# NEURAL NETWORKS

Before specifying how neural nets are parametrised, we will talk about the optimisation algorithm that is used to train them, which is the same we use for the logistic regression model
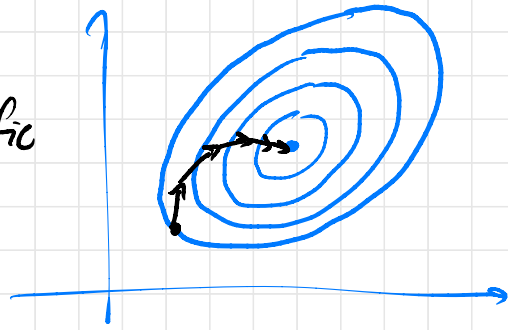
## Gradient Descent

An iterative algorithm to find the minimum of a function $C(\vec{w})$ such that, at iteration $\tau+1$, the parameters $\vec{w}$ are updated as

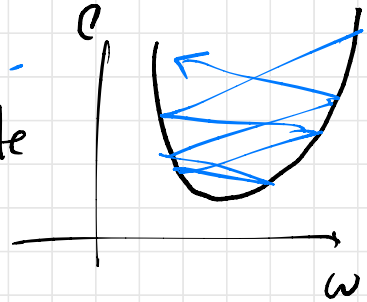$$\vec{w}_{\tau+1} = \vec{w}_{\tau} - \eta \frac{dC(\vec{w})}{d\vec{w}}$$

↳ when arrived at the minimum, the derivative is zero and the parameters don't evolve any more

↳ called the "learning rate" and it could depend on $\tau$

The learning rate has to
be tuned for each specific

case :

- If too large, the
algorithm may miss the minimum.
- If too small, there is a waste
of computational time.

- The state-of-the-art in optimisation of
neural nets is given by variations of
what it's called "Stochastic Gradient Descent".

The motivation is two-fold:

- Standard G.D. requires evaluation of
the full $C(\vec{w}) = \frac{1}{N} \sum_{i=1}^{N} C_i(\vec{w})$. So it
may be costly to evaluate for
large datasets
- Standard G.D. gets easily stucked

in local minima, and thus unable to
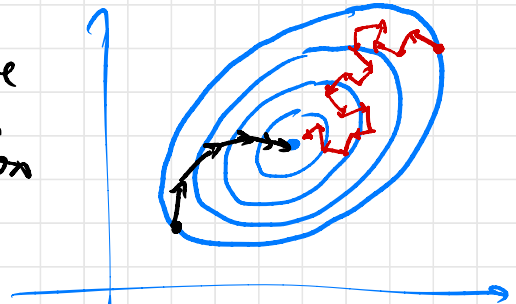find the global minimum

The idea of Stochastic G.D. is to
approximate

$$\frac{dC(\vec{w})}{d\vec{w}} \approx \frac{dC_i(\vec{w})}{d\vec{w}}$$

where "i" is a data point (or a subset of
data points, a.k.a. "batch") randomly
chosen

- The trajectory will have
thus a random direction
at each iteration, <u><u>but</u></u>

on average it will converge to the minimum,
because

$$\langle C_i(\vec{w}) \rangle = \frac{1}{N} \sum_i C_i(\vec{w}) = C(\vec{w})$$

Also, because of this stochastic behaviour,
it is easier for the algorithm to jump
out of a local minimum.

# Neural nets

Let's start with the shape of a linear
model:

$$f(\vec{x}, \vec{v}) = f\left( \sum_{h=1}^{H} v_h \phi_h (\vec{x}) \right)$$

parameters

$$f(a) = \begin{cases} a & \text{in regression problems} \\ sigm(\sigma) & \text{in classification} \\ softm(\sigma) \end{cases}$$

The idea of neural nets is to have
more flexibility / expressivity by letting the
$\phi_h (\vec{x})$ themselves depend on other
parameters

