



Deep Learning



BITS Pilani
Pilani Campus

Dr. Sugata Ghosal
sugata.ghosal@pilani.bits-pilani.ac.in



BITS Pilani
Pilani Campus

Worked Out Problems Optimization

These slides are assembled by the instructor with grateful acknowledgement of the many others who made their course materials freely available online.

Optimal Learning Rate: Multivariate Diagonal Quadratic Error Function

Error surface is given by $E(x,y,z) = 3x^2 + 2y^2 + 4z^2 + 6$. What is the optimal learning rate that leads to fastest convergence to the global minimum?

$$E = \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{w}^T \mathbf{b} + c = \frac{1}{2} \sum_i (a_{ii} w_i^2 + b_i w_i) + c$$

$$\eta_{x,\text{opt}} = 1/6$$

$$\eta_{y,\text{opt}} = 1/4$$

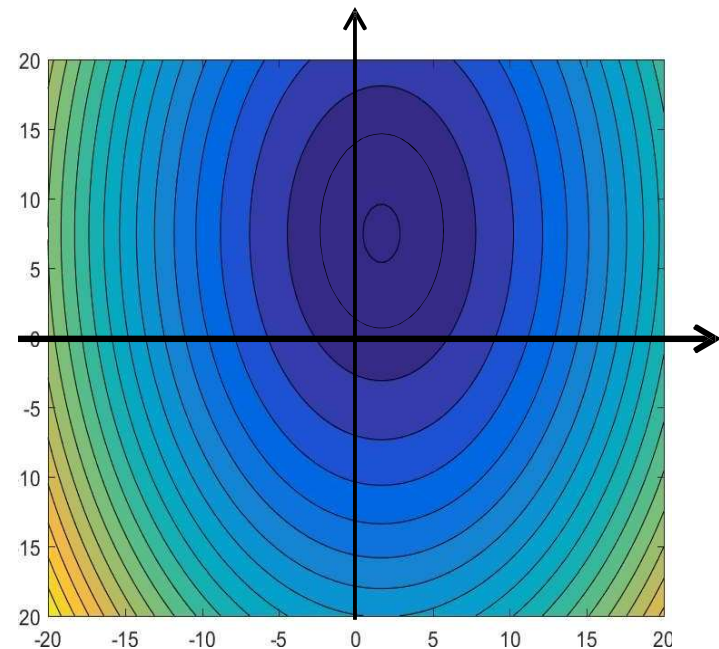
$$\eta_{z,\text{opt}} = 1/8$$

Optimal learning rate =

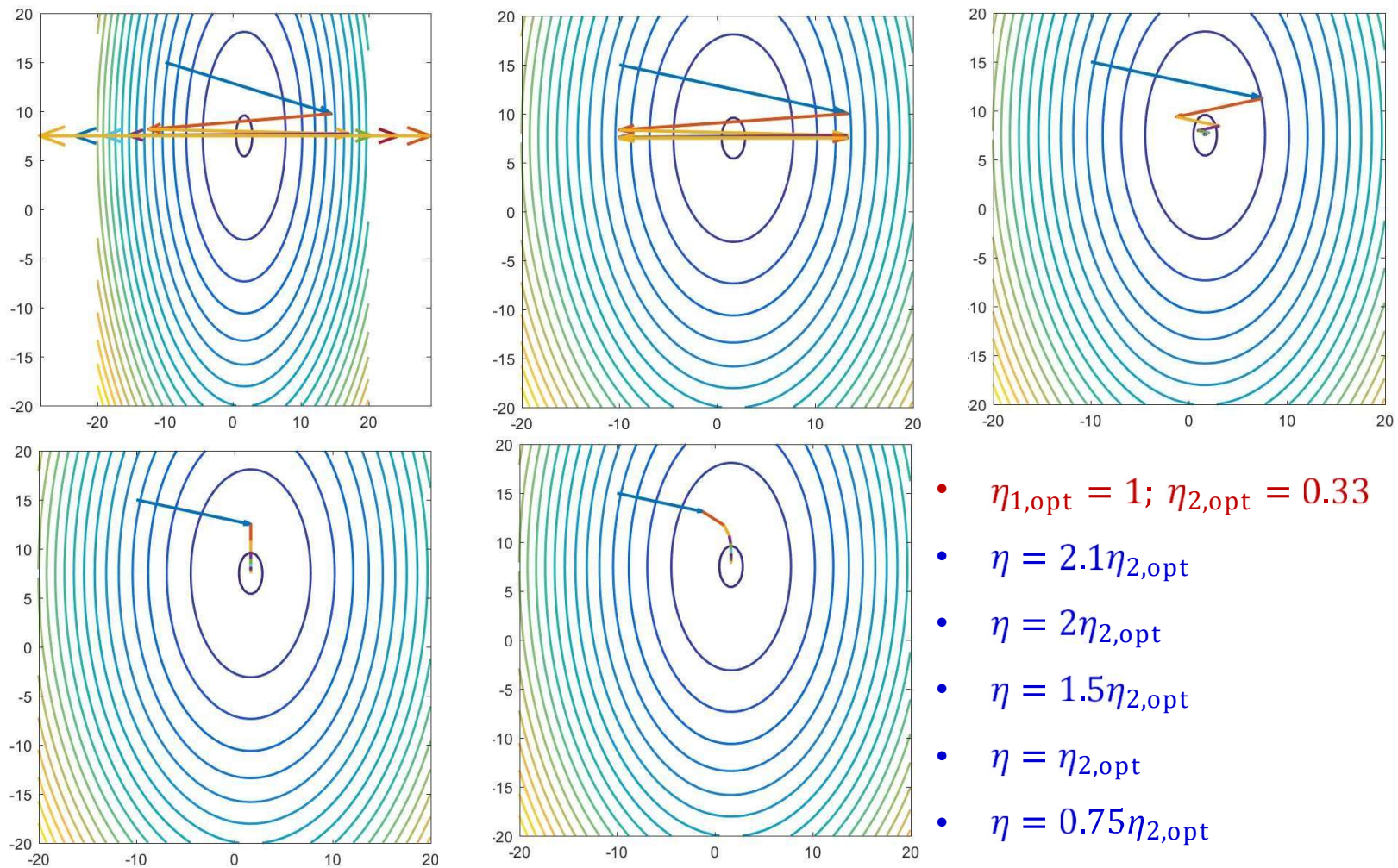
$$\min(\eta_{x,\text{opt}}, \eta_{y,\text{opt}}, \eta_{z,\text{opt}}) = 0.125$$

Largest learning rate for convergence =

$$\min(2\eta_{x,\text{opt}}, 2\eta_{y,\text{opt}}, 2\eta_{z,\text{opt}}) = 0.333$$



Dependence on learning rate – Error Minimization



Minimization of Quadratic Error Function

Consider an error function $E(w_1, w_2) = 0.05 + \frac{(w_1-3)^2}{4} + \frac{(w_2-4)^2}{9} - \frac{(w_1-3)(w_2-4)}{6}$. Different variants of gradient descent algorithm can be used to minimize this error function w.r.t. w_1, w_2 . Assume $(w_1, w_2) = (1, 1)$ at time $(t-1)$ and after update $(w_1, w_2) = (1.5, 2.0)$ at time t . Assume, learning rate $\eta = 0.3$ and momentum update rate $\beta = 0.9$.

What is the value of (w_1, w_2) that minimizes this error function? What is the minimum possible value of E ?

$w_1=3, w_2=4. E_{min} = 0.05$

Weight Updates – Ordinary Gradient Descent

Consider an error function $E(w_1, w_2) = 0.05 + \frac{(w_1-3)^2}{4} + \frac{(w_2-4)^2}{9} - \frac{(w_1-3)(w_2-4)}{6}$. Different variants of gradient descent algorithm can be used to minimize this error function w.r.t. w_1, w_2 . Assume $(w_1, w_2) = (1, 1)$ at time $(t-1)$ and after update $(w_1, w_2) = (1.5, 2.0)$ at time t . Assume, learning rate $\eta = 0.3$ and momentum update rate $\beta = 0.9$.

What is the value of (w_1, w_2) at time $(t+1)$ if standard gradient descent is used?

Weight Updates – Ordinary Gradient Descent

Consider an error function $E(w_1, w_2) = 0.05 + \frac{(w_1-3)^2}{4} + \frac{(w_2-4)^2}{9} - \frac{(w_1-3)(w_2-4)}{6}$. Different variants of gradient descent algorithm can be used to minimize this error function w.r.t. w_1, w_2 . Assume $(w_1, w_2) = (1, 1)$ at time $(t-1)$ and after update $(w_1, w_2) = (1.5, 2.0)$ at time t . Assume, learning rate $\eta = 0.3$ and momentum update rate $\beta = 0.9$.

What is the value of (w_1, w_2) at time $(t+1)$ if standard gradient descent is used?

$$\delta E / \delta w_1 = 0.5 \cdot (w_1 - 3) - (w_2 - 4) / 6 \quad \text{and} \quad \delta E / \delta w_2 = 2/9 \cdot (w_2 - 4) - (w_1 - 3) / 6$$

$$\text{So, } w_1(t+1) = 1.5 - 0.3 \cdot 0.5(1.5 - 3) + 0.3 \cdot (2 - 4) / 6 = 1.625, \text{ and}$$

$$w_2(t+1) = 2.0 - 0.3 \cdot 2 \cdot (2 - 4) / 9 + 0.3 \cdot (1.5 - 3) / 6 = 2.058$$

Nestorov

$$w_{1_int} = 1.5 + 0.9 \cdot 0.5 = 1.95 \quad w_{2_int} = 2.0 + 0.9 \cdot 1 = 2.9$$

$$dE/dw_1(w_{1_int}, w_{2_int}) = 0.5 \cdot (1.95 - 3) - (2.9 - 4) / 6 = -0.342$$

$$dE/dw_2(w_{1_int}, w_{2_int}) = 2/9 \cdot (2.9 - 4) - (1.95 - 3) / 6 = -0.0694$$

$$w_1(t+1) = 1.95 + 0.3 \cdot 0.342 = 2.0526 \quad w_2(t+1) = 2.9 + 0.3 \cdot 0.0694 = 2.92$$

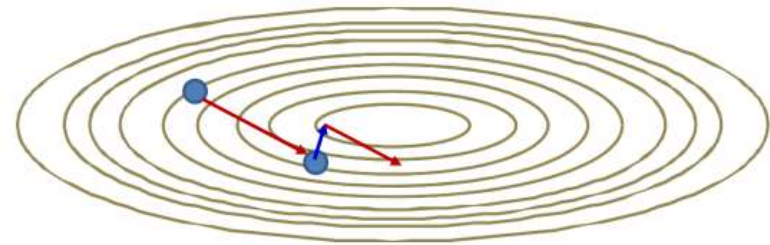
Weight Updates – Momentum Method

Consider an error function $E(w_1, w_2) = 0.05 + \frac{(w_1-3)^2}{4} + \frac{(w_2-4)^2}{9} - \frac{(w_1-3)(w_2-4)}{6}$. Different variants of gradient descent algorithm can be used to minimize this error function w.r.t. w_1, w_2 . Assume $(w_1, w_2) = (1, 1)$ at time $(t-1)$ and after update $(w_1, w_2) = (1.5, 2.0)$ at time t . Assume, learning rate $\eta = 0.3$ and momentum update rate $\beta = 0.9$.

What is the value of (w_1, w_2) at time $(t+1)$ if momentum based gradient descent is used?

$$w_1(t+1) = 1.625 + (1.5 - 1.0) * 0.9 = 2.075$$

$$w_2(t+1) = 2.058 + 0.9 * (2.0 - 1.0) = 2.958$$



The momentum method

$$\Delta W^{(k)} = \beta \Delta W^{(k-1)} - \eta \nabla_W \text{Loss}(W^{(k-1)})^T$$

Nestorov

$$w_{1_int} = 1.5 + 0.9 \times 0.5 = 1.95 \quad w_{2_int} = 2.0 + 0.9 \times 1 = 2.9$$

$$dE/dw_1(w_{1_int}, w_{2_int}) = 0.5 \times (1.95 - 3) - (2.9 - 4)/6 = -0.342$$

$$dE/dw_2(w_{1_int}, w_{2_int}) = 2/9 \times (2.9 - 4) - (1.95 - 3)/6 = -0.0694$$

$$w_1(t+1) = 1.95 + 0.3 \times 0.342 = 2.0526 \quad w_2(t+1) = 2.9 + 0.3 \times 0.0694 = 2.92$$

Weight Updates – RProp

Consider an error function $E(w_1, w_2) = 0.05 + \frac{(w_1-3)^2}{4} + \frac{(w_2-4)^2}{9} - \frac{(w_1-3)(w_2-4)}{6}$. Different variants of gradient descent algorithm can be used to minimize this error function w.r.t. w_1, w_2 . Assume $(w_1, w_2) = (1, 1)$ at time $(t-1)$ and after update $(w_1, w_2) = (1.5, 2.0)$ at time t . Assume, $\alpha = 1.5, \beta = 0.6$

What will be (w_1, w_2) at $(t+1)$?

At time $t-1$,

$$dE/dw_1 = 0.5 \cdot (1-3) - (1-4)/6 = -0.5$$

$$dE/dw_2 = 2/9 \cdot (1-4) - (1-3)/6 = -0.333$$

At time t ,

$$dE/dw_1 = 0.5 \cdot (1.5-3) - (2.0-4)/6 = 0.4167$$

$$dE/dw_2 = 2/9 \cdot (2-4) - (1.5-3)/6 = -0.194$$

$$\Delta w_1 = 1.5 - 1 = 0.5$$

$$\Delta w_2 = 2 - 1 = 1$$

$$w_1(t+1) = 1 + 0.5 \cdot 0.6 = 1.3, \text{ sign of derivation became different}$$

$$w_2(t+1) = 2.0 + 1.5 \cdot 1 = 3.5, \text{ sign of derivative remained same}$$



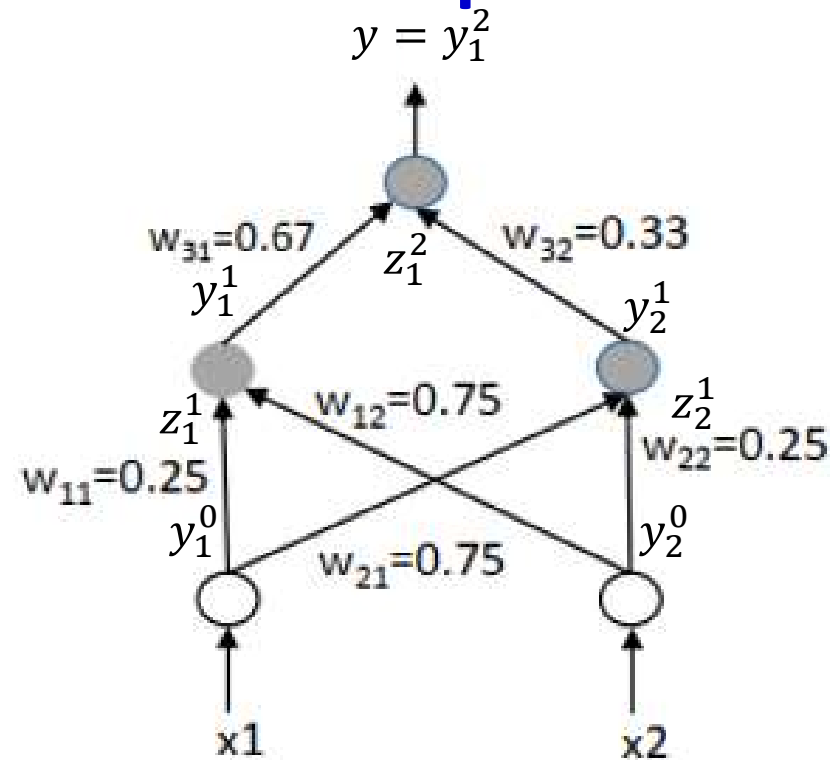
BITS Pilani
Pilani Campus



Solved Examples

These slides are assembled by the instructor with grateful acknowledgement of the many others who made their course materials freely available online.

Weight Updates in Deep Networks

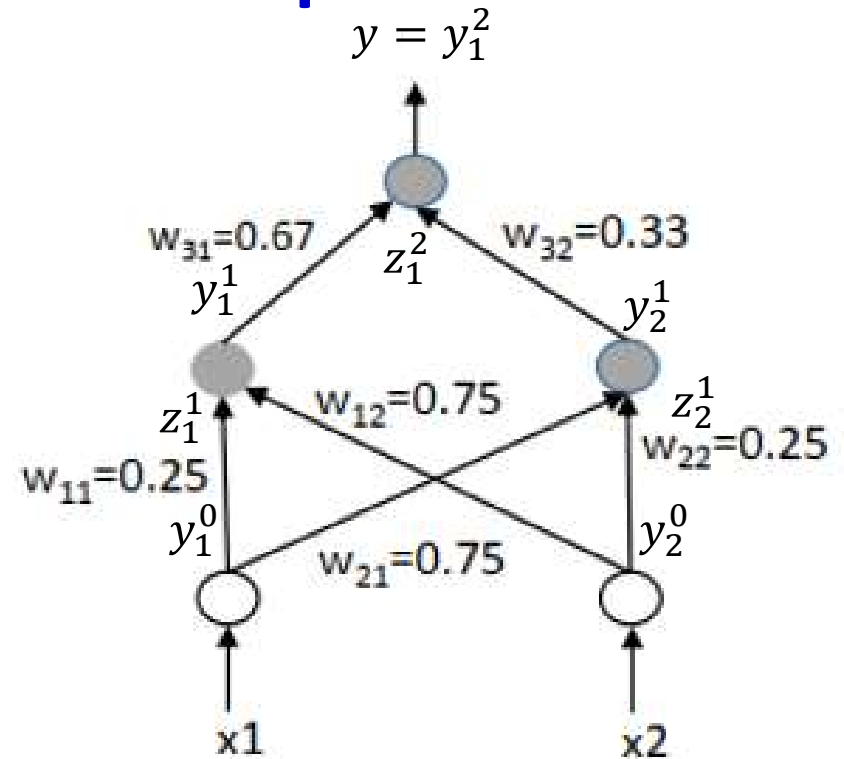


- Training Input: $(x_1, x_2) = (1, 1)$ Target Output: $d = 0.0$
- *div* function: square error
- Activation function $f(\cdot)$: ReLU
- Bias: 0 at all nodes
- $\eta = 0.1$
- What are the values of w_{31} and w_{12} in next iteration?

Weight Updates with ReLU and square error

- Training Input: $(x_1, x_2) = (1, 1)$
- Target Output: $d = 0.0$

$$w_{31}(t+1) = ?$$



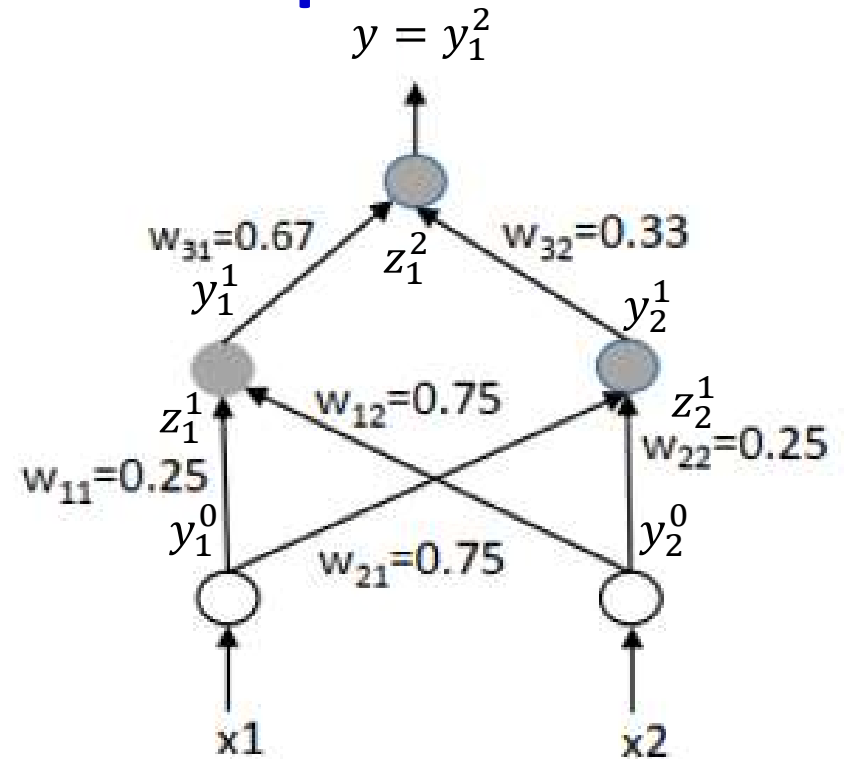
$$w_{31} = w_{31} - \eta * \frac{\delta div}{\delta w_{31}} = 0.67 - 0.1 = 0.57$$

Weight Updates with ReLU and square error

... Training Input: $(x_1, x_2) = (1, 1)$

- Target Output: $d = 0.0$

$$w_{12}(t+1) = ?$$



$$w_{12} = w_{12} - \eta * \frac{\delta div}{\delta w_{12}} = 0.75 - .067 = 0.683$$

Weight Updates with ReLU. square error

$$z_1^1 = 0.25 * 1 + 0.75 * 1 = 1.0$$

$$z_2^1 = 0.75 * 1 + 0.25 * 1 = 1.0$$

$$y_1^1 = \text{ReLU}(z_1^1) = 1 \quad y_2^1 = \text{ReLU}(z_2^1) = 1$$

$$z_1^2 = 0.67 * 1 + 0.33 * 1 = 1.0$$

$$y_1^2 = \text{ReLU}(z_1^2) = 1.0$$

