# Neural Networks

Neural nets use composition $\left( f \circ g (x) = f(g(x)) \right)$
to build non-linear functions from linear ones & simple
non-linear functions.

$(s_i, 0, ..., w_{i,j}, ..., 0)$

Like Adaboost: stumps $g_i(x) = \mathbb{1}\{w_i^T x > 0\}$

ensemble $\sum_i \alpha_i g_i(x) = f(x)$



linear        non-linear        linear
              activation
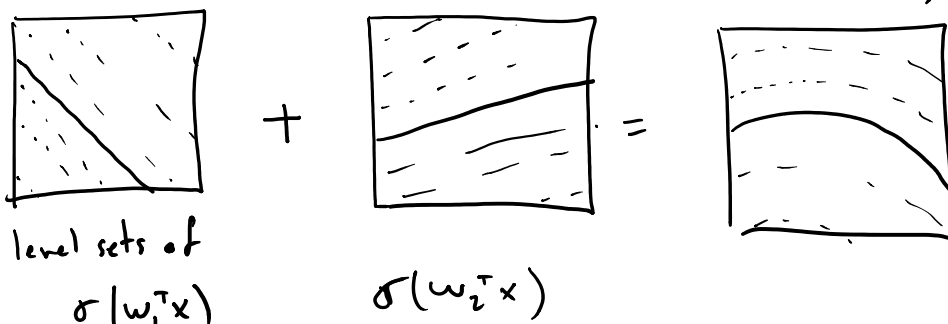
Sigmoid: $\sigma(z) = \dfrac{1}{1+e^{-z}}$

is a smooth surrogate for $\mathbb{1}\{z>0\}$

Hidden units: weight $w_1$    constuct feature    $h_1 = \sigma(w_1^T x)$
vector

$\qquad\qquad\qquad w_2 \qquad \cdots \qquad\qquad\qquad h_2 = \sigma(w_2^T x)$

$\qquad\qquad\qquad \vdots$

$\qquad\qquad\qquad w_H \qquad \cdots\cdots \qquad\qquad h_H = \sigma(w_H^T x)$

Combine $h_1, ..., h_H$ we get $f(x) \parallel \sum_{j=1}^{H} \beta_j h_j(x)$ can look non-linear

(non-linear level sets)



level sets of
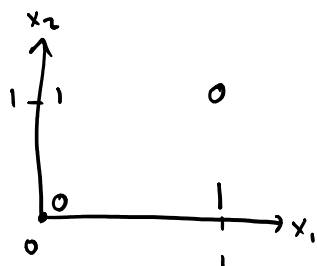$\sigma(w_1^T x)$        $\sigma(w_2^T x)$

Why not $\sigma(z) = z$ (identity)

$$f(x) = \sum_{i=1}^{H} \beta_i h_i(x) = \sum_{i=1}^{H} \beta_i w_i^T x = \left( \overbrace{\sum_{i=1}^{H} \beta_i w_i}^{\tilde{\beta}} \right)^T x$$

Other common option: rectified linear unit, ReLU

$$ReLU(z) = z_+ = \begin{cases} z, & z > 0 \\ 0, & \dots \end{cases}$$

$x_2$

$1 - 1 \qquad 0$

$0 \qquad 1$

$0 \qquad 1 \qquad \to x_1$

$(x_1)$

$(x_2) \xrightarrow{W} (h_1)$
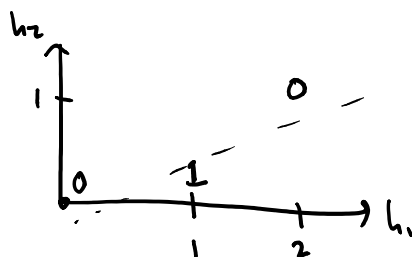
$\xrightarrow{\beta} (f)$

$(h_2)$

$(1)$

$$W^T x = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & -1 \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} = \begin{pmatrix} x_1 + x_2 \\ x_1 + x_2 - 1 \end{pmatrix}$$

$h_1 = (x_1 + x_2)_+ \quad - \text{use } ReLU$

$h_2 = (x_1 + x_2 - 1)_+ \qquad \beta = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$

$h_2$

$1 \qquad \qquad 0$

$0 \qquad 1$

$1 \qquad 2 \qquad \to h_1$

$\beta^T h \underset{\substack{\uparrow \\ \text{vectorized}}}{} = (x_1 + x_2)_+ - 2(x_1 + x_2 - 1)_+$

