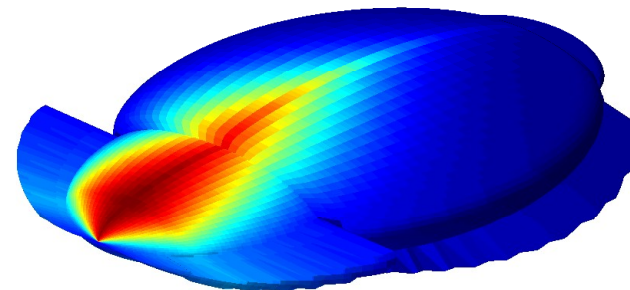
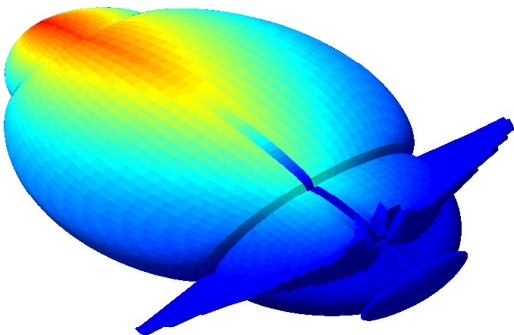


Lecture 5

Gradient Descent Learning



The Supervised Learning Problem

Input patterns $\bar{x}^k \in \mathbf{X} \subset \mathcal{R}^{n+1} \quad k = 1, \dots, M$

$$\bar{x}^k = \begin{bmatrix} x_0^k & x_1^k & x_2^k & \dots & x_n^k \end{bmatrix}^T$$

Output patterns $\bar{y}^k \in \mathbf{Y} \subset \mathcal{R}^{l+1} \quad k = 1, \dots, M$

$$\bar{y}^k = \begin{bmatrix} y_1^k & y_2^k & \dots & y_l^k \end{bmatrix}^T$$

Training Set

Fit a model:

$$\hat{\bar{y}} = f(\bar{x}; \bar{w}) \quad \bar{w} = \begin{bmatrix} w_1 & w_2 & \dots & w_N \end{bmatrix} \text{ are the parameters of the model}$$

to minimize the loss function

$$J(\bar{w}) = \sum_k J^k(\bar{w}) = \sum_k e(\hat{\bar{y}}^k, \bar{y}^k; \bar{w}) \quad e(\hat{\bar{y}}, \bar{y}; w) \text{ is an error function}$$

by adapting the parameters \bar{w}

Loss Functions

The three most commonly used loss functions are:

Absolute Error:

$$J_1(\bar{w}) = \sum_{k=1}^M \sum_{i=1}^l |y_i^k - \hat{y}_i^k| \quad e_i^k \equiv y_i^k - \hat{y}_i^k \quad J^k \equiv \sum_i |e_i^k| = \sum_i |y_i^k - \hat{y}_i^k|$$

Mean-Squared Error (MSE):

$$J_2(\bar{w}) = \sum_{k=1}^M \sum_{i=1}^l (y_i^k - \hat{y}_i^k)^2 \quad e_i^k \equiv y_i^k - \hat{y}_i^k \quad J^k \equiv \sum_i (e_i^k)^2 = \sum_i (y_i^k - \hat{y}_i^k)^2$$

Cross-Entropy Loss:

Used for binary (0/1) desired outputs and sigmoid actual outputs

$$J_{CE}(\bar{w}) = - \sum_{k=1}^M \sum_{i=1}^l (y_i^k \log \hat{y}_i^k + (1 - y_i^k) \log(1 - \hat{y}_i^k))$$