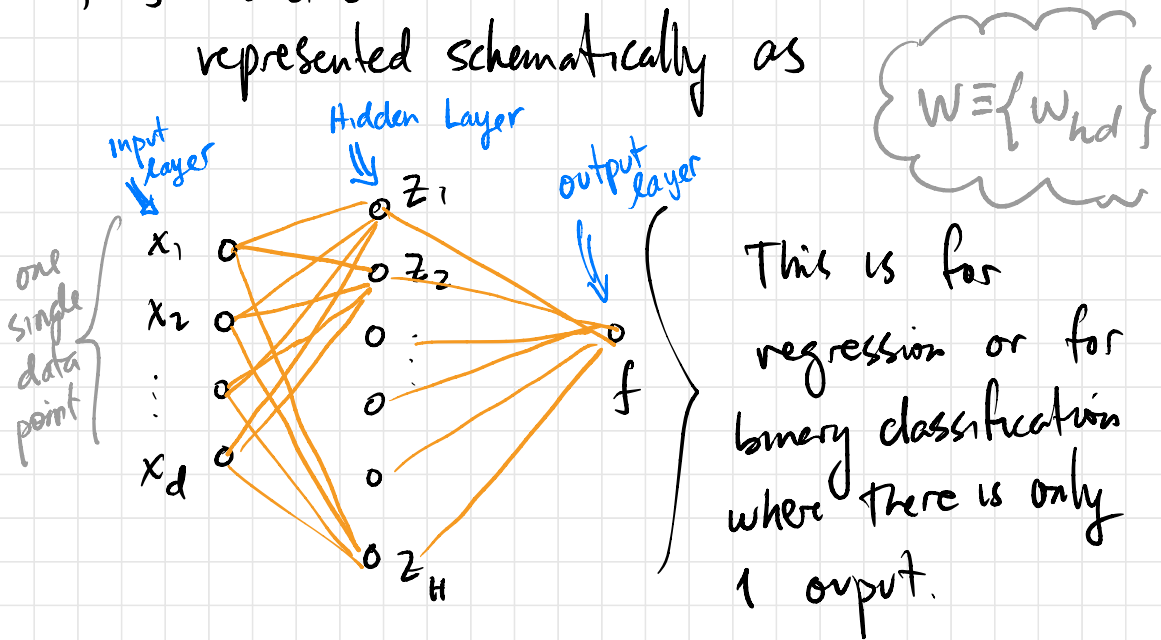$$\phi_h(\vec{x}, \vec{w_h}) = g\left(\vec{w_h}^T \cdot \vec{x} + w_{oh}\right)$$

$\hookrightarrow$ a non-linear "activation" function

So the model becomes

$$f(\vec{x}; \vec{v}, W) = f\left(\sum_{h=1}^{H} v_h \underbrace{g\left(\vec{w_h}^T \cdot \vec{x} + w_{oh}\right)}_{\equiv Z_h \Rightarrow \text{"neurons"}} + v_o\right)$$

This model can be represented schematically as



$\{W \equiv \{w_{hd}\}$

Input layer — Hidden Layer — output layer

one single data point

$x_1$, $x_2$, ..., $x_d$ — $Z_1$, $Z_2$, ..., $Z_H$ — $f$

This is for regression or for binary classification where there is only 1 output.

For multiclass classification we had

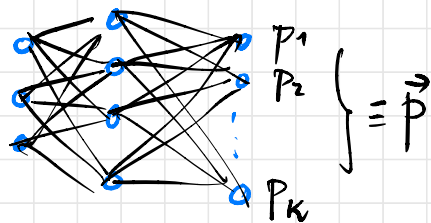$$p(C_k | \vec{x}) = \sigma_k(\vec{a}) \rightsquigarrow \text{softmax}$$

So we have $k-1$ outputs ( since the last one is obtained by

$$p(C_k \mid \vec{x}) = 1 - \sum_{k=1}^{k-1} p(C_k \mid \vec{x})$$

$$P_k \equiv p(C_k \mid \vec{x})$$

$$= f\left( \sum_{h=1}^{H} v_{hk} \, g\left( \sum_{d=1}^{D} w_{nd} \, x_d + w_{ho} \right) + v_{ko} \right)$$

