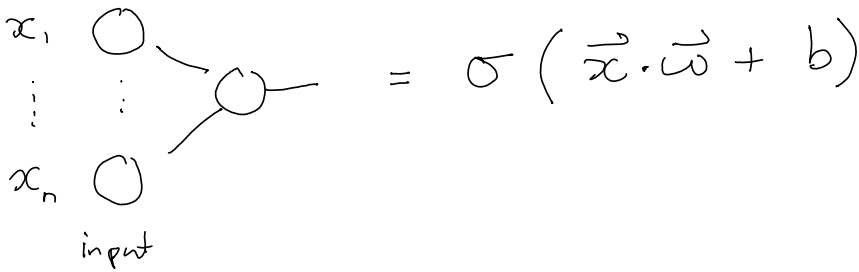


Classification

1.1) Logistic regression

The perceptron we introduced before represents a linear classifier



where $\sigma(z) = \Theta(z)$ for a hard classifier and (e.g.) $\sigma(z) = \frac{1}{1 + e^{-z}}$ for a soft classifier

hard \rightarrow either 0 or 1

soft \rightarrow probability the result is 0 or 1

it's convenient for notation to write

$$\vec{X} = (1, \vec{x}) \quad \text{and}$$

$$\vec{W} = (b, \vec{w})$$

such that the output is simply

$$\sigma(\vec{X} \cdot \vec{W})$$

Concretely we may write the probability that a result y (also called score) is 1 for a point \vec{x} with the model \vec{W} as

$$P(y=1 \mid \vec{x}, \vec{W}) = \frac{1}{1 + e^{-\vec{x} \cdot \vec{w}}}$$

and

$$P(y=0 \mid \vec{x}, \vec{W}) = 1 - P(y=1 \mid \vec{x}, \vec{W})$$

This can be seen clearly by analogy with a Boltzmann distribution of a two state (non-interacting) system

$$P_i = Q^{-1} \exp \left(- E_i / (kT) \right)$$

$\underbrace{P_i}_{\text{probability for state } i=0,1}$
 \uparrow energy of state i
 \nwarrow Boltzmann constant times Temp.

where $Q = \sum_{i=0}^1 \exp \left(- E_i / (kT) \right)$

$$P_1 = \frac{e^{-E_1/(kT)}}{e^{-E_0/(kT)} + e^{-E_1/(kT)}}$$

$$= \frac{1}{1 + \exp \left(- \underbrace{(E_0 - E_1)}_{\Delta E} / (kT) \right)}$$

and $P_0 = 1 - P_1$. So we see that

$\vec{X} \cdot \vec{\omega}$ plays the rôle of

$$\Delta E / (kT) .$$

We can write the likelihood that a model \vec{w} matches a dataset

$$\mathcal{D} = \{ \vec{x}_i, y_i \} \quad i=1, \dots, n$$

as

$$P(\{ \vec{x}_i, y_i \} | \vec{w}) =$$

$$\prod_{i=1}^n \left(\sigma(\vec{x}_i \cdot \vec{w}) \right)^{y_i} \left(1 - \sigma(\vec{x}_i \cdot \vec{w}) \right)^{1-y_i}$$

the loss function for this model is then $-\log(\text{likelihood})$

$$\mathcal{L}(\vec{w}) = - \sum_{i=1}^n y_i \log(\sigma(\vec{x}_i \cdot \vec{w})) + (1-y_i) \log(\sigma(\vec{x}_i \cdot \vec{w}))$$

which is known as the cross entropy.

As always the minimum of the loss function gives the optimal set of weights for our model. It is easy to compute

$$\frac{\partial J}{\partial \vec{w}} = \sum_i \vec{x}_i (\sigma(\vec{x}_i \cdot \vec{w}) - y_i)$$

$$= \sum_i \vec{x}_i (\hat{y}_i - y_i)$$

using $\sigma(\vec{x}_i \cdot \vec{w}) = \hat{y}_i$.

From where we may now turn to gradient descent methods. NB

$\frac{\partial J}{\partial \vec{w}} = 0$ is a complicated set of transcendental equations so we must turn to numerical methods.

1.2) SoftMax regression

An obvious generalisation of the binary classifier is a discrete set of possible outcomes, $y = 0, 1, \dots, M$ (often referred to as labels)

The generalisation can be seen obtained by considering a Boltzmann distribution for a multi-state system.

$$P_I = Q^{-1} \exp(-E_I / (kT))$$

↖
probability of a
state I

$$Q = \sum_{I=1}^M \exp(-E_I / (kT))$$

$$\Rightarrow \sum_{I=1}^M p_I = 1$$

translating to a model of weights \vec{w}_I
 we write the probability of a state I
 as

$$P(y=I \mid \vec{x}, \{\vec{w}_J\}) = \frac{e^{-\vec{x} \cdot \vec{w}_I}}{\sum_J e^{-\vec{x} \cdot \vec{w}_J}}$$

$$= \frac{1}{1 + \sum_{J \neq I} \exp(-\vec{x} \cdot (\vec{w}_I - \vec{w}_J))}$$

we may think of this as a vector
 of probabilities $\vec{p}(\vec{x}, \{\vec{w}_J\})$.

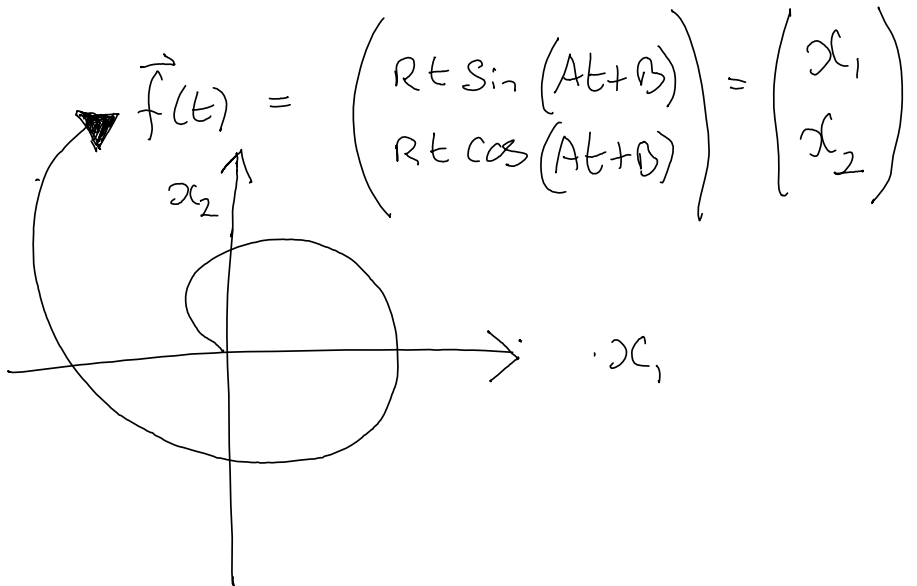
We can understand the notation
 best by following an example.

1.2.1) Soft-Max regression example

^{Jaupyter}

[see notebook] We consider the case

of a 2d distribution of points, each of which may have a colour red, blue or yellow. We generate data points that live on a spiral starting at the origin and generate 3 spiral "arms" of different colours,



data input - $x_{\alpha i}$,

data output - y_{α} ,

where $\alpha = 1, \dots, N$

$i = 0, 1$

$y_{\alpha} \in \{0, 1, 2\}$

Let's first try a linear model

$$z_I = x \cdot w_I + b_I$$

where $I = (0, 1, 2)$

$$\sigma_I(z_I(x)) = \frac{e^{-z_I}}{\sum_{j=0,1,2} e^{-z_j}}$$

The probabilities of an input value

z corresponding to an output value
 $I = (0, 1, 2)$ is just

$$p_I(x) = \sigma_I(z_I(x))$$

For our dataset for training, x_i ,
the probabilities for the model \vec{w}_I, b_I
are

$$p_{\alpha I} = \sigma_I(z_I(x_\alpha))$$

where the scores z are

$$z_{\alpha I} = x_{\alpha i} w_{iI} + b_I$$

$$\Rightarrow p_{\alpha I} = \frac{\exp(z_{\alpha I})}{\sum_j \exp(z_{\alpha j})}$$

There's a quick shortcut to the loss function, \mathcal{L} , from here

$$\mathcal{L}_\alpha = -\log \left(p_\alpha(I=y_\alpha) \right)$$

we may also add a regularisation loss coming from the model weights w_{iI} to prevent preferring large weights.

$$\mathcal{L}^{\text{reg}} = \frac{r}{2} \sum_{i,I} w_{iI}^2$$

(therefore
reducing
overfitting)

such that

$$\langle \mathcal{L} \rangle = \frac{1}{N} \sum_{\alpha=1}^N \mathcal{L}_\alpha + \mathcal{L}^{\text{reg}}$$

The gradient is computed from

$$d\mathcal{L}_{\alpha I} = \left(p_{\alpha I} - \delta_{\alpha(I=y_\alpha)} \right) \frac{1}{N}$$

where the updated parameters are computed from

$$dW_{iI} = x_{ai} dZ_{aI}$$

$$dB_I = \sum_a Z_{aI}$$

and $dW_{iI}^{\text{reg}} = r W_{iI} \leftarrow \frac{d}{dw_{iI}} \left(\frac{r}{2} w_{iI}^2 \right)$

Now update via

$$W_{iI}^{\text{new}} = W_{iI} - \eta (dW_{iI} + dW_{iI}^{\text{reg}})$$

$$B_I^{\text{new}} = B_I - \eta dB_I$$

See the Jupyter notebook for the implementation with a hidden layer.

Application : Supersymmetry and Searches at Colliders

Supersymmetry (SUSY) was (is) a popular model for physics beyond the Standard Model (BSM).

This proposal introduces partners for every particle in the SM through symmetry between bosons (integer spin) and fermions ($\frac{1}{2}$ -integer spin).

It may seem a fairly innocuous extension but the underlying principle is a unique extension of the

Poincaré group for spacetime.

Poincaré algebra in $SO(1,3)$

P^μ - translations

$M^{\mu\nu}$ - Lorentz boosts + rotations

satisfy:

$$[P^\mu, P^\nu] = 0$$

$$-i [M_{\mu\nu}, P_\rho] = g_{\mu\nu} P_\rho - g_{\nu\rho} P_\mu$$

$$\begin{aligned} -i [M_{\mu\nu}, M_{\rho\sigma}] &= g_{\mu\rho} M_{\nu\sigma} - g_{\nu\rho} M_{\mu\sigma} \\ &\quad + g_{\nu\sigma} M_{\mu\rho} - g_{\mu\sigma} M_{\nu\rho} \end{aligned}$$

Supersymmetry is represented by an operators Q, \bar{Q} that relate particles differing by $\frac{1}{2}$ integer spin. Q and \bar{Q} are Weyl spinors and the central objects in the anti-commutation relations

involving the (extended) Pauli matrices $\sigma_{\alpha\dot{\alpha}}^{\mu}$

$$\{Q_{\alpha}, \bar{Q}_{\dot{\alpha}}\} = 2\sigma_{\alpha\dot{\alpha}}^{\mu} p_{\mu}$$

In addition

$$[M^{\mu\nu}, Q_{\alpha}] = \sigma_{\alpha}^{\mu\nu\beta} Q_{\beta}$$

$$[M^{\mu\nu}, \bar{Q}_{\dot{\alpha}}] = (\bar{\sigma}^{\mu\nu})_{\dot{\alpha}}^{\dot{\beta}} Q^{\dot{\beta}}$$

where $\sigma_{\alpha}^{\mu\nu\beta} = \frac{i}{4} (\sigma^{\mu}\bar{\sigma}^{\nu} - \sigma^{\nu}\bar{\sigma}^{\mu})_{\alpha}^{\beta}$

and

$$[Q_{\alpha}, p^{\mu}] = 0$$

$$\{Q_{\alpha}, Q_{\beta}\} = 0$$