# Logistic Regression in Python
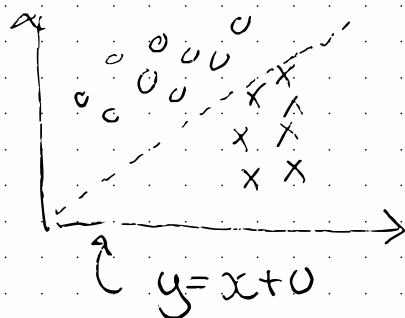
Classification → supervised learning
                        (supervised)

decide between classes
   ↳ cat vs dog, numbers from pics

2d classification ⤳

Classes can be sep.
by a straight line



$y = x + 0$
$0 = x - y$

Linear line can also be written as ↰

If $h(x,y) = x - y$ then $h(x,y) \begin{cases} > 0 = \text{'o'} \\ < 0 = \text{'x'} \end{cases}$

difficulty: if $h(x,y) = 0$ then how do we define
                point?

$h()$ is a linear combo. of weights $\theta_1 \theta_2 \theta_3 \ldots = \underset{D \times 1}{\Theta}$
inputs $x_1 x_2 x_3 = \underset{u \times D}{X}$ and targets $y_1 y_2 y_3 = \underset{n \times 1}{\bar{y}}$
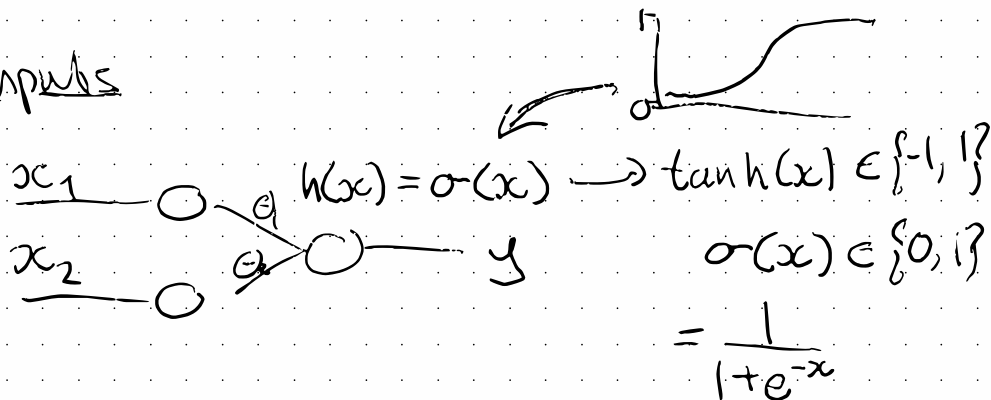   ↳ $\underset{u \times 1}{\hat{y}} = X\Theta$

⌐ note how in 3d, decision boundary becomes a line.

⌐ in 4d; 3d-hyperplane

• Can we calculate the output of a neuron?

  ↳ neurons, machine until threshold is passed & becomes active.

inputs



$x_1$ —O
         a    $h(x) = \sigma(x) \longrightarrow \tanh(x) \in \{-1, 1\}$
$x_2$ —O  O— y        $\sigma(x) \in \{0, 1\}$

                              $= \dfrac{1}{1+e^{-x}}$

thus, if $\theta^T x \uparrow\uparrow\uparrow$ then closer to 1.
      if $\theta^T x \downarrow\downarrow\downarrow$ then closer to 0.

note that $\sigma(\theta^T x) \in \{0, 1\} = P(y=1 \mid x)$

where $P(y=1 \mid x) + P(y=0 \mid x) = 1$

      ↳ we predict the class w/ higher likelihood.

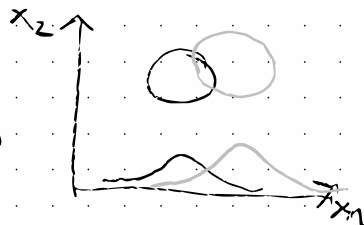Assumption of Logistic Regression
= data can be separated by a plane.
= "linearly separable"

How do we find weights, $\Theta$?