• Now there is a compact way of writing this function containing the predictor for all data points $\vec{x}_i$



$$\left.\begin{array}{l} P_1 \\ P_2 \\ \vdots \\ P_k \end{array}\right\} = \vec{P}$$

$$P = \left\{ \vec{P}(\vec{x}_1), \vec{P}(\vec{x}_2), \ldots, \vec{P}(\vec{x}_N) \right\} \qquad \dim P = N \times k$$

$$X = \left\{ \vec{x}_1, \vec{x}_2, \ldots, \vec{x}_N \right\} \qquad \dim X = N \times D$$

$$W = \left\{ w_{dh} \right\} \qquad \dim W = D \times H$$

$$W_0 = \left\{ w_{ho} \right\} \rightarrow \text{all rows are identical} \qquad \dim W_0 = N \times H$$

$$V = \{ v_{hk} \} \qquad \text{dim } V = H \times K$$

(identical rows) →

$$V_0 = \{ v_{k0} \} \qquad \text{dim } V_0 = N \times K$$

$$\left[\!\!\left[ \; P = f\left( \; g(X \cdot W + W_0) \cdot V + V_0 \right) \; \right]\!\!\right]$$

↳ 1-Hidden Layer Neural network.

What about a 2-hidden layer network?

$$P = f\left( g_2 \left[ \underbrace{g_1(X W_1 + W_{10})}_{\text{has this instead of just } X} W_2 + W_{20} \right] V + V_0 \right)$$

and similarly for neural nets with many hidden layers.

There are many possibilities for the shape of the activation functions (counted 30+ in Wikipedia)



sigmoid          tanh x          "ReLU"

All of them are able to achieve very good
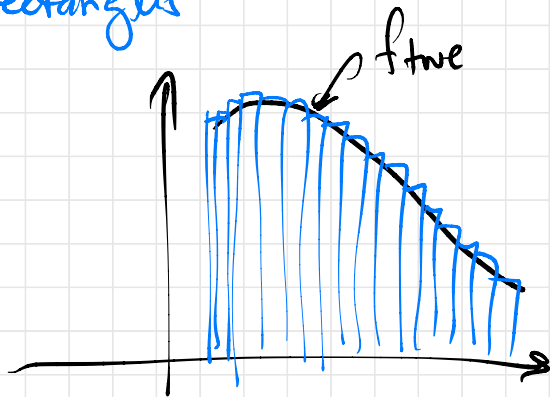results (modulo training problems with some
        of them)

This has to do with:

The Universal Approximation Theorem (Cybenko, '89)

↳ " A feedforward, single hidden layer NN
containing finite # (#) of neurons can
approximate any continuous function under
    mild assumptions on the activation function."

In practice: all it requires is to be able
to approximate any function using an
arbitrary # of rectangles

{ see supplementary

    material }

# Training a 1HLNN

This is just the MLE of the parameters of the network, where the likelihood is the usual cross-entropy (in the case of multiclass classification).

$$C(\vec{w}) = -\sum_{i=1}^{N}\sum_{k=1}^{k} t_{ik} \ln y_k(\vec{x}_i, \vec{w})$$

observed output $\underbrace{\quad}$ prediction $\longleftarrow$

$$D = \{\vec{x}_i, \vec{t}_i\}$$

Differentiating $C(\vec{w})$ wrt all the parameters is straightforward (but long), so I only present here the results:

$$\frac{\partial C}{\partial v_{ok}} = \sum_{i}(Y-T)_{ik}$$

$$\frac{\partial C}{\partial v} = Z^T \cdot (Y-T)$$

$$Z = \{z_{ih}\}$$

$$\frac{\partial C}{\partial w_{oh}} = \sum_i \left[ (Y-T) V^T * Z * (\mathbf{1} - Z) \right]_{ih}$$

element-by-element product

$$\frac{\partial C}{\partial w} = X^T \left[ (Y-T) V^T * Z * (\mathbf{1} - Z) \right]$$