The Error Surface

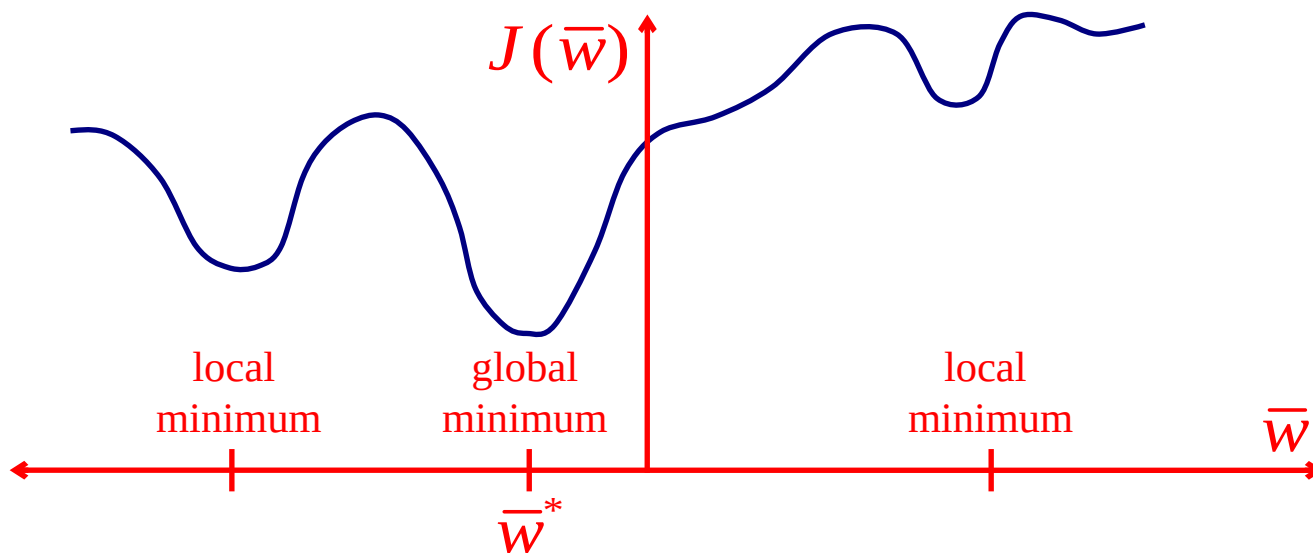
$J(\bar{w})$ defines a surface on the N -dimensional space of parameters \bar{w}

→ **error surface**

Optimization \leftrightarrow Finding the lowest point on this error surface

→ **global minimum.**

For most non-trivial problems, the error surface is very complicated and has many local minima.



Iterative Optimization

1. Present data points from the training set one at a time to the model.
2. Calculate the loss for each data point.
3. Modify the parameters:
 - a) After each step (on-line mode)
 - or
 - b) after going through all the data points (batch mode)

such that:

$$J(\bar{w}(T+1)) < J(\bar{w}(T)) \quad \text{where } T \text{ is the learning time step}$$

Training Epoch:

Usually, one pass through all the training data. In some cases, only a subset of the data may be seen in any given epoch.

On-Line Mode:

- a. Initialize $\bar{w} = \bar{w}(0)$
- b. For $t = 1$ to L_{train} (some large integer)
 - Present $\bar{x}(t) = \bar{x}^k \in X$ (thus k is the index of $\bar{x}(t)$ in X)
 - Obtain $\hat{y}(t) = f(\bar{x}(t); \bar{w})$
 - Calculate error $J(t) = J^k$
 - Calculate parameter change $\Delta\bar{w}(t)$ based on $J(t)$
 - Modify \bar{w} : $\bar{w}(t+1) = \bar{w}(t) + \Delta\bar{w}(t)$
- c. Evaluate error. If ok, end else, repeat from b.

An alternative is to use a repeat...until instead of for so that the # of learning steps depends on the performance rather than a preset value.

Batch Mode:

- a. Initialize: $\bar{w} = \bar{w}(0)$ $T = 0$ t = fast dynamics (data points)
 T = slow dynamics (epoch)
- b. $\bar{\Delta}_{acc} = \bar{0}$ (same dimension as \bar{w})
- c. For $t = 1$ to L_{epoch} (number of data points in the epoch)
- Present $\bar{x}(t) = \bar{x}^k \in X$ (thus k is the index of $\bar{x}(t)$ in X)
 - Obtain $\hat{y}(t) = f(\bar{x}(t); \bar{w})$
 - Calculate error $J(t) = J^k$
 - Calculate parameter change $\Delta \bar{w}(t)$ based on $J(t)$
 - $\bar{\Delta}_{acc} = \bar{\Delta}_{acc} + \Delta \bar{w}(t)$
- d. Modify \bar{w} : $\bar{w}' = \bar{w}(T) + \bar{\Delta}_{acc}$
- e. $T = T + 1$, $\bar{w}(T + 1) = \bar{w}'$
- f. Repeat from b. until done (or exhausted)

Gradient Descent Optimization

The main question in iterative optimization is: How to determine $\Delta \bar{w}$?