consider 2 gaussian dist. clouds

$$P(x) = \frac{1}{\sqrt{2(\pi)^D |Z|}} e^{-\frac{1}{2}(x-\mu)^T Z^{-1}(x-\mu)}$$



↳ from bayes rule $P(Y=1|X) = \dfrac{P(X|y=1)\,P(y=1)}{P(X)}$

$$\& \ P(Y=0|X) = \dfrac{P(X|y=0)\,P(y=0)}{P(X)}$$

if we combine the two probabilities...

$$P(y=1|X) = \cfrac{1}{1 + \cfrac{P(X|y=0)\,P(y=0)}{P(X|y=1)\,P(y=1)}}$$

⎫ already
⎬ looks a
⎭ little like
  logistic reg

where $\ln\left(\dfrac{P(X|y=0)\,P(y=0)}{P(X|y=1)\,P(y=1)}\right) = -\Theta^T x + b$

Set $p(y=1) = 1-\alpha$ & $p(y=0) = \alpha$

then $= \ln(P(x|y=0) + \ln(\alpha) - \ln(P(x|y=1)) - \ln(1-\alpha)$

Gaussian dist.

If we add in gaussians then we get...

$$= \frac{-1}{2}(x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x + \mu_0^T \Sigma^{-1} \mu_0)$$

$$\frac{+1}{2}(x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1)$$

$$+ \ln\left(\frac{\alpha}{1-\alpha}\right)$$

$$\rightarrow w^T = (\mu_1^T - \mu_0^T)\Sigma^{-1}$$

$$b = \frac{1}{2}\mu_0^T \Sigma^{-1}\mu_0 - \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 - \ln\frac{\alpha}{\alpha-1}$$

Notes:

· Linear Discriminant analysis performed above.

... This is totally useless!

↳ gradient descent used instead.

→ Loss function, cross entropy

$$J = -\sum_{n=1}^{N} y_n \log(\hat{y}_n) + (1-y_n)\log(1-\hat{y}_n)$$

note that when $y=1$, ② $=0$ & when $y=0$, ① $=0$

$\hookrightarrow$ if $y=1, \hat{y}=1 \cdots J=0$

$y=0, \hat{y}=0 \cdots J=0$ — small error

$y=1, \hat{y}=.99 \cdots J=0.11$ — small error

$y=1, \hat{y}=.5 \cdots J=0.69$ — larger error

Maximum Likelihood

$\hookrightarrow$ $p(y=1)=P \qquad p(y=0)=1-P$

say $N=10, N_1=7, N_0=3$

$L = P^7(1-P)^3 \qquad \leftarrow$ maximize $L$ w.r.t. $P$

$\log L = \log[P^7(1-P)^3]$

$\qquad = 7 \log P + 3 \log(1-P)$

$\dfrac{\partial \log L}{\partial P} = \dfrac{7}{P} + \dfrac{-3}{1-P} = 0 \quad \cdots \quad \dfrac{1-P}{P} = \dfrac{3}{7}$

$$\cdots \quad \frac{1}{p} - 1 = \frac{3}{7}$$

$$P = \frac{7}{10} = P(H) \quad \checkmark$$

↳ Now that we understand this w/ 1s & 0s, apply to logistic reg.

$$P(y=1|x) = \sigma(\Theta^T x) = y$$

$$L = \prod_{n=1}^{N} y_n^{t_n} (1-y_n)^{1-t_n} \quad \leftarrow \text{ target } \hat{y}$$

↳ take log likelihood, get cross entropy
- Loss function

## How to optimize weights

↳ LDA shown before only works when:

　　1. data is normally dist.

　　2. equal cov.

Try LR approach? can't solve for the
　　　　　　　　weights

Solution, gradient descent!

↳ start w/ loss function

$$J = -\sum_{n=1}^{N} y_n \log(\hat{y}_n) + (1-y_n)\log(1-\hat{y}_n)$$

↳
$$\frac{\partial J}{\partial \Theta_i} = \sum_{n=1}^{N} \frac{\partial J}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial a_n} \frac{\partial a_n}{\partial \Theta_i} \quad \leftarrow a_n = \Theta^T x_n$$

$$\boxed{\frac{\partial J}{\partial \hat{y}_n} = -y_n \frac{1}{\hat{y}_n} - (1-y_n)\frac{1}{1-\hat{y}_n}}$$

$$y_n = \sigma(a_n) \qquad \boxed{\frac{\partial \hat{y}_n}{\partial a_n} = \hat{y}_n(1-\hat{y}_n)} \xleftarrow{\text{deriv of } \sigma}$$

$$\boxed{\frac{\partial a_n}{\partial \Theta_i} = x_{ni}}$$

assemble...

$$\frac{\partial J}{\partial \Theta_i} = -\sum_{n=1}^{N} \frac{y_n}{\hat{y}_n} \hat{y}_n(1-\hat{y}_n)x_{ni} - \frac{1-y_n}{1-\hat{y}_n} \hat{y}_n(1-\hat{y}_n)x_{ni}$$

$$\therefore \quad \boxed{\frac{\partial J}{\partial \Theta} = \sum_{n=1}^{N} (y_n-\hat{y}_n)x_n}$$

vectorize save more : $\boxed{\dfrac{\partial J}{\partial \Theta} = X^T(Y - \hat{Y})}$ ✓

bias term? add column of 1s to X

↳ remember from Lin. Reg. class :

$$\Theta := \Theta - \eta \frac{\partial J}{\partial \Theta} J(\Theta)$$

where $\dfrac{\partial J}{\partial \Theta}$ $J(\Theta)$ = boxed equation above

pseudocode:

· import data
· z = XΘ
· learning rate = 0.1 (or else)
· steps
↳ for i in steps
    $\Theta := \Theta - \eta * [X^T(Y - \sigma(X\Theta)]$

# Practical Issues

So far... training① & predicting②

↳ problems w/ overfitting
↳ problems where ideal weights are inf
↳ problems w/ non-linear decision boundaries

# Interpreting the weights

↳ Still similar to Lin. reg. except features w/ higher weights contribute more to $P(y=1|x)$

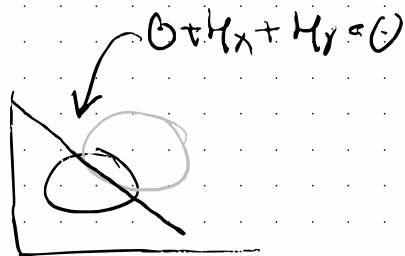you can also look @ odds $\boxed{\dfrac{P(y=1|x)}{P(y=0|x)}}$

$$= exp(\omega^T x)$$

← linear regression

log odds ∴ $= \omega^T x$

→ L2 regularisation

⤷ prevent overfitting

Why use regularisation?

⤷ cloud problem

$0 + M_x + M_y = 0$



for the line, $x = -y$, do weights matter?

⤷ due to the math, the best weights are
  $\inf \, x = -\inf \, y$ ⟵ difficult to calculate

use regularisation to penalize large weights

re. $J_{reg} = J + \lambda \frac{1}{2} w^T w = J + J_{L2}$

⤷ smoothing param chosen manually.

$\dfrac{\partial J_{L2}}{\partial w} = \lambda w$

∴ $\dfrac{\partial J_{reg}}{\partial w} = x^T(\hat{y} - y) + \lambda w$

- Probabilistic perspective
  - ↳ normally, maximize likelihood
    regularization, maximize posterior

$$\underbrace{\text{max. a posteriori}}$$
$$\text{(MAP)}$$

→ L1 regularization

↳ normally we want $X$ to be tall
but what if $X$ is wide?

- Goal: select features w/ most importance
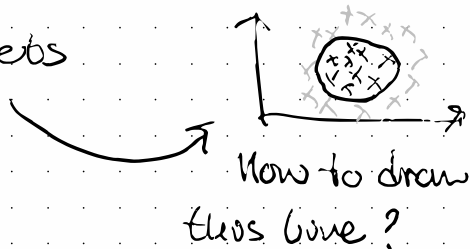  set others to 0

once again $J_{reg} = J + \lambda w$

$$\qquad \qquad \qquad \qquad ↳ \frac{\partial}{\partial w} = \lambda \, \text{sign}(w)$$

$$\therefore \quad \frac{\partial J_{reg}}{\partial w} = X^T(\hat{Y} - Y) + \lambda \, \text{sign}(w)$$

→ Donut Problem

↳ Logistic Regression is good at classifying linearly seperable data.

Problem: non-linear sets

How to draw this line?

Similar to how we solve for non linear problems w/ linear regression

↳ in dataset add bias term & radius of each point

↳ grad desc based on L2 reg.