

Regressão Logística

Classificação -

online transaction \rightarrow fraudulento
 \rightarrow Não

e-mail \rightarrow spam
 \rightarrow não

tumor \rightarrow maligno
 \rightarrow benigno

Saída: $y \in \{0, 1\}$ classif. binária

negative class,
benign

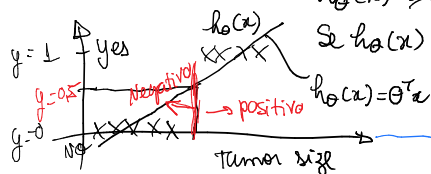
positive class,
malignant

$y \in \{0, 1, 2, 3\}$

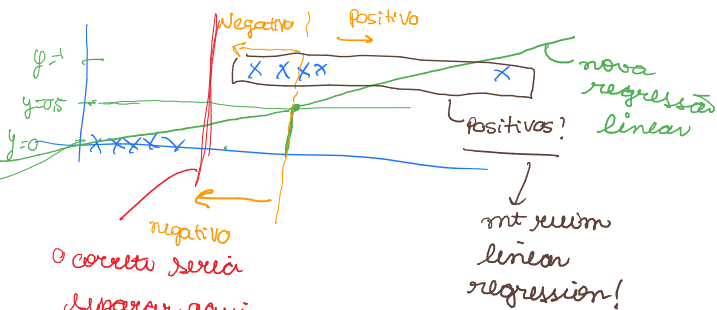
\hookrightarrow multiclass classif. problem

Se eu usasse linear regression, poderia estabelecer: Se $h_0(x) \geq 0,5$, $y=1$

Se $h_0(x) < 0,5$, predict $y=0$



\rightarrow um linear reg. com esse exemplo, a linha fica pior



O correto seria
 Supraer aqui,
 mas adicionar
 um training example à direita alterou
 significativamente a seta.

Logistic Regression

$0 \leq h_0(x) \leq 1$
 classification algorithm,
 Apesar do nome,
 output é
 DISCRETE VALUE,
 ao contrário de
 uma regressão.

Binary Classification

$x^{(i)} \rightarrow y^{(i)}$ será o label p1
 aquele training example.

Qual hipótese usaremos?

- Queremos nosso classificador de hipóteses

Entre zero e 1.

$$0 \leq h_0(x) \leq 1$$

linear reg:

$$h_0(x) = \theta^T x$$

$$0 \leq g(z) \leq 1$$

logistic:

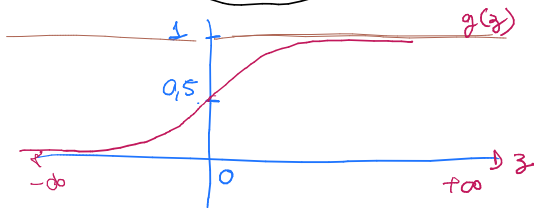
$$h_0(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid function
or
logistic function

$$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$

preciso ajustar θ aos
nossos dados.
Essa hipótese
permite as
predições



Interpretação do output da hipótese:

$h_0(x) \rightarrow$ probabilidade estimada de que
 $y=1$ para a entrada x .

Ex:

$$x = \begin{bmatrix} x_c \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tamanho do tumor} \end{bmatrix}$$

$$h_\theta(x) = 0, \neq$$

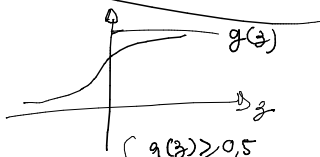
a prob. de ser maligno é $\neq 0$:

$$h_\theta(x) = P(y=1 | x; \theta)$$

Prob. de $y=1$, dado x , parametrizado por θ .

Decision Boundary

$$h_\theta(x) = g(\theta^T x)$$



Queremos tentar entender melhor quando essa hipótese prevê que $y=1$ ou $y=0$.

$$h_\theta(x) = g(\theta^T x)$$

$\left\{ \begin{array}{l} g(z) \geq 0,5 \\ \text{quando } z \geq 0 \end{array} \right.$

$g(\theta^T x) \geq 0,5$ se $\theta^T x \geq 0 \rightarrow$ prevê 1

Previzão $y=0$ se $h_\theta(x) < 0,5$

\hookrightarrow se $\theta^T x < 0$.