# Optimization for Deep Learning

Recall empirical risk: $R_n = \frac{1}{n} \sum_i l(y_i, f(x_i))$

Full gradient: $\frac{\partial R_n}{\partial w_j} = \frac{1}{n} \sum_i \frac{\partial}{\partial w_j} l(y_i, f(x_i))$

stochastic: $\frac{\partial}{\partial w_j} l(y_i, f(x_i))$
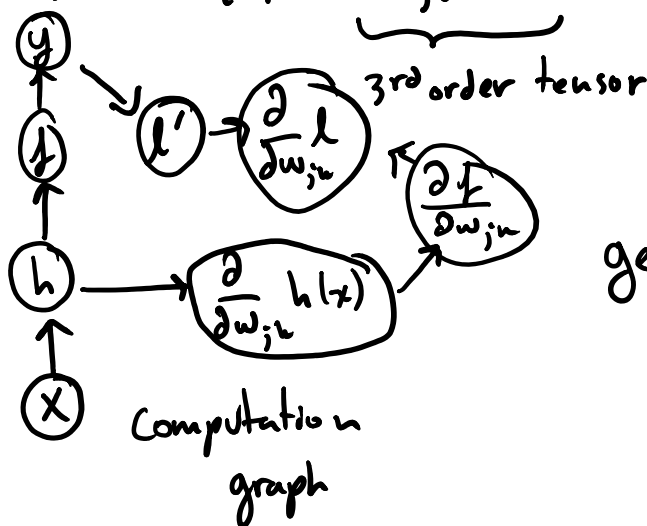
mini-batch: sample $i = 1, \dots, m$

$$\frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial w_j} l(y_i, f(x_i))$$

$$\frac{\partial}{\partial w_{jk}} l(y_i, f(x_i)) = l'(y_i, f(x_i)) \cdot \frac{\partial f}{\partial w_{jk}}(x_i) \quad \text{(chain)}$$

$$\overset{\curvearrowleft \frac{\partial}{\partial z} l(y_i, z)}{}$$

$$\frac{\partial f(x_i)}{\partial w_{jk}} = \sum_{l=1}^{H} \beta_l \underbrace{\frac{\partial}{\partial w_{jk}} h_l(x_i)}_{\text{3rd order tensor}}$$
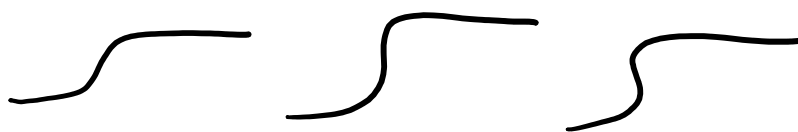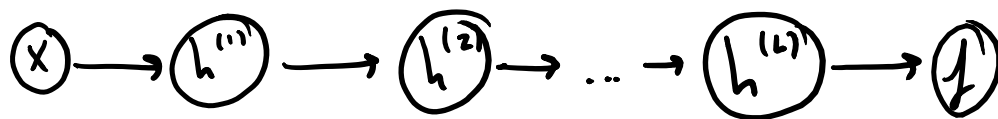
$h_l(x) = \sigma(w_l^T x)$

$\frac{\partial}{\partial w_{jk}} h_l(x) = \sigma'(w_l^T x) \cdot$
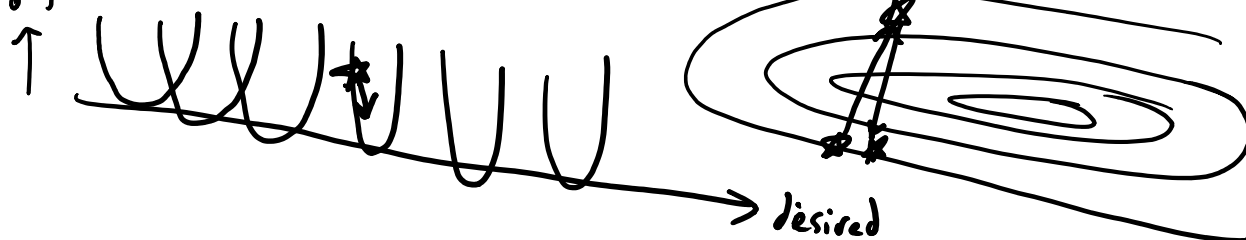$\cdot \begin{cases} x_k, & l=j \\ 0 & l \neq j \end{cases}$



Computation graph

general procedure: backprop

Deep Nnet: $H^{(l)}$ dim

$$X \rightarrow h^{(1)} \rightarrow h^{(2)} \rightarrow \ldots \rightarrow h^{(L)} \rightarrow \hat{y}$$

objective composing sigmoid $\rightarrow$ cliffs



$\rightarrow$ desired direction

acceleration "remembers" past directions in a momentum term

$$g \leftarrow \alpha g - \varepsilon \nabla_\theta \left( \frac{1}{m} \sum_{i=1}^{m} \ell(y_i, f_\theta(x_i)) \right)$$

$\uparrow$ any parameter

$$\theta \leftarrow \theta + g$$

Nesterov acceleration:

idea: gradient is at $\theta$ but applied to $\theta + \alpha g$

$$g \leftarrow \alpha g - \varepsilon \nabla_\theta \left[ \frac{1}{m} \sum_{i=1}^{m} \ell(y_i, f_{\theta + \alpha g}(x_i)) \right]$$

Adagrad / RMSprop: still have to choose $\eta_t$
(learning schedule) track gradient in $j^{th}$ direction
accumulate variance of grad. make the step

size be inversely prop. to var.

Adam: combines Nesterov accel. with RMSprop

Deep learning pipeline:

1. Specify Nnet "architecture": # and configuration of the layer, loss

2. Automatic differentiation: make gradients in closed form

3. Backprop: evaluate gradients

4. Use them in SGD solver (Adam)

★ tensor operations are ideal for GPUs