Assume that $J(\bar{w})$ is differentiable

Gradient of $J(\bar{w})$ is

$$\nabla J(\bar{w}) = \left[\frac{\partial J}{\partial w_1} \quad \frac{\partial J}{\partial w_2} \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \frac{\partial J}{\partial w_N} \right]^T$$

where N = total # of parameters = dimension of parameter space

$$\text{Optimum: } J(\bar{w}^*) \leq J(\bar{w}) \quad \forall \bar{w} \quad \Rightarrow \quad \nabla J(\bar{w}^*) = 0$$

Gradient Descent Rule: Change \bar{w} in the direction of the negative gradient.

$$\bar{w}(T+1) = \bar{w}(T) + \Delta \bar{w}(T) = \bar{w}(T) - \eta \nabla J(\bar{w}(T)) \quad \eta = \text{learning rate}$$

$$\Delta w_i(t) = -\eta \frac{\partial J}{\partial w_i}$$

To see why this works, expand $J(\bar{w}(T+1))$ in a Taylor series.

$$J(\bar{w}(T+1)) \approx J(\bar{w}(T)) + [\nabla J(\bar{w}(T))]^T \Delta \bar{w}(T) + \text{higher order terms}.$$

To lower LHS. $(\nabla J)^T \Delta \bar{w} < 0$

Since $(\nabla J)^T \Delta \bar{w} = \|\nabla J\| \|\Delta \bar{w}\| \cos \theta$ θ is the angle between ∇J and $\Delta \bar{w}$

Loss is reduced most when $\cos = -1 \Rightarrow \Delta \bar{w} \propto -\nabla J$

Intuitive Explanation:

- If you are trying to get to the bottom of a valley, the fastest way down is to go down the steepest slope at every step.
- The gradient vector points up the steepest slope at every point.
- Thus, going in the opposite direction of the gradient is the steepest way down.

Stochastic Gradient Descent (Online)

Note that, ideally, gradient descent training of a network requires:

$$\bar{w}(T+1) = \bar{w}(T) - \eta \nabla J(\bar{w}(T))$$

where $\nabla J(\bar{w}(T)) = \nabla \sum_{k=1}^m J^k(\bar{w}(T))$ i.e., the gradient must be calculated over **all** data summed together, using weights from epoch T .

This can be problematic for very large data-sets

Solution: Train in on-line mode, selecting data points **randomly** from the data-set at each step, and updating immediately.

- a. Initialize $\bar{w} = \bar{w}(0)$
- b. For $t = 1$ to M (some large integer)
 - Choose $\bar{x}(t) = x^k \in X$ randomly using some policy
 - Obtain $\hat{y}(t) = f(\bar{w}^T \bar{x}(t))$
 - Calculate error $e(t)$ and local gradient $\nabla J^k(\bar{w}(t))$
 - Modify \bar{w} : $\bar{w}(t+1) = \bar{w}(t) - \eta \nabla J^k(\bar{w}(t))$
- c. Evaluate error. If ok, end else, repeat from b.

Stochastic gradient descent is used widely in applications of supervised learning to **big data**, including **deep learning**.

Stochastic Gradient Descent (Minibatch)

- Divide up each epoch into several minibatches of data (random or regular).
- Average gradient over each minibatch and adjust weights.

a. Initialize $\bar{w} = \bar{w}(0)$

b. For $T = 1$ to N_{epochs}

Set minibatch gradient $\nabla J_{mb} = 0$

For $t = 1$ to q

- Choose $\bar{x}(t) = x^k \in X$ randomly using some policy
- Obtain $\hat{y}(t) = f(\bar{w}^T \bar{x}(t))$
- Calculate error $e(t)$ and local gradient $\nabla J^k(\bar{w}(t))$
- Accumulate minibatch gradient $\nabla J_{mb} = \nabla J_{mb} + \nabla J^k(\bar{w}(t))$

Modify \bar{w} : $\bar{w}(t+1) = \bar{w}(t) - \eta \frac{1}{q} \nabla J_{mb}$

Obviously, this can also be run in *repeat until* mode instead of a *for* loop with a pre-specified number of epochs.