Se temos a f.c.:

$$h_0(x) = g(\underbrace{\theta_0 + \theta_1 x_1 + \theta_2 x_2}_z)$$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

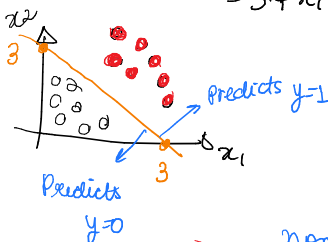
Ele vai prever $y=1$ se

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

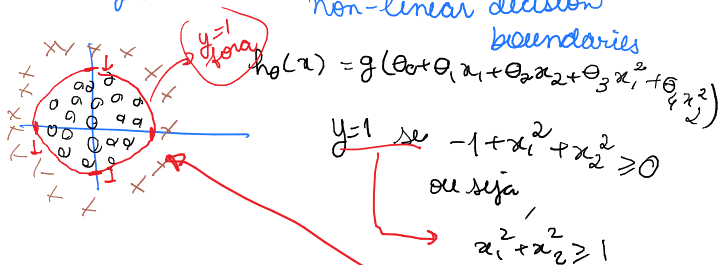
E prever $y=0$ se

$$-3 + x_1 + x_2 < 0$$

$$x_1 + x_2 \geq 3$$



non-linear decision boundaries



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \dots)$$

Pergunta:

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

How do we fit parameters θ ?

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$x_0 = 1, y \in \{0, 1\}$$

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

Essa cost-function
é não-convexa.

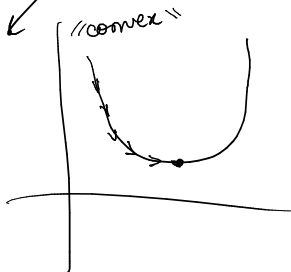
$$\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$

non-convex function - many local optima



non-convex

(many local optima)



Com GD, n' é garantida
a convergência p/ mínimo
global.

O problema dessa não-linearidade que aparece
aqui é ^{per} uma de não-convexa.

Preciso de uma cost-function
convexa e que possamos aplicar gradient
descent.

métodos preditivos

- Classif. / ^{previs} ^{notáveis}
- Regressão _(previs)

métodos descritivos

- associação
- agrupamento
- detecção de desvios
- padrões sequenciais
- sumarização

Relação entre a idade e se vai ou não pagar um crédito
Pagar (atributo-classe)

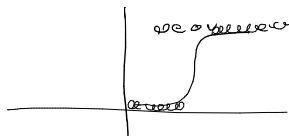
Idade
(atrib. previsor)

Base de
treinam.

- naive bayes - análise a probabilidade
 - ↳ risco alto, baixo, moderado
- decision trees
 - ↳ classificação
- aprend. por exemplos
 - ↳ não gera um modelo!
- KNN - baseado em instâncias
 - ↳ faz cálculos da distância e vê o q tem a menor distância

Regressão logística

↳ encontrar a melhor função que vai desenhar o "S" do gráfico



Avaliar o algoritmo:

- base de teste \neq base de treinamento
↳ com isso, tenho erros de acerto do meu modelo.

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x + \dots)}} = \textcircled{p}$$

Para eu construir o algoritmo \rightarrow transformação LOGIT

$$\log\left(\frac{p}{1-p}\right) = b_0 + b_1 x$$

\rightarrow o objetivo é encontrar a melhor linha

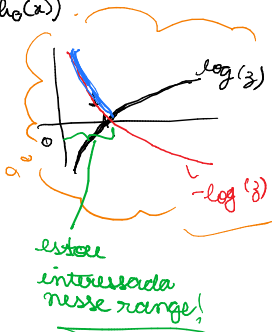
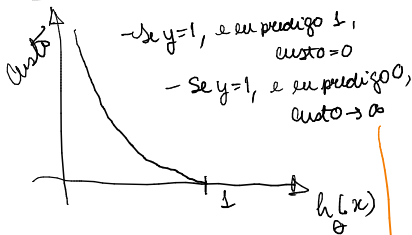
os coeficientes correspondem a pesos p / cada atributo!

Preço de uma outra cost-function por a da regressão linear fica altamente η -convexa!

Se a regressão dá um valor $h_\theta(x)$,
nosso custo será

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{se } y=1 \\ -\log(1-h_\theta(x)) & \text{se } y=0 \end{cases}$$

Se $y=1$: cost-function será $-\log(h_\theta(x))$.



Se $h_\theta(x)=0$

$P(y=1 | x; \theta) = 0$ mas $y=1$,

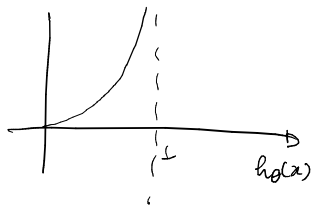
penalizo com custo no alto!

hipótese tá dizendo

que a chance de y ser 1 = 0.

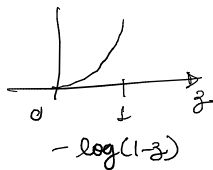
("a chance do seu tumor ser maligno é zero" mas na verdade $y=1$)

Se $y=0$, o custo será $-\log(1-h_\theta(x))$



Se o predigo $y=1$, vai
p/ o custo.

Se predigo $y=0$,
custo será zero.



Com uma escolha particular de
cost-function, isso nos dará um problema de
otimiz. convexa.

$J(\theta)$ será convexa e livre de ótimos locais

Gradientes Conjugados, BFGS e L-BFGS

↳ formas + sofisticadas de otimizar θ
q podem substituir o gradient-descent.

Primeiro temo que escrever a fç que avalia
as seguintes funções p/ um dado θ .

$$J(\theta)$$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

Multiclass classification problems

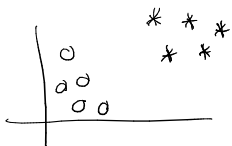
- one vs. all

↳ e-mail em folders diferentes são automaticamente tag e-mails

work $y=1$
family $y=2$
friends... $y=3$
hobby $y=4$

Classificação

Binária

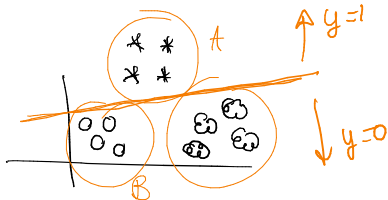


multiclass
multiclass classification
classification



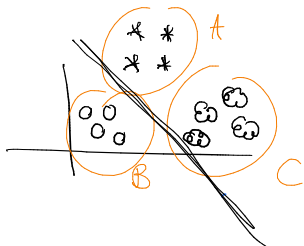
⇒ I turn it into 3 separate classification problems.

A é o positive e o resto é negative



Pergunto:
é da classe A?

trato B como
positivo



trato C
como
positivo



onde
trato i como a classe positiva

to summarize, we fit 3 classifiers

$$h_{\theta}^{(i)}(x) = P(y=1 \parallel x, \theta)$$

qual a probabilidade de y
ser da classe i , dado
 x parametrizado por θ

train a logistic regressor classifier $h_{\theta}^{(i)}(x)$
for each class (i) to predict the probability
that $y=i$

- on a new input x , to make a prediction,
pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

Assim, se tenho k classes e uso one vs. all,
tenho que treinar k classificadores.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{se } y=1 \\ -\log(1-h_{\theta}(x)) & \text{se } y=0 \end{cases}$$

Mas y é sempre 0 ou 1.

Isso equivale a escrever:

$$\text{cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

Esta cost function foi derivada da estatística
usando Princípio da MLE

↓
Serve para achar eficientemente
parâmetros para diversos modelos
Além de ser convexa.

To fit parameters θ :

$$\min_{\theta} J(\theta) \rightarrow \text{get } \theta$$

To make a prediction given new x :
Output $h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$

isso vai
medar
 $P(y=1|x;\theta)$

Quer minimizar cost - função por GDsc.

$$\text{Quero } \min_{\theta} J(\theta)$$

Repeat {

$$\theta_j := \theta_j - \alpha \underbrace{\frac{\partial}{\partial \theta_j} J(\theta)}$$

}

$$\frac{1}{m} \sum_{i=1}^m (h(\theta^{(i)}) - y^{(i)}) x^{(i)}$$

Algoritmo idêntico a reg. linear!

O que muda da linear p/ a logística
é que a definição da nossa hipótese
mudou.

Antes	Agora:
$h(x) = \theta^T x$	$h(x) = \frac{1}{1 + e^{-\theta^T x}}$

~ esquecer do feature scaling!

Vectorized implementation:

$$h = g(X\theta) \quad \begin{matrix} \text{real} \\ \nearrow \text{predito} \end{matrix}$$
$$J(\theta) = \frac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h))$$
$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$$

Optimization w/ logistic regression
even much more quickly than GD.

→ E os algoritmos escalam melhor
w/ very large ML problems.

W/ cada iteração, o GD calcula

$$J(\theta) \text{ e } \frac{\partial}{\partial \theta_j} J(\theta) \quad \forall j=0,1,\dots,n$$

Optimization Algorithms

→ GD

→ Conjugate Gradients
→ BFGS
→ L-BFGS

- n precisa pegar manualmente um α

- often faster than GD

- MORE COMPLEX!

daí
as derivadas parciais

function [jval, gradient] = costFunction(theta)

$$jval = (\theta(1) - 5)^2 + \dots + (\theta(2) - 5)^2;$$
$$gradient(1) = 2 * (\theta(1) - 5)$$

daí $J(\theta)$

isso em
implemento!

fminunc - function minimization unconstrained

@ \Rightarrow pontos p/a cost function

ver se tem no matlab

fminunc

help fminunc

Se é large ML problem, use esses
algoritmos em vez de GD