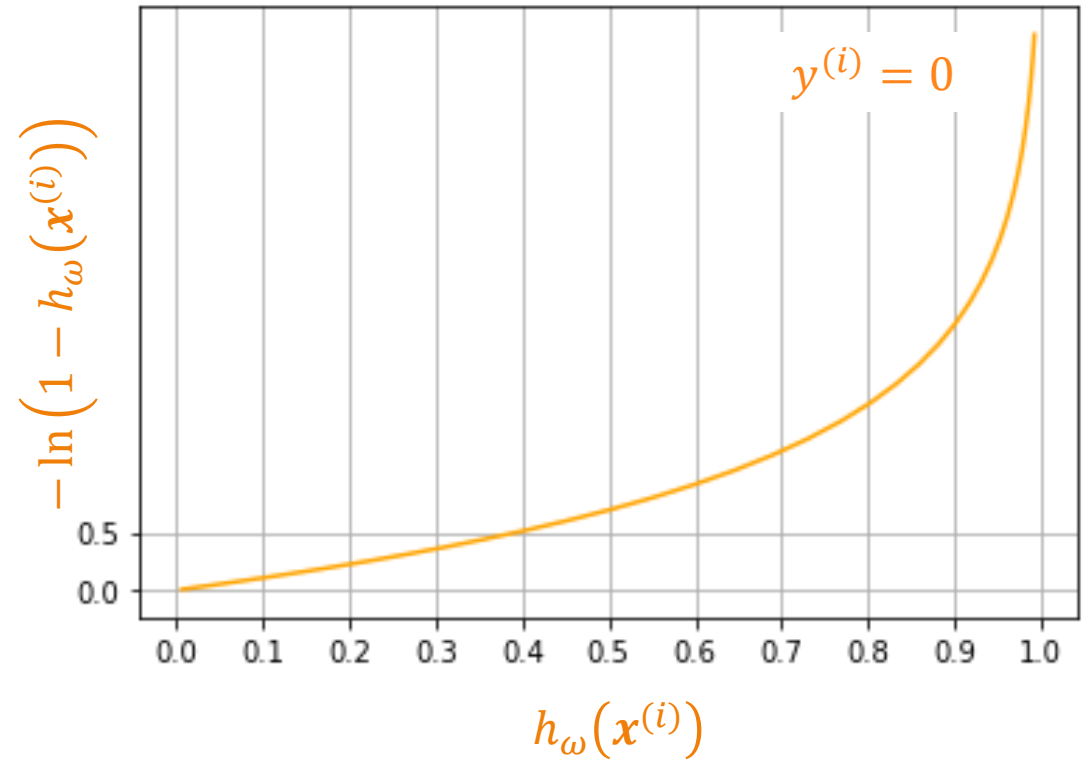
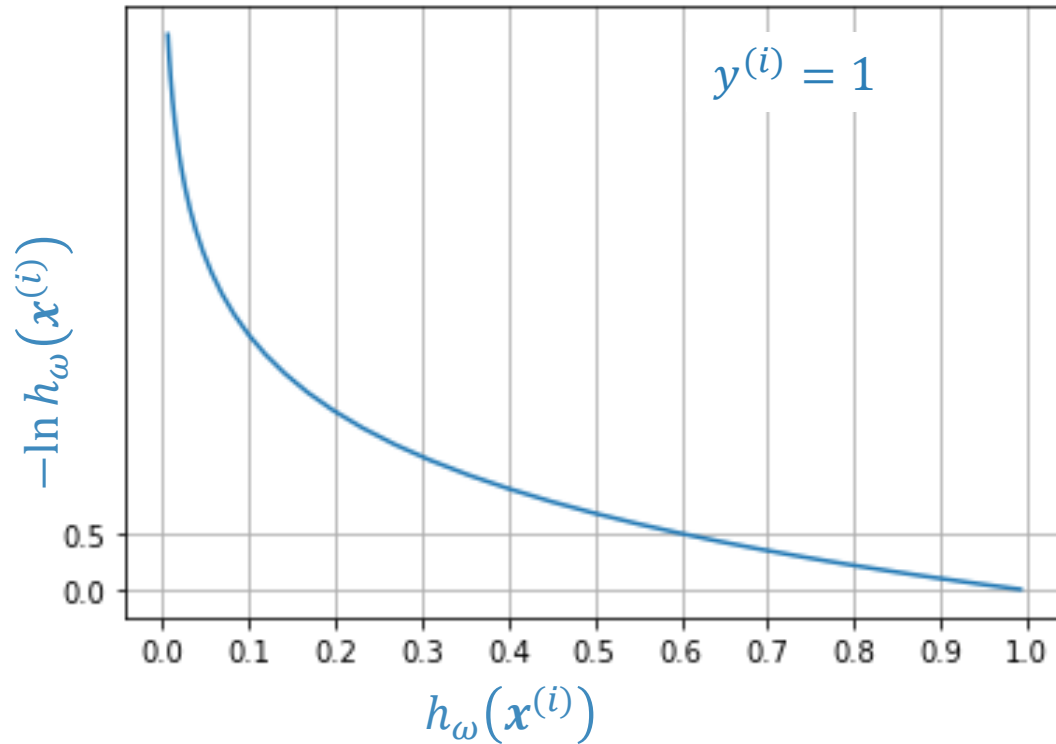
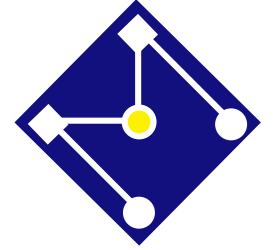
$$\ln p(\mathbf{y}|\mathbf{x}; \boldsymbol{\omega}) = \sum_{i=1}^m y^{(i)} \ln h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln[1 - h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)})]$$

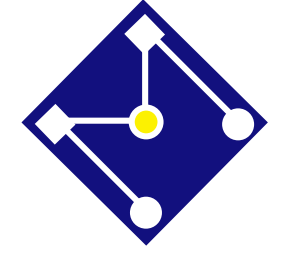
Portanto, achar o valor mais provável equivale a

$$\hat{\boldsymbol{\omega}} = \arg \max_{\boldsymbol{\omega}} \sum_{i=1}^m y^{(i)} \ln h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln[1 - h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)})]$$

$y^{(i)}$ vale 0 ou 1. Portanto,

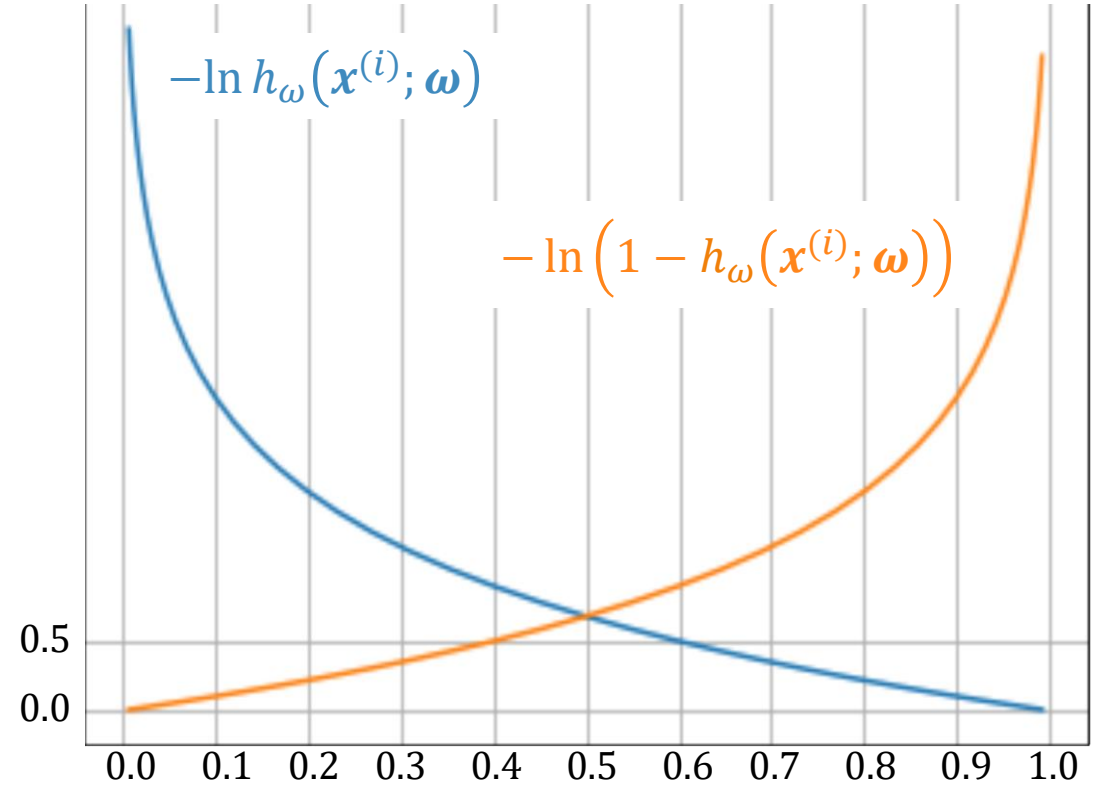
$$y^{(i)} \ln h_{\omega}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \ln[1 - h_{\omega}(\mathbf{x}^{(i)})] = \begin{cases} -\ln h_{\omega}(\mathbf{x}^{(i)}) & \text{se } y^{(i)} = 1 \\ -\ln(1 - h_{\omega}(\mathbf{x}^{(i)})) & \text{se } y^{(i)} = 0 \end{cases}$$

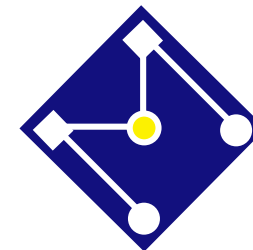




OU SEJA...

$$Custo(\omega) = \begin{cases} -\ln h_{\omega}(x^{(i)}; \omega) & \text{se } y^{(i)} = 1 \\ -\ln(1 - h_{\omega}(x^{(i)}; \omega)) & \text{se } y^{(i)} = 0 \end{cases}$$





FUNÇÃO PERDA (LOSS FUNCTION)

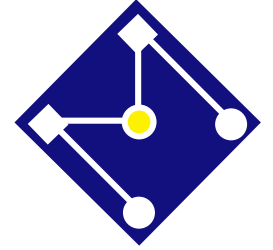
Fazendo o somatório em todo meu conjunto de dados, tenho:

$$J(\omega) = \frac{1}{m} \sum_{i=1}^m \begin{cases} -\ln[h_{\omega}(\mathbf{x}^{(i)})] & \text{se } y^{(i)} = 1 \\ -\ln[1 - h_{\omega}(\mathbf{x}^{(i)})] & \text{se } y^{(i)} = 0 \end{cases}$$

entropia cruzada

$$J(\omega) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \ln[h_{\omega}(\mathbf{x}^{(i)})] + (1 - y^{(i)}) \ln[1 - h_{\omega}(\mathbf{x}^{(i)})]$$

Entropia cruzada como função perda (Cross Entropy Loss function). Também é conhecido como **perda de log** (log loss).



GRADIENTE DESCENDENTE

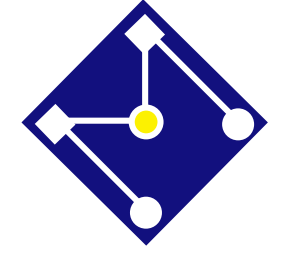
$$J(\boldsymbol{\omega}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \ln[h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)})] + (1 - y^{(i)}) \ln[1 - h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)})]$$

Queremos $\min_{\boldsymbol{\omega}} J(\boldsymbol{\omega})$

Repetir até convergência {

$$\omega_{j+1} := \omega_j - \alpha \frac{\partial}{\partial \omega_j} J(\boldsymbol{\omega}) \quad \frac{\partial}{\partial \omega_j} J(\boldsymbol{\omega}) = \frac{1}{m} \sum_{i=1}^m [h_{\boldsymbol{\omega}}(\mathbf{x}^{(i)}) - y^{(i)}] x_j^{(i)}$$

}



Controle
convergência
plotando gráfico
de J pelo número
de iterações

Algoritmo

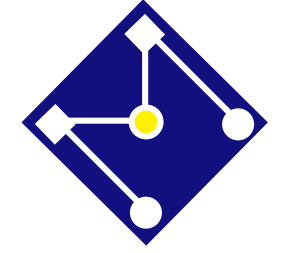
$$J(\omega) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \ln[h_{\omega}(\mathbf{x}^{(i)})] + (1 - y^{(i)}) \ln[1 - h_{\omega}(\mathbf{x}^{(i)})]$$

Queremos $\min_{\omega} J(\omega)$

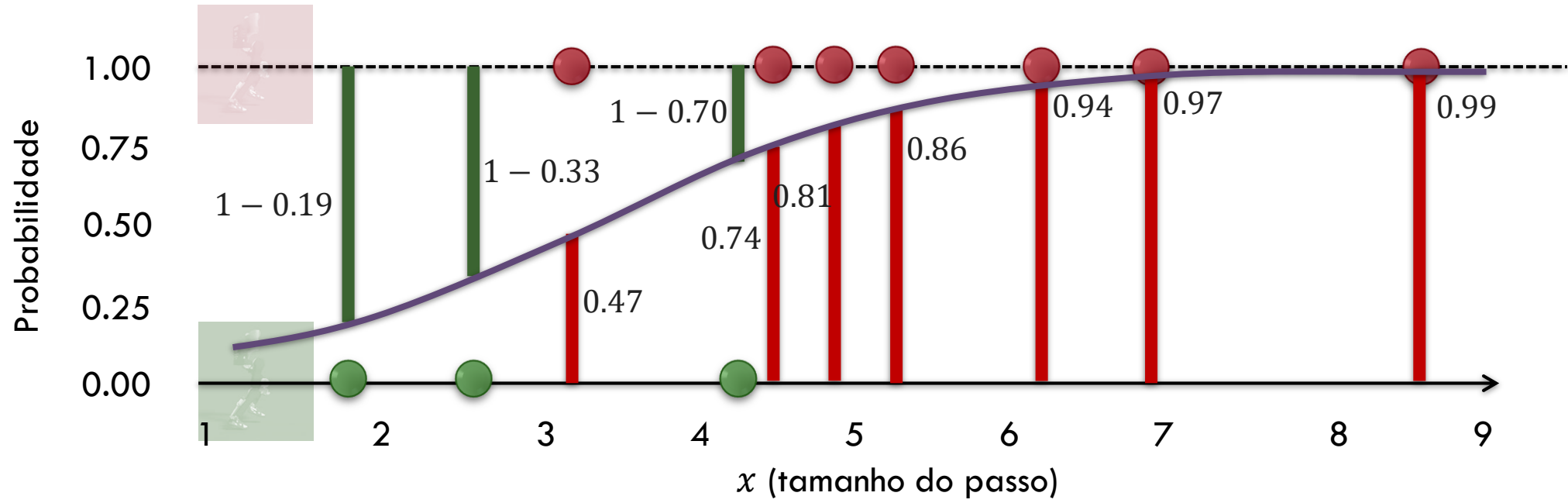
Repetir até convergência {

$$\omega_{j+1} := \omega_j - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\omega}(\mathbf{x}^{(i)}) - y^{(i)}] \mathbf{x}_j^{(i)}$$

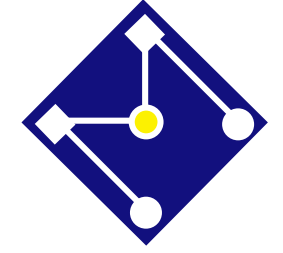
}



EM NOSSO EXEMPLO



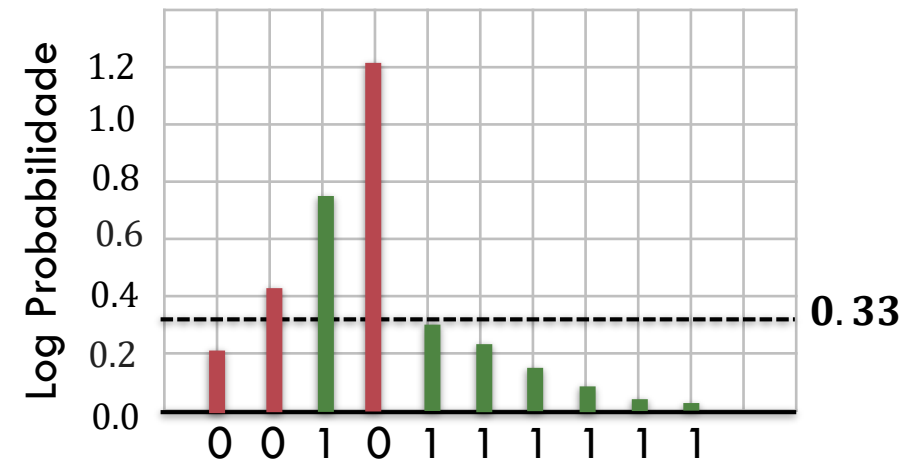
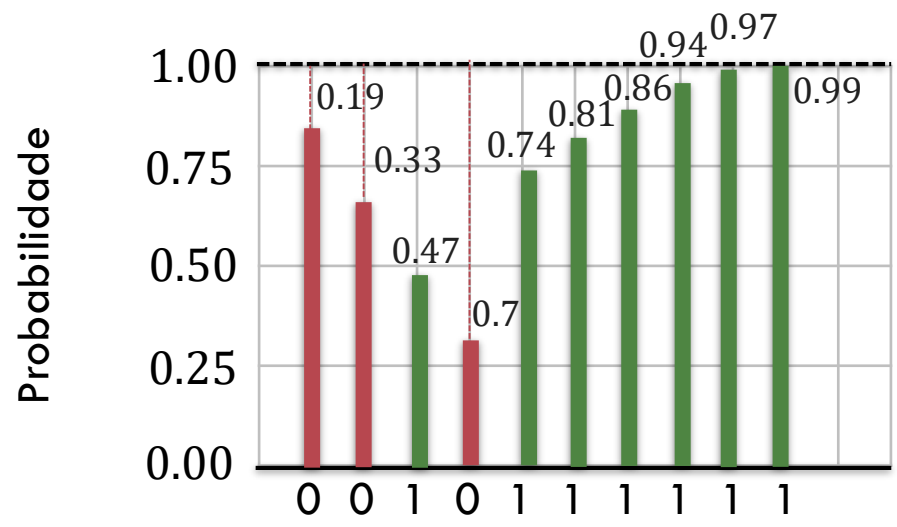
As barras representam as **probabilidades previstas** correspondentes à **verdadeira classe** de cada ponto!

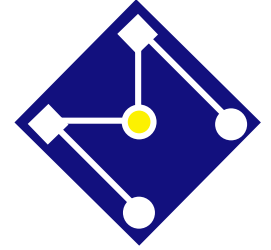


$$J(\omega) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \ln[h_{\omega}(x^{(i)})] + (1 - y^{(i)}) \ln[1 - h_{\omega}(x^{(i)})]$$

$$J(\omega) = -\frac{1}{10} \left[\begin{array}{l} \ln(1 - 0.19) + \ln(1 - 0.33) + \ln(0.47) + \ln(1 - 0.7) + \ln(0.74) + \\ \ln(0.81) + \ln(0.86) + \ln(0.94) + \ln(0.97) + \ln(0.99) \end{array} \right] = 0.3329$$

Quão próxima é a distribuição prevista da distribuição verdadeira?
É isso que determina o erro de entropia cruzada.

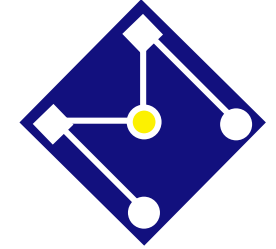




MÉTRICAS

Matriz de confusão é uma medida de desempenho para o problema de classificação de aprendizado de máquina em que a saída pode ser duas ou mais classes.

↓Prediction/Target →	Positivo (1)	Negativo (0)
Positivo (1)	VP	FP
Negativo (0)	FN	VN



Acurácia

$$A = \frac{VP + VN}{VP + VN + FP + FN}$$

performance geral do modelo

Precisão

$$P = \frac{VP}{VP + FP}$$

dentre as classificações positivas que o modelo fez, quantas estão corretas

Revocação (Recall)

$$R = \frac{VP}{VP + FN}$$

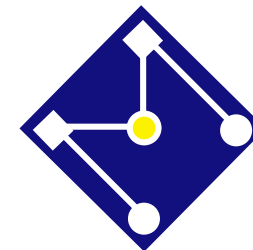
Ou sensibilidade, dentre todas as situações de classe positiva como valor esperado, quantas estão corretas

F1

$$F1 = \frac{2PR}{P + R}$$

média harmônica entre **precisão e revocação**.

↓Prediction/Target →	Positivo (1)	Negativo (0)
Positivo (1)	VP	FP
Negativo (0)	FN	VN



TRADING OFF ENTRE **PRECISÃO** E **RECALL**

$$P = \frac{\text{Verdadeiros Positivos}}{\text{N. de positivos previstos}} = \frac{VP}{VP + FN}$$

$$R = \frac{\text{Verdadeiros Positivos}}{\text{N. de positivos reais}} = \frac{VP}{VP + FP}$$

$$h_{\omega}(\mathbf{x}) = \frac{1}{1 + e^{-\omega^T \mathbf{x}}} = g(\omega^T \mathbf{x})$$

Prevemos 1 se $h_{\omega}(\mathbf{x}) \geq 0.3$

Prevemos 0 se $h_{\omega}(\mathbf{x}) < 0.3$

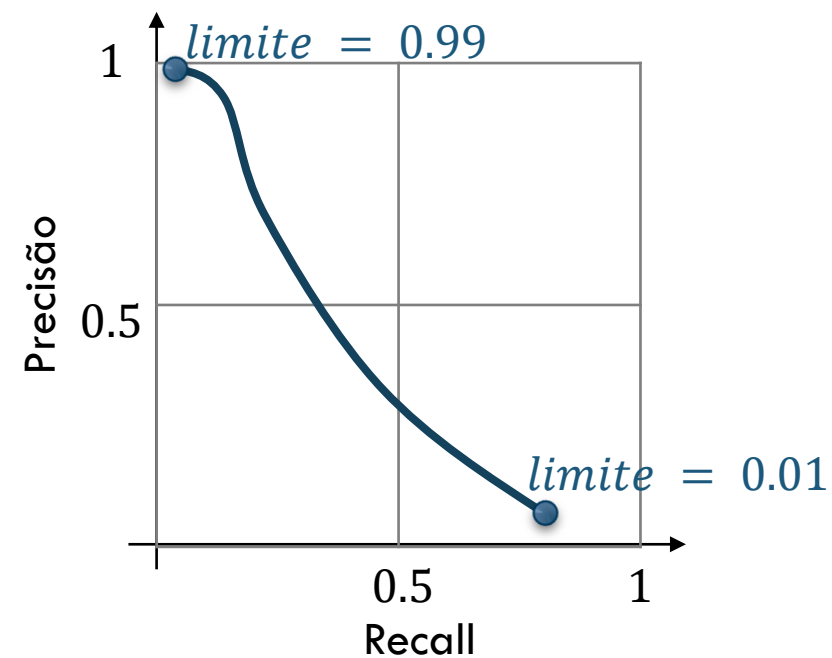
Supondo que queremos prever $y = 1$ somente se tivermos bastante certeza:

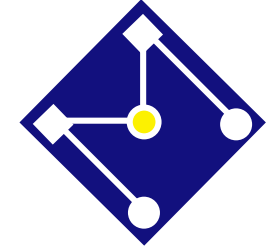
Alta precisão, baixo recall

Supondo que queremos prever $y = 1$ somente se tivermos bastante certeza

Alto recall, baixa precisão

De forma general: $h_{\omega}(\mathbf{x}) \geq \text{limite}$





F1

	Precisão (P)	Recall (R)
Modelo 01	0.5	0.4
Modelo 02	0.7	0.1
Modelo 03	0.02	1.0

$$M = \frac{P + R}{2}$$

$$F1 = \frac{2PR}{P + R}$$

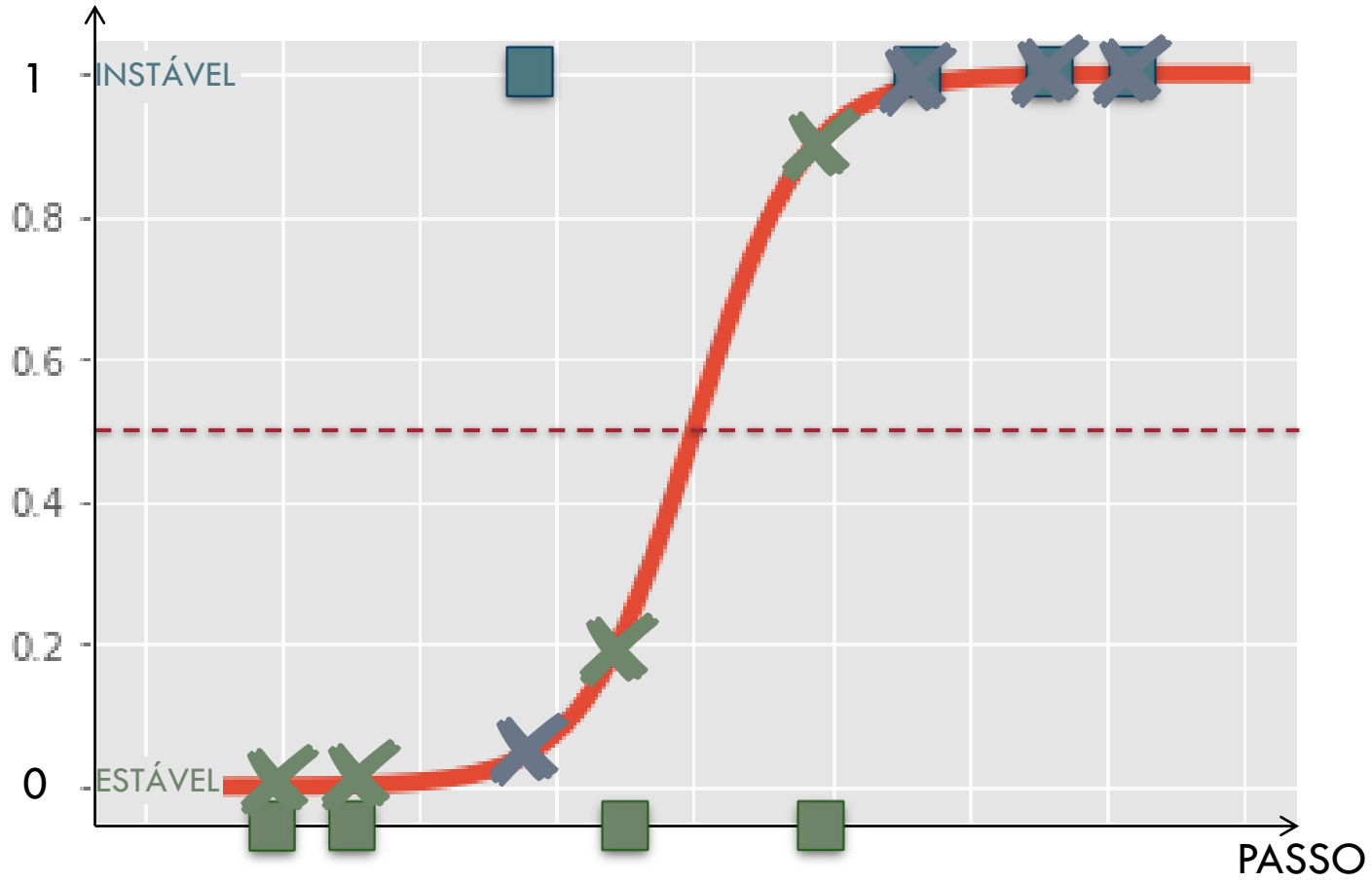
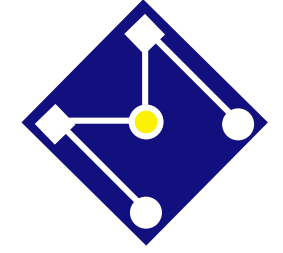
$$P = 0, R = 0 \rightarrow F1 = 0$$



$$P = 1, R = 1 \rightarrow F1 = 1$$

Exemplo extraído de:

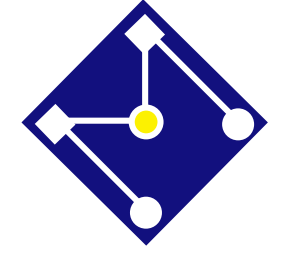
<https://www.coursera.org/learn/machine-learning/lecture/CuONQ/trading-off-precision-and-recall>



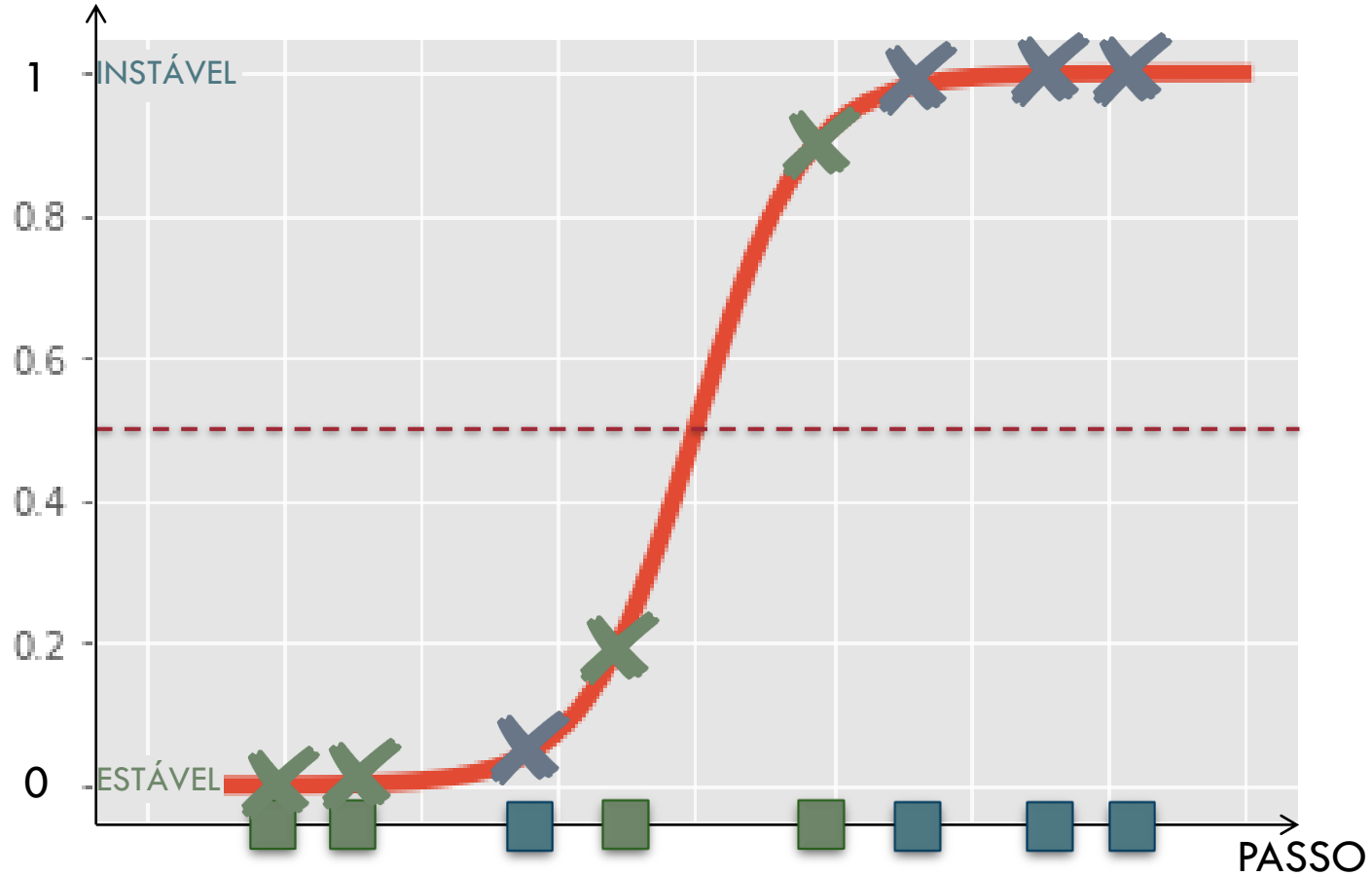
		Target		
		VP	0	
Predicted	1	3	1	FP
	0	1	3	
		FN	VN	

$$R = \frac{VP}{VP + FN} = \frac{3}{4} = 0.75$$

$$FPR = \frac{VP}{VP + VN} = \frac{3}{4}$$



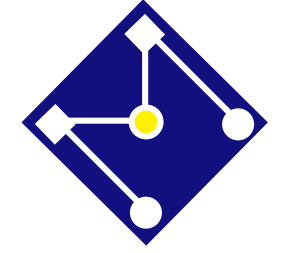
E se colocarmos o limite de modo a classificarmos SEMPRE como instável?



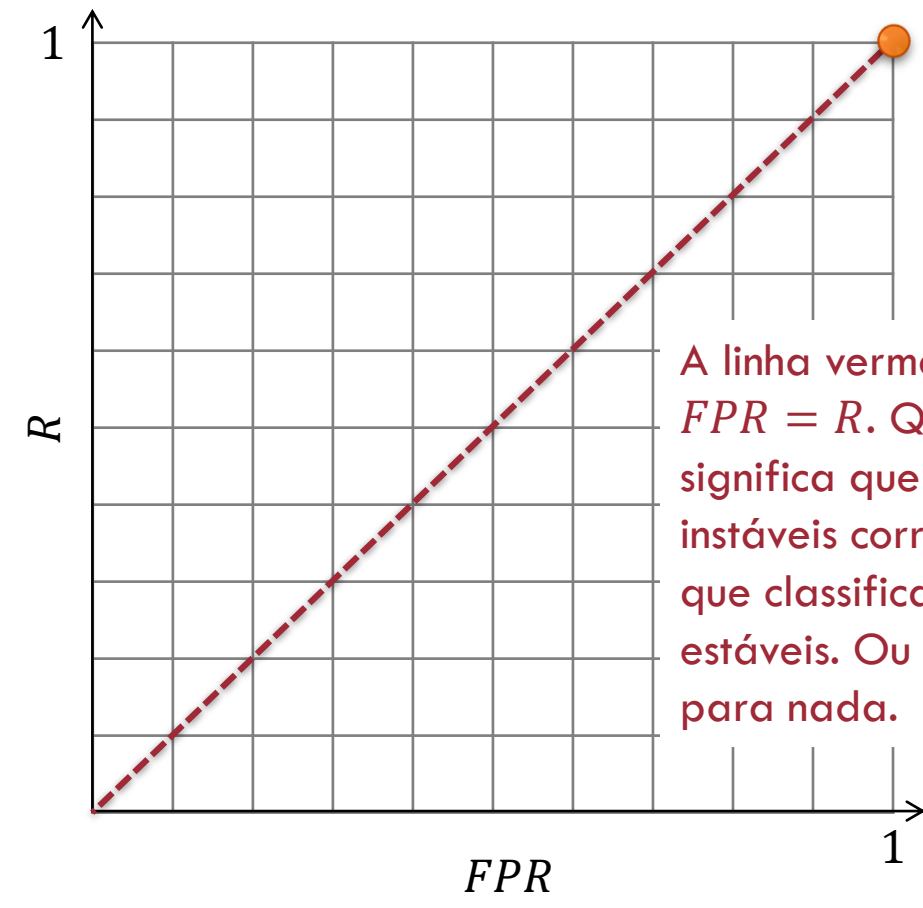
		Target		
		VP	0	
Predicted	1	4	4	FP
	0	0	0	
		FN	VN	

$$R = \frac{VP}{VP + FN} = \frac{4}{4 + 0} = 1$$

$$FPR = \frac{FP}{FP + VN} = \frac{4}{4 + 0} = 1$$

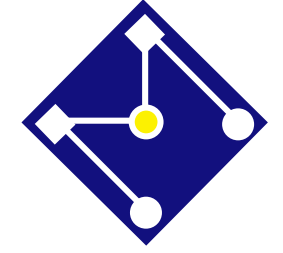


ROC: R vs FPR

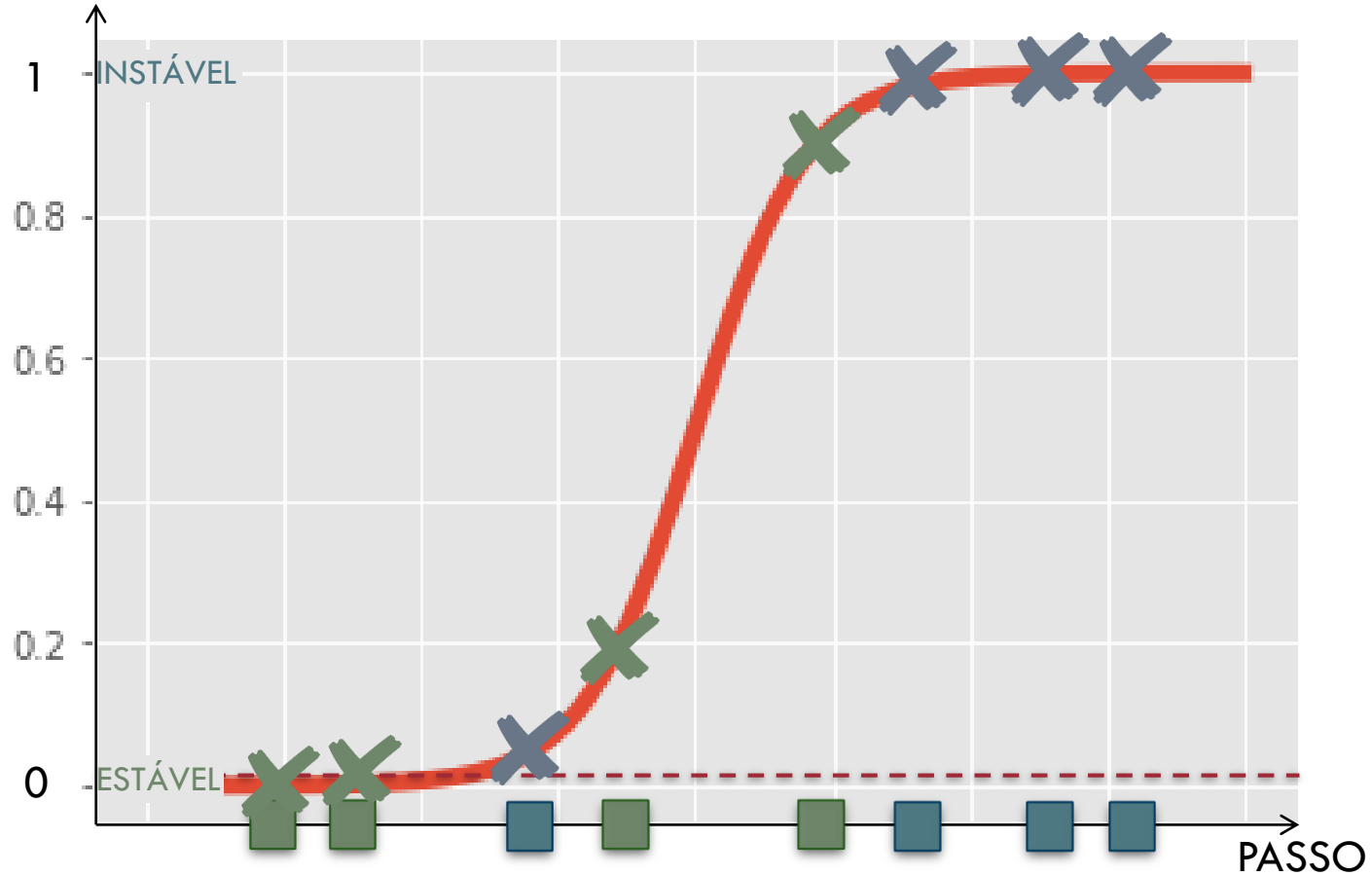


No ponto (1,1) significa que, apesar de **classificar corretamente todos os passos instáveis**, classificou **incorretamente todos os passos estáveis**.

A linha vermelha diagonal significa que $FPR = R$. Qualquer ponto nessa linha significa que a proporção de classificados instáveis corretamente é a mesma proporção que classifica incorretamente os passos estáveis. Ou seja, o classificador não serve para nada.



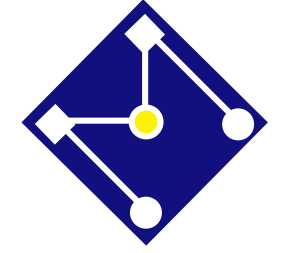
E se colocarmos o limite de modo que apenas o menor de todos os passos é classificado como estável?



		Target		
		VP	0	
Predicted	1	4	3	FP
	0	0	1	

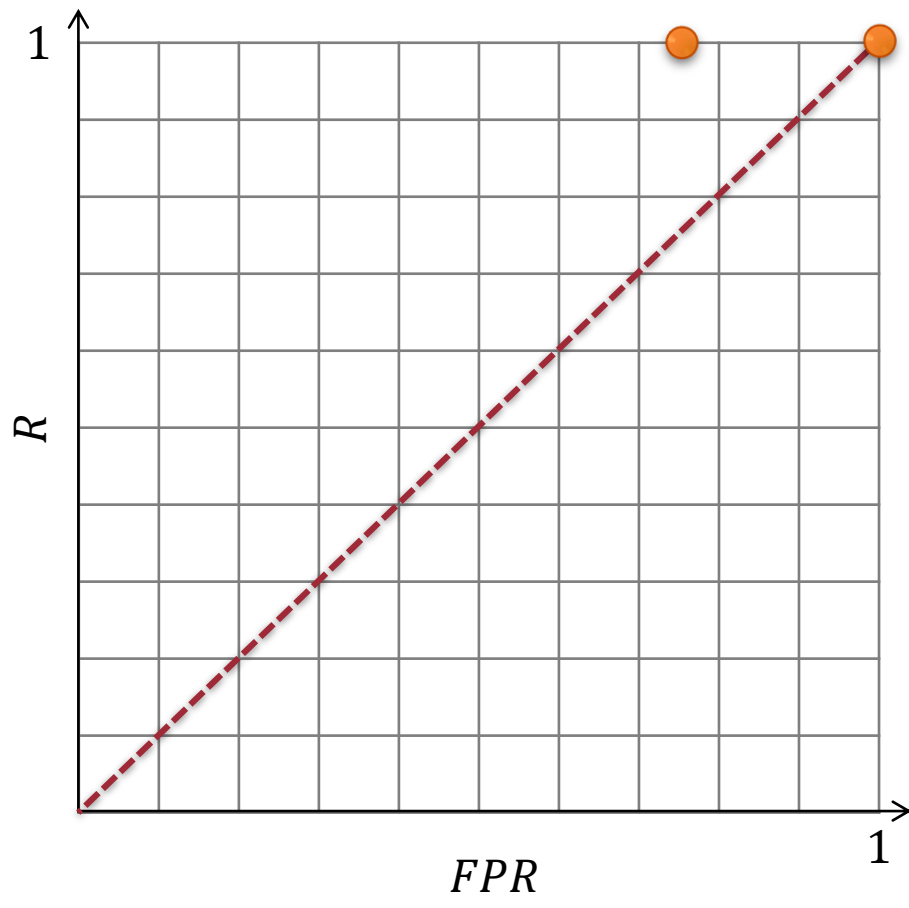
$$R = \frac{VP}{VP + FN} = \frac{4}{4 + 0} = 1$$

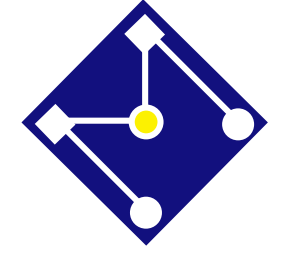
$$FPR = \frac{FP}{FP + VN} = \frac{3}{3 + 1} = 0.75$$



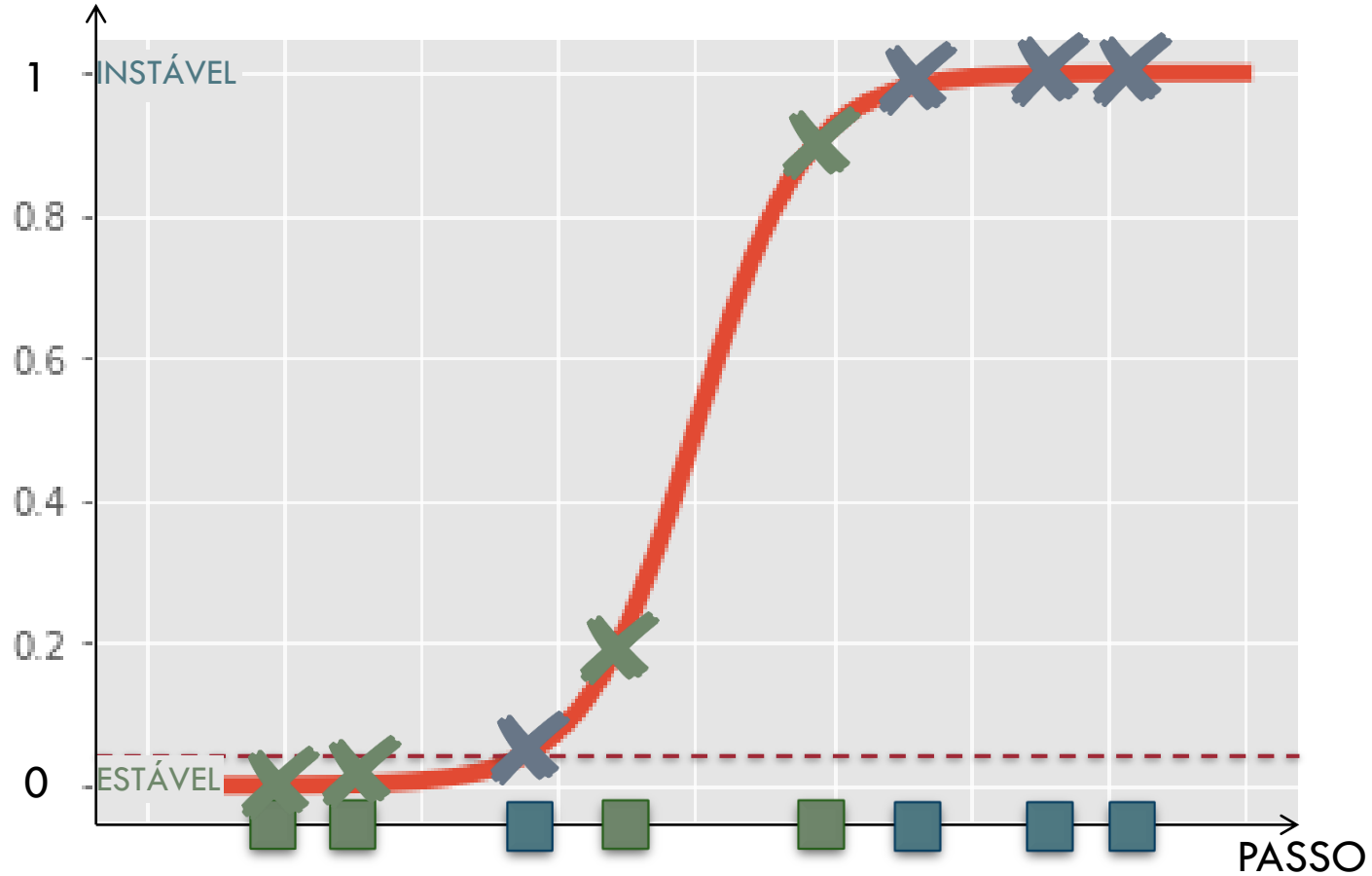
Ou seja,

ROC: R vs FPR O segundo ponto é melhor que o primeiro...





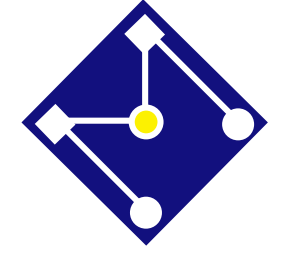
E se colocarmos o limite de modo que os dois menores passos são classificados como estáveis?



		Target		
		VP		
Predicted	1	4	2	FP
	0	0	2	
		FN	VN	

$$R = \frac{VP}{VP + FN} = \frac{4}{4 + 0} = 1$$

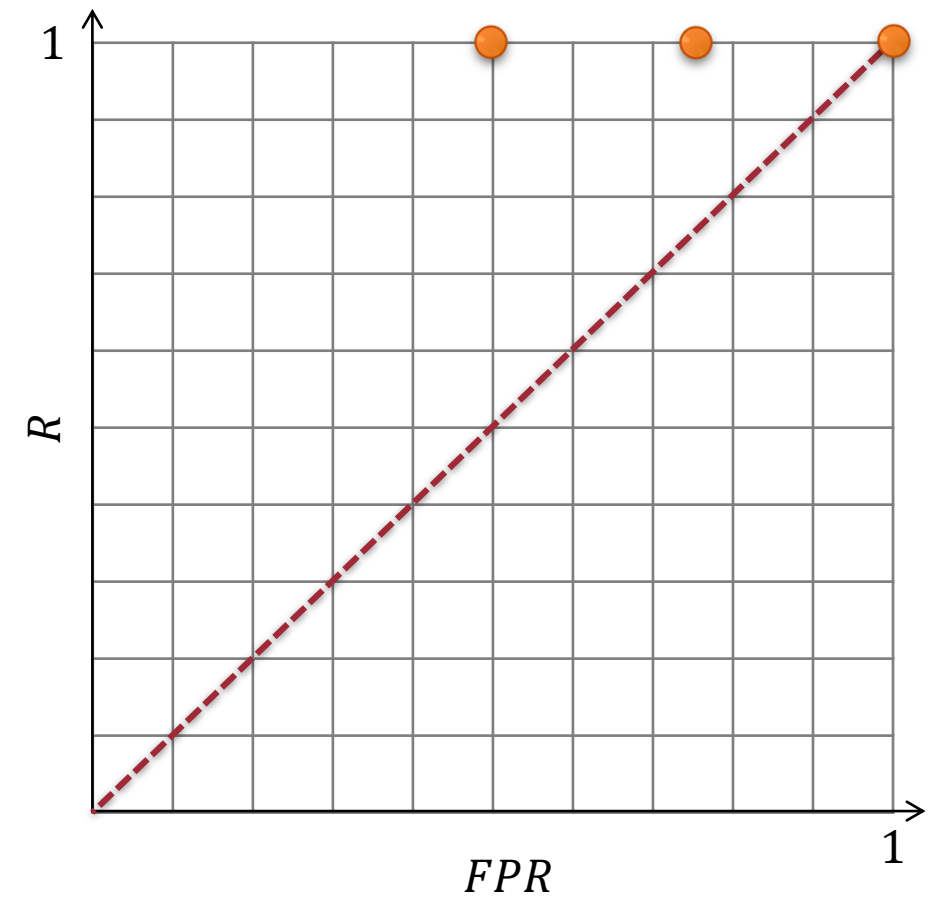
$$FPR = \frac{FP}{FP + VN} = \frac{2}{2 + 2} = 0.5$$

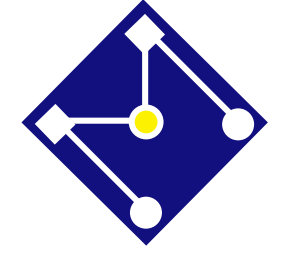


Ou seja,

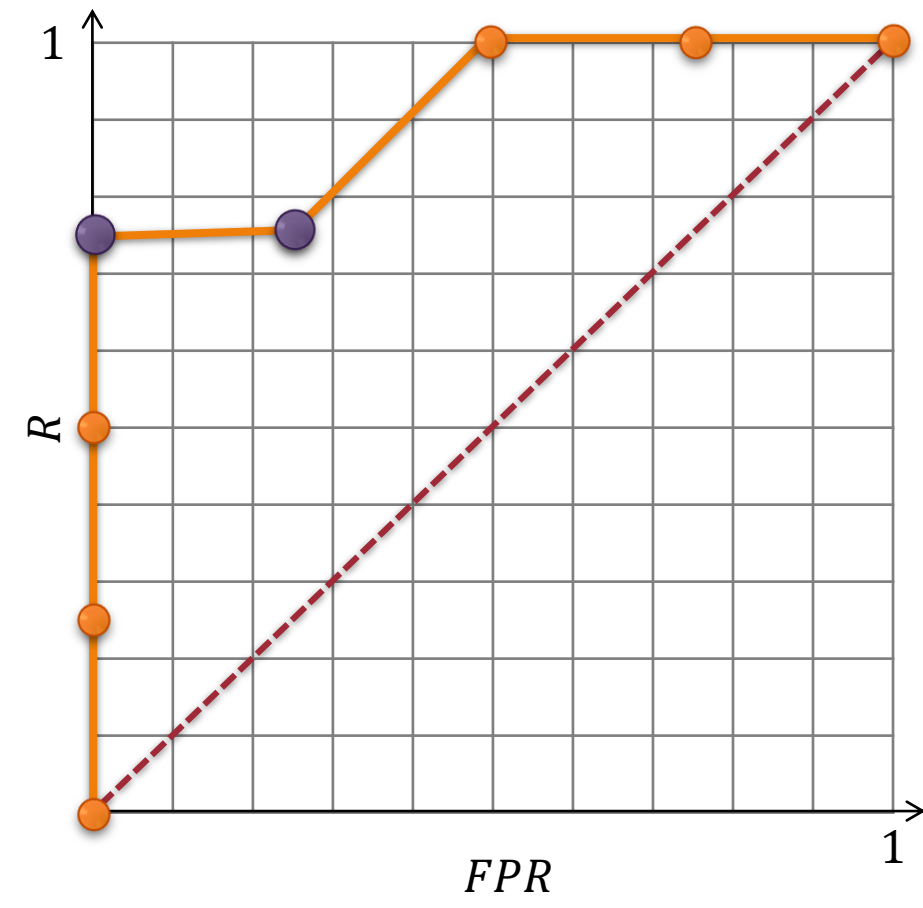
ROC: R vs FPR

O terceiro ponto é melhor que o segundo...

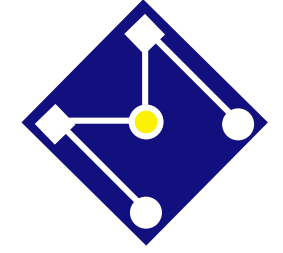




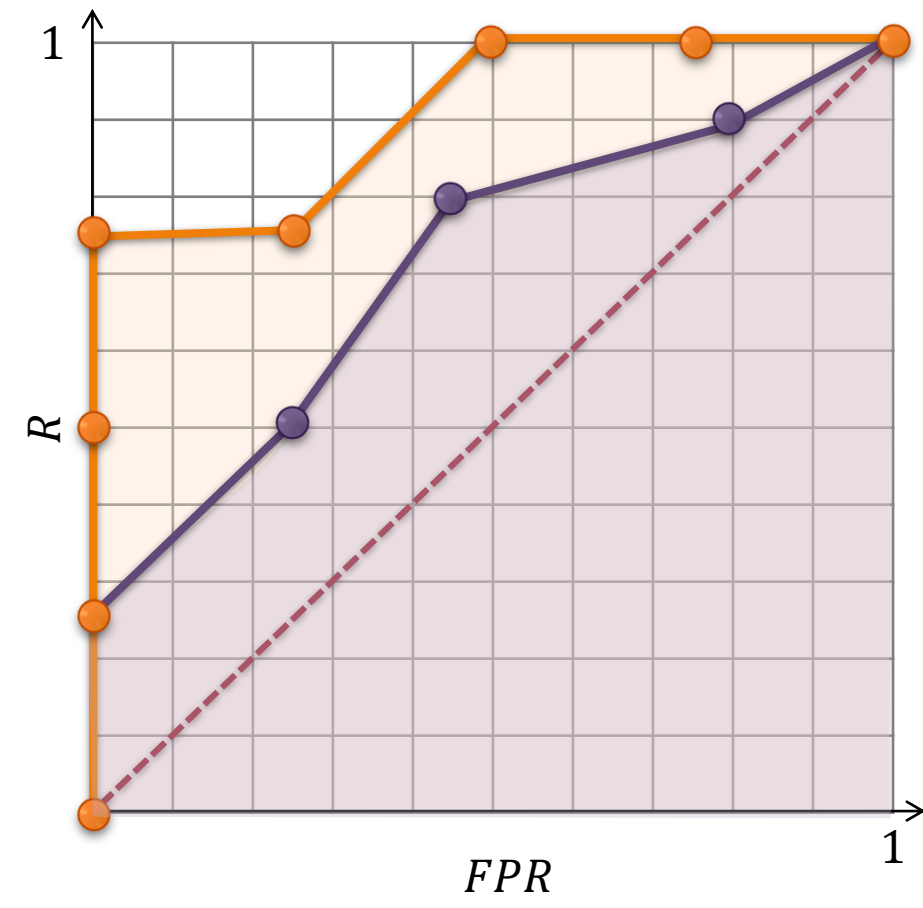
ROC: R vs FPR



Curva ROC resume a matriz de confusão de cada limite escolhido.

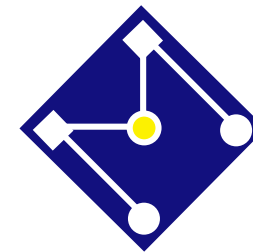


AUC



A AUC torna fácil a comparação entre diferentes ROC.

A curva laranja tem AUC maior que a curva roxa, sugerindo que é melhor.



OBSERVAÇÃO

Se as amostras são desbalanceadas (por exemplo, o número de passos estáveis é muito maior que não estáveis) então PRECISÃO pode ser mais útil que FPR. • • •

Isso ocorre porque a precisão, definida como,

$$P = \frac{VP}{VP + FP}$$

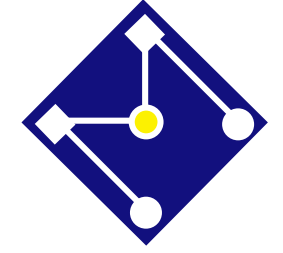
não inclui o número de VN em seu cálculo e, portanto, não é afetada pelo desbalanceamento.

$$FPR = \frac{FP}{FP + VN}$$



REGRESSÃO LOGÍSTICA PARA MAIS DE DUAS CLASSES





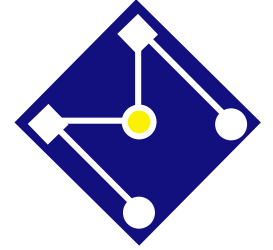
REGRESSÃO SOFTMAX

$$h_{\omega^{(k)}}(\mathbf{x}) = \frac{e^{\omega^{(k)} \cdot \mathbf{x}}}{\sum_{i=1}^K e^{\omega^{(i)} \cdot \mathbf{x}}}$$

$$h_{\omega}(\mathbf{x}) = \begin{bmatrix} P(y = 1 | \mathbf{x}; \omega) \\ P(y = 2 | \mathbf{x}; \omega) \\ \vdots \\ P(y = K | \mathbf{x}; \omega) \end{bmatrix} = \frac{1}{\sum_{j=1}^K \exp(\omega^{(j)\top} \mathbf{x})} \begin{bmatrix} \exp(\omega^{(1)\top} \mathbf{x}) \\ \exp(\omega^{(2)\top} \mathbf{x}) \\ \vdots \\ \exp(\omega^{(K)\top} \mathbf{x}) \end{bmatrix}$$

$$J(\omega) = - \left[\sum_{i=1}^m \sum_{k=1}^K 1 \{ y^{(i)} = k \} \log \frac{\exp(\omega^{(k)\top} \mathbf{x}^{(i)})}{\sum_{j=1}^K \exp(\omega^{(j)\top} \mathbf{x}^{(i)})} \right]$$

$1\{\cdot\}$ → Função indicativa



ONE-VS-ALL (ONE-VS-REST)

$$y \in \{0, 1, 2, \dots, K\}$$

Treinamos K classificadores binário separado para cada classe e executamos todos esses classificadores. Para qualquer novo exemplo \mathbf{x} que desejamos prever escolhemos a classe com a pontuação máxima:

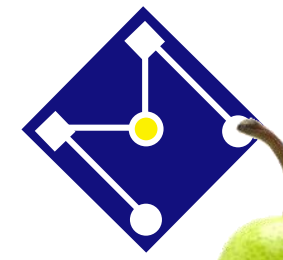
$$\hat{y} = \arg \max_{k \in \{1, 2, \dots, K\}} h_{\omega}^{(k)}(\mathbf{x}).$$

$$h_{\omega}^{(0)}(\mathbf{x}) = P(y = 0 | \mathbf{x}; \omega)$$

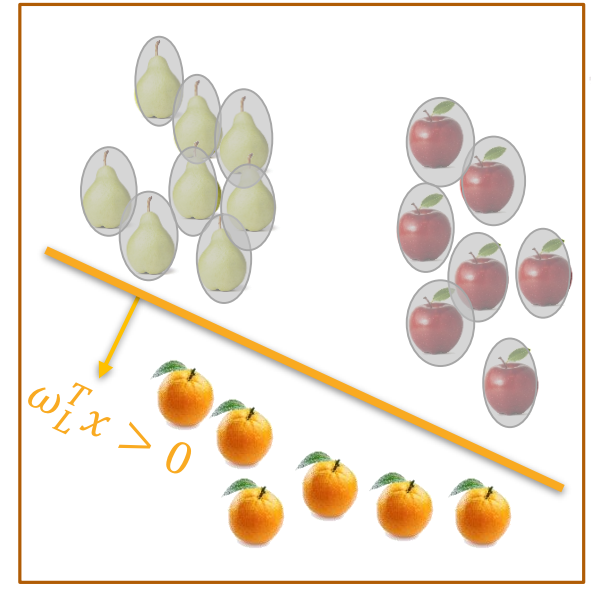
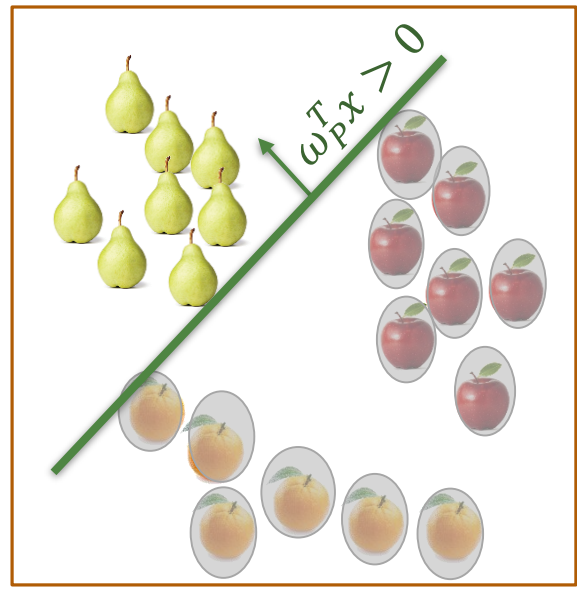
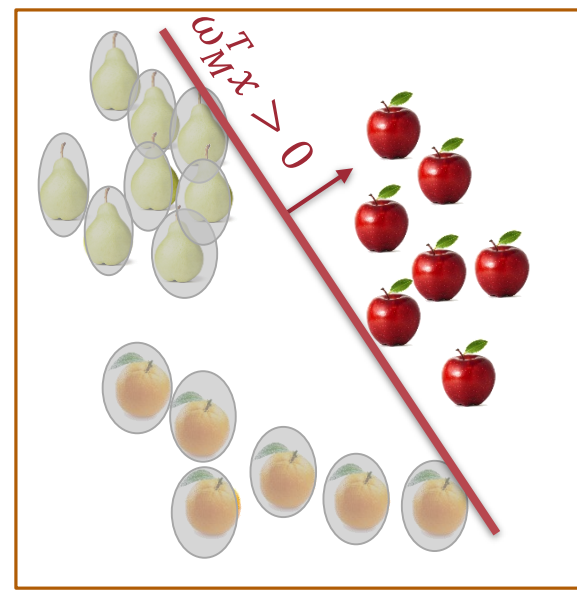
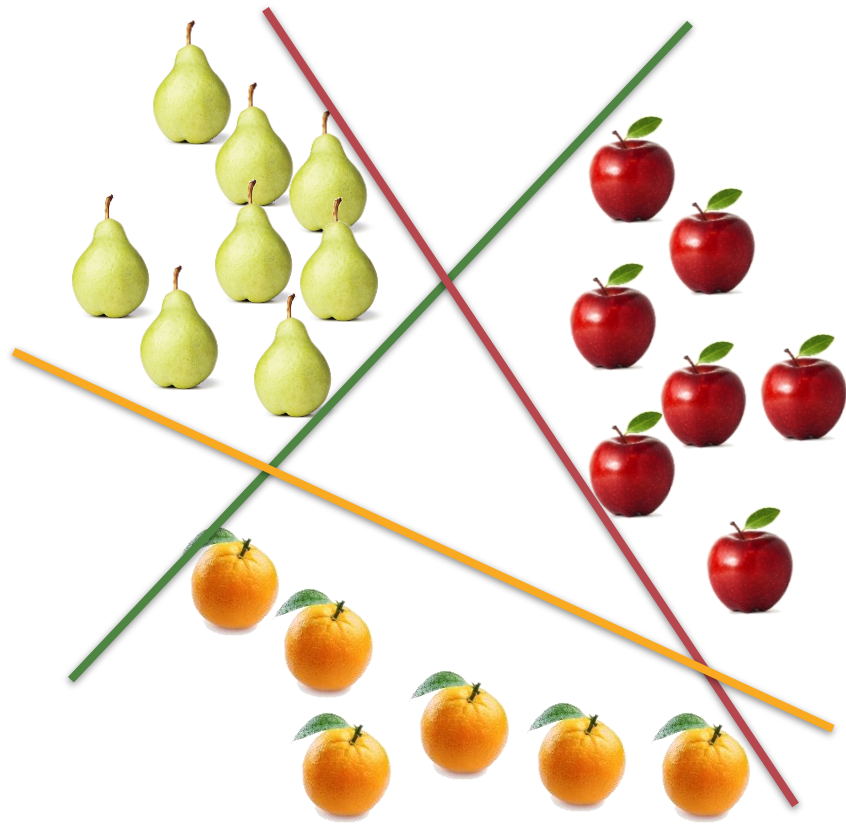
$$h_{\omega}^{(1)}(\mathbf{x}) = P(y = 1 | \mathbf{x}; \omega)$$

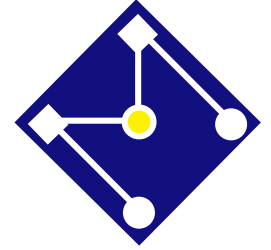
...

$$h_{\omega}^{(C)}(\mathbf{x}) = P(y = C | \mathbf{x}; \omega)$$



ONE-VS-ALL





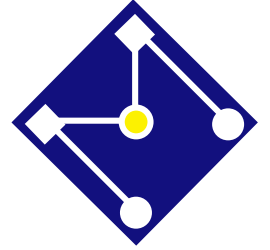
ONE-VS-ONE

Treinamos

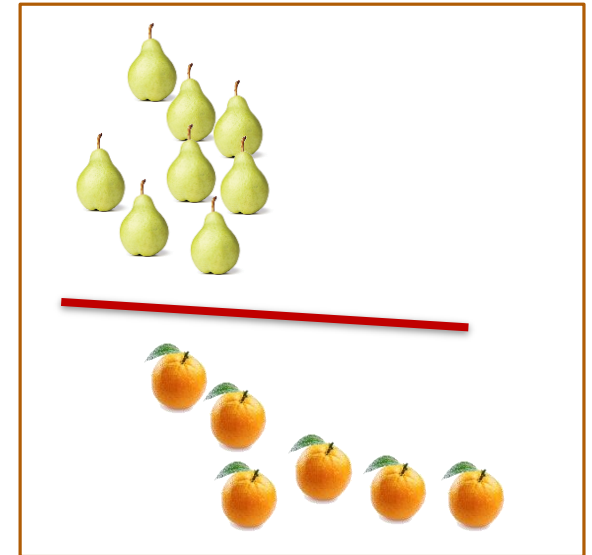
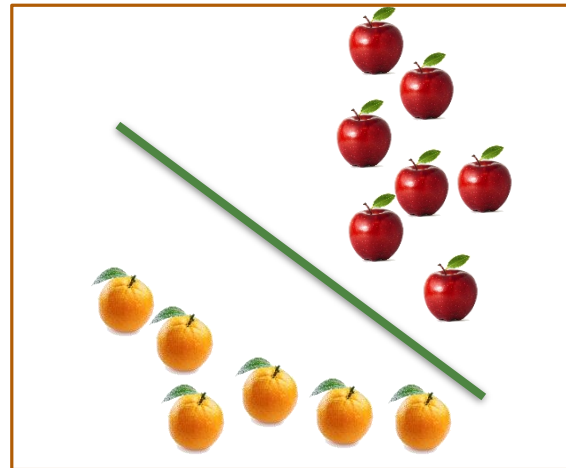
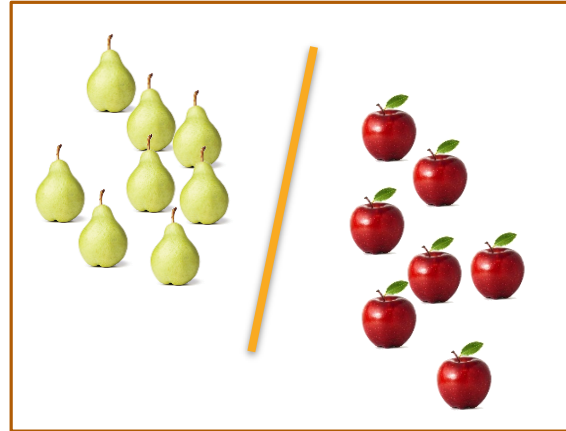
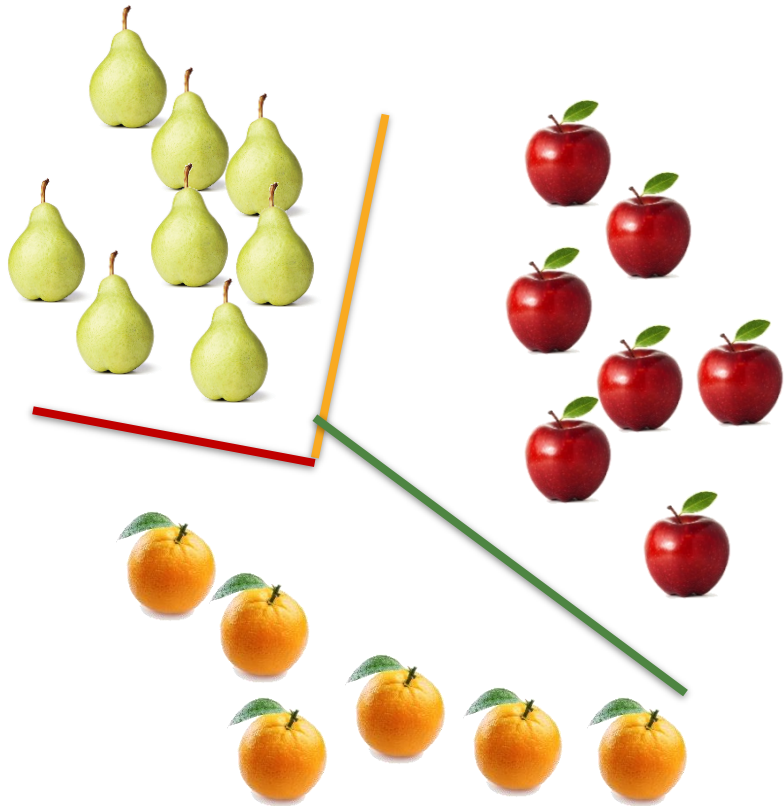
$$\binom{K}{2} = \frac{K(K-1)}{2}$$

modelos de classificação binária separados. Para qualquer novo exemplo \mathbf{x} que desejamos prever escolhemos a classe com a pontuação máxima:

$$\hat{y} = \arg \max_{k \in \{1, 2, \dots, K\}} h_{\omega}^{(k)}(\mathbf{x})$$



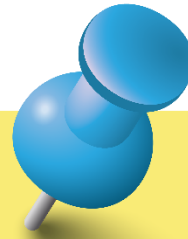
ONE-VS-ONE





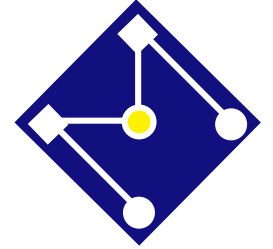
REVISÃO





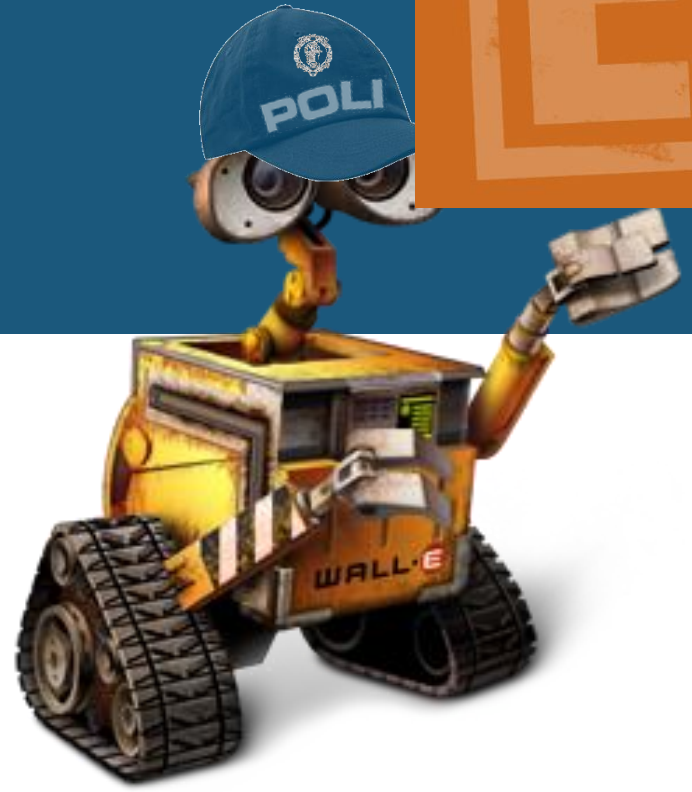
Lição de casa

Estude os slides da aula de hoje e refaça os Notebooks.





ENOUGH IS ENOUGH!



ACABOU... |