Active Management of Training Data

- Use examples on which error is greater.
- Use examples that are very different from those already used (to cover new regions of sample space)
- Randomize the order of presentation.
- Weight probability of presentation for different classes.
- Watch out for outliers that can really skew the weights catastrophically.
- Use multiple, overlapping training and testing sets.

Gradient Descent with Momentum

Original gradient descent update:

$$\Delta w_{ij}(t) = -\eta \frac{\partial J}{\partial w_{ij}}$$

Gradient descent with momentum:

$$\Delta w_{ij}(t) = -\eta \frac{\partial J}{\partial w_{ij}} + \alpha \Delta w_{ij}(t-1)$$

$$0 \leq |\alpha| < 1$$

If $|\alpha| \geq 1$, the equation can become unstable.

What does momentum do?

if Δw_i has the same sign for several consecutive iterations

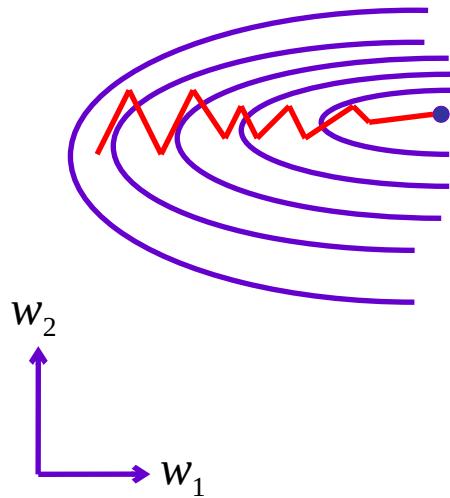
\Rightarrow stable downhill direction for w_i

\rightarrow momentum increases descent rate

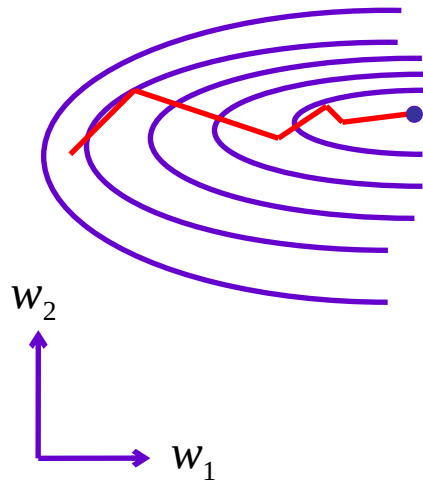
if Δw_i has alternating signs over consecutive iterations

\Rightarrow w_i is not contributing much to descent

\rightarrow momentum damps movement in w_i



Without momentum, Δw_2 is not contributing to progress. This is indicated by its oscillating sign.



With momentum, Δw_2 is damped and Δw_1 increased to find minimum faster

Momentum magnifies progress in directions of sustained descent and diminishes progress in other directions

The Adam (Adaptive Moment) Algorithm

Kingma & Ba (2015) ADAM: A Method for Stochastic Optimization, Proc. ICLR 2015

Parameters: Learning rate η , Biases β_1, β_2

Variables: 1st moment $\bar{m}(t), \bar{M}(t)$, 2nd moment $\bar{v}(t), \bar{V}(t)$

Weights: $\bar{w}(t) = [w_{ij}(t)]$

Initialize: $\bar{w}(0), \bar{m}(0) = 0, \bar{v}(0) = 0, t = 0$

while \bar{w} not converged **do**

$t = t + 1$

$\bar{g}(t) = \nabla J(\bar{w}(t - 1))$ current loss gradient

$\bar{m}(t) = \beta_1 \bar{m}(t - 1) + (1 - \beta_1) \bar{g}(t)$ update biased 1st moment

$\bar{v}(t) = \beta_2 \bar{v}(t - 1) + (1 - \beta_2) (\bar{g}(t) \odot \bar{g}(t))$ update biased 2nd moment

$\bar{M}(t) = \bar{m}(t) / (1 - \beta_1^t)$ compute unbiased 1st moment

$\bar{V}(t) = \bar{v}(t) / (1 - \beta_2^t)$ compute unbiased 2nd moment

$\Delta \bar{w}(t) = -\eta \bar{M}(t) / (\sqrt{\bar{V}(t)} + \varepsilon)$ update weights

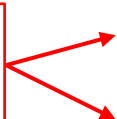
end while

return $\bar{w}(t)$

Elementwise
product



raised to
power t

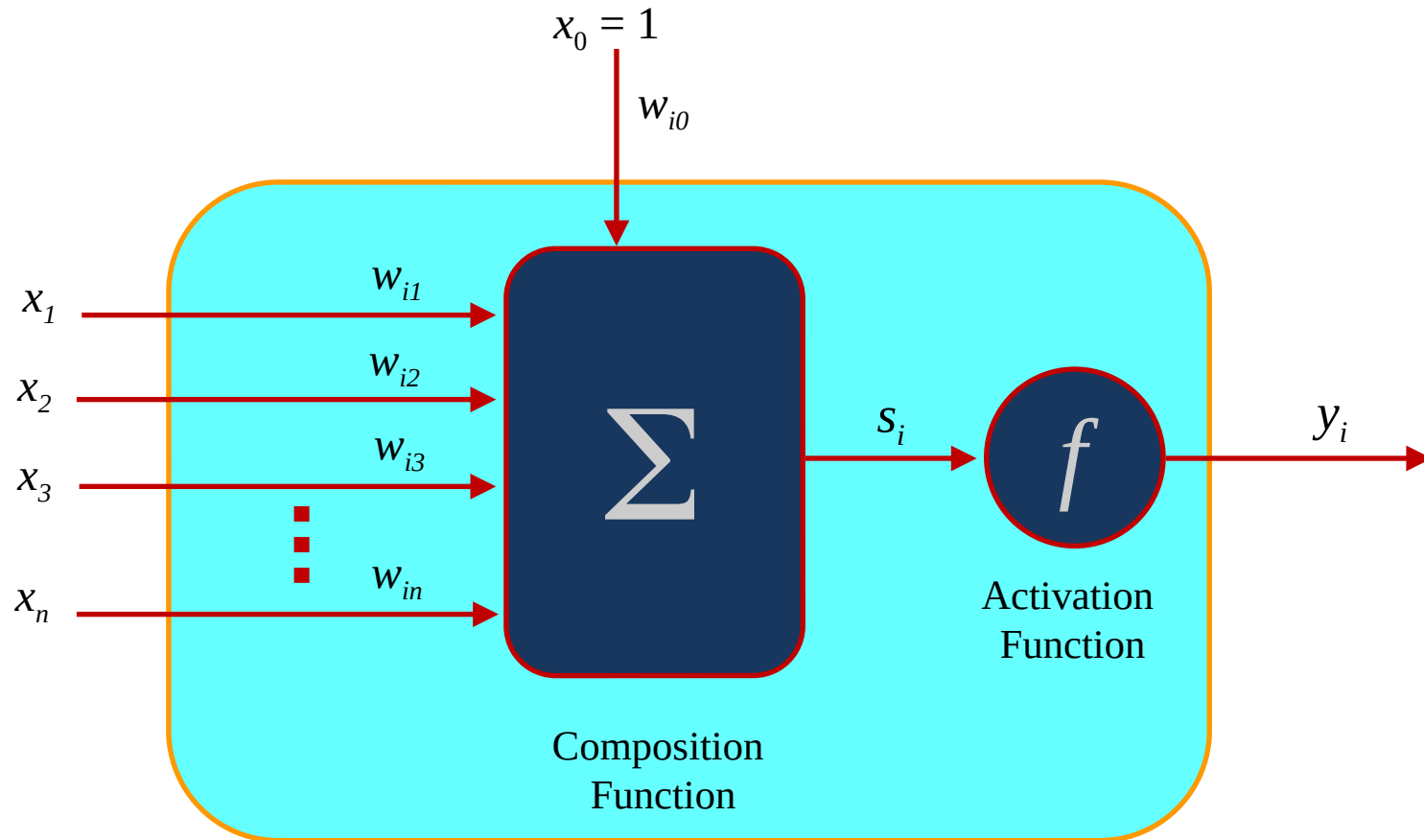


Elementwise
ratio



Recommended choices: $\beta_1 = 0.9, \beta_2 = 0.999, \varepsilon = 10^{-8}$

Simple Neuron Model



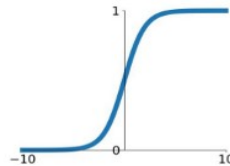
$$s_i = \sum_j w_{ij} x_j$$

$$y_i = f(s_i)$$

Activation Functions

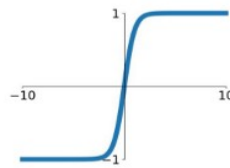
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



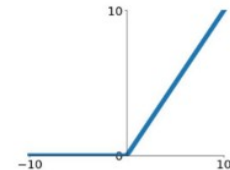
tanh

$$\tanh(x)$$



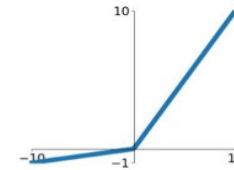
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

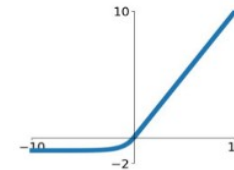


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Source: <https://www.quora.com/>

Why use ReLU?

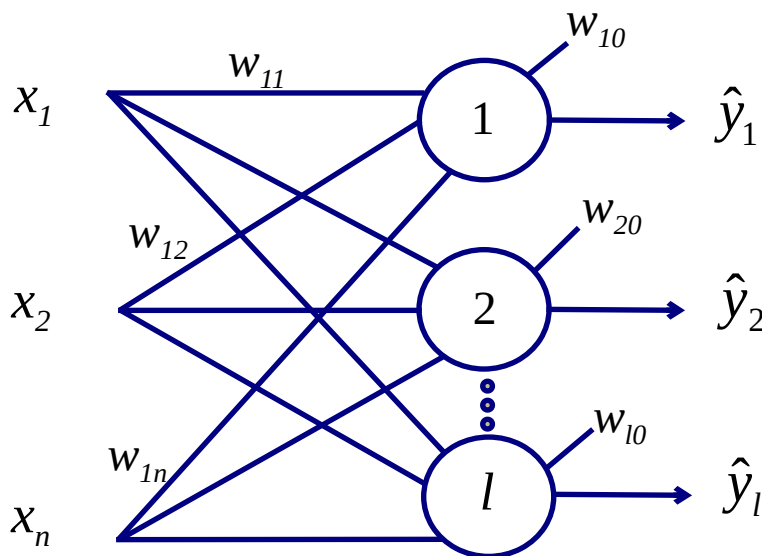
- It is nonlinear and monotonic (both useful/necessary)
- Its gradient does not vanish on the positive side (unlike sigmoids).
- It gives infinite range of output values (0 to $+\infty$)

But

- ReLU gradient vanishes on negative side \rightarrow Leaky ReLU
- ReLU is non-differentiable at 0 \rightarrow ELU

The Least Mean Square (LMS) Algorithm

Consider a single layer of neurons



For data point (\bar{x}^q, \bar{y}^q)

Define:

$$J^q = \frac{1}{2} \sum_{i=1}^l (e_i^q)^2 = \frac{1}{2} \sum_{i=1}^l (y_i^q - \hat{y}_i^q)^2$$

$$J = \sum_{q=1}^M J^q$$

$$J = \frac{1}{2} \sum_{q=1}^M \sum_{i=1}^l \left[y_i^q - f \left(\sum_{j=1}^n w_{ij} x_j^q \right) \right]^2$$

Input patterns $\bar{x}^q \in X \subset \mathbb{R}^{n+1}$ $q = 1, \dots, M$

$$\bar{x}^q = \begin{bmatrix} 1 & x_1^q & x_2^q & \dots & x_n^q \end{bmatrix}^T$$

bias

Output patterns $\bar{y}^q \in Y \subset \mathbb{R}^l$

Note that y_i are not necessarily 0,1

Weight matrix $\bar{w} = [\bar{w}_1 \ \bar{w}_2 \ \dots \ \bar{w}_l]^T$

$$\bar{w}_i = \begin{bmatrix} w_{i0} & w_{i1} & w_{i2} & \dots & w_{in} \end{bmatrix}^T$$

bias weight

$$\bar{w}_i \in \mathbb{R}^{n+1}$$

$$\bar{w} = \begin{bmatrix} w_{10} & w_{11} & w_{12} & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & w_{1n} \\ w_{20} & w_{21} & w_{22} & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & w_{2n} \\ \cdot & & & & & & & & & \\ \cdot & & & & & & & & & \\ \cdot & & & & & & & & & \\ \cdot & & & & & & & & & \\ w_{l0} & w_{l1} & w_{l2} & & & & & & & w_{ln} \end{bmatrix} = l \times (n+1) \equiv N$$

$\bar{x}(t) \equiv$ pattern presented at time step t

$$x(t) \in X$$

$\bar{y}(t) \equiv$ desired response at step t

$$\bar{y}(t) \in Y$$

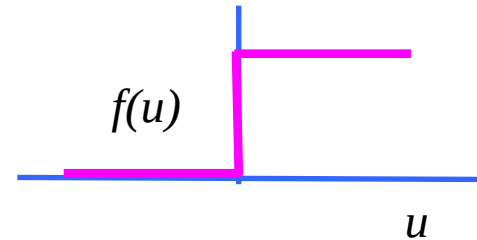
$\bar{w}(t) \equiv$ weight matrix at step t

$\hat{\bar{y}}(t) \equiv$ actual output at step t

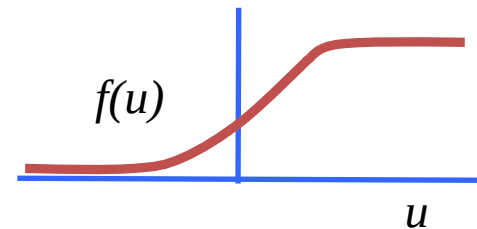
$$\hat{\bar{y}}(t) = [\hat{y}_1(t) \quad \hat{y}_2(t) \quad \dots \quad \hat{y}_l(t)]^T$$

$$\hat{y}_i(t) = f(\bar{w}_i^T(t) \bar{x}(t))$$

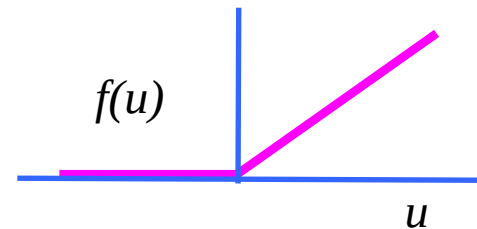
Typical $f(u)$



Threshold



Sigmoid



ReLU

if \bar{x}^q is presented to the network
 $\Delta \bar{w} = -\eta \nabla J^q$

$$\nabla J^q = \begin{bmatrix} \frac{\partial J^q}{\partial w_{10}} \\ \cdot \\ \cdot \\ \cdot \\ \frac{\partial J^q}{\partial w_{l_n}} \end{bmatrix}$$

so

$$\Delta w_{ij} = -\eta \frac{\partial J^q}{\partial w_{ij}}$$

$$\frac{\partial J^q}{\partial w_{ij}} = \frac{\partial \frac{1}{2} \sum_p (e_p^q)^2}{\partial w_{ij}}$$

$$= \frac{1}{2} \sum_p \frac{\partial (e_p^q)^2}{\partial w_{ij}}$$

$$= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial w_{ij}}$$

p = index of units.
 $p = 1, \dots, l$

superposition
of derivatives

since w_{ij} affects only e_i^k

$$\begin{aligned}
 \frac{\partial J^k}{\partial w_{ij}} &= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial w_{ij}} \\
 &= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial e_i^q} \frac{\partial e_i^q}{\partial w_{ij}} \\
 &= e_i^q \frac{\partial e_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial w_{ij}} \\
 &= e_i^q \cdot \frac{\partial e_i^q}{\partial \hat{y}_i^q} \cdot \frac{\partial \hat{y}_i^q}{\partial s_i^q} \cdot \frac{\partial s_i^q}{\partial w_{ij}} \\
 &= e_i^q \cdot (-1) \cdot f'(s_i^q) \cdot x_j^q
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial J^q}{\partial w_{ij}} &= -e_i^q f'(s_i^q) x_j^q \\
 &= -(y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q
 \end{aligned}$$

$$\Delta w_{ij} = \eta (y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q$$

$$\frac{\partial J^k}{\partial w_{ij}} = \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial w_{ij}}$$

$$= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial e_i^q} \frac{\partial e_i^q}{\partial w_{ij}}$$

$$= e_i^q \frac{\partial e_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial w_{ij}}$$

$$= e_i^q \cdot \frac{\partial e_i^q}{\partial \hat{y}_i^q} \cdot \frac{\partial \hat{y}_i^q}{\partial s_i^q} \cdot \frac{\partial s_i^q}{\partial w_{ij}}$$

$$= e_i^q \cdot (-1) \cdot f'(s_i^q) \cdot x_j^q$$

$$\frac{\partial J^q}{\partial w_{ij}} = - e_i^q f'(s_i^q) x_j^q$$

$$= - (y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q$$

$$\Delta w_{ij} = \eta (y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q$$

Note that the error terms are dissociated from others, so we can use any differentiable error metric without changing the algorithm.

$$\begin{aligned}
 \frac{\partial J^k}{\partial w_{ij}} &= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial w_{ij}} \\
 &= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial e_i^q} \frac{\partial e_i^q}{\partial w_{ij}} \\
 &= e_i^q \frac{\partial e_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial w_{ij}} \\
 &= e_i^q \cdot \frac{\partial e_i^q}{\partial \hat{y}_i^q} \cdot \frac{\partial \hat{y}_i^q}{\partial s_i^q} \cdot \frac{\partial s_i^q}{\partial w_{ij}} \\
 &= e_i^q \cdot (-1) \cdot \boxed{f'(s_i^q)} \cdot x_j^q
 \end{aligned}$$

And any activation function – as long as it is differentiable

$$\begin{aligned}
 \frac{\partial J^q}{\partial w_{ij}} &= -e_i^q f'(s_i^q) x_j^q \\
 &= -(y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q
 \end{aligned}$$

$$\Delta w_{ij} = \eta (y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q$$

$$\begin{aligned}\frac{\partial J^k}{\partial w_{ij}} &= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial w_{ij}} \\&= \frac{1}{2} \frac{\partial (e_i^q)^2}{\partial e_i^q} \frac{\partial e_i^q}{\partial w_{ij}} \\&= e_i^q \frac{\partial e_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial w_{ij}} \\&= e_i^q \cdot \frac{\partial e_i^q}{\partial \hat{y}_i^q} \cdot \frac{\partial \hat{y}_i^q}{\partial s_i^q} \cdot \frac{\partial s_i^q}{\partial w_{ij}} \\&= e_i^q \cdot (-1) \cdot f'(s_i^q) \cdot x_j^q\end{aligned}$$

$$\begin{aligned}\frac{\partial J^q}{\partial w_{ij}} &= -e_i^q f'(s_i^q) x_j^q \\&= -(y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q\end{aligned}$$

$$\Delta w_{ij} = \eta (y_i^q - \hat{y}_i^q) f'(s_i^q) x_j^q$$

$-\nabla J(\bar{w})$

Sigmoid $f(u)$

$$f(u) = \frac{1}{1 + e^{-u}}$$

$$\frac{df}{du} = \frac{e^{-u}}{(1 + e^{-u})^2} = \frac{e^{-u}}{1 + e^{-u}} \cdot \frac{1}{1 + e^{-u}}$$

$$= (1 - f(u)) f(u)$$

$$\frac{\partial J^q}{\partial w_{ij}} = - (y_i^q - \hat{y}_i^q) \hat{y}_i^q (1 - \hat{y}_i^q) x_j^q$$

$$\text{defining } \delta_i^q = \hat{y}_i^q (1 - \hat{y}_i^q) (y_i^q - \hat{y}_i^q)$$

$$\Delta w_{ij} = \eta \delta_i^q x_j^q$$

Linear $f(u)$

$$\hat{y}_i^q = s_i^q = \sum_{j=0}^n w_{ij} x_j^q$$

$$f'(s_i^q) = 1$$

$$\frac{\partial J^q}{\partial w_{ij}} = - (y_i^q - \hat{y}_i^q) x_j^q$$

\Rightarrow LMS for linear neurons = perceptron learning

$$\Delta w_{ij} = \eta (y_i^q - \hat{y}_i^q) x_j^q$$

which is the perceptron learning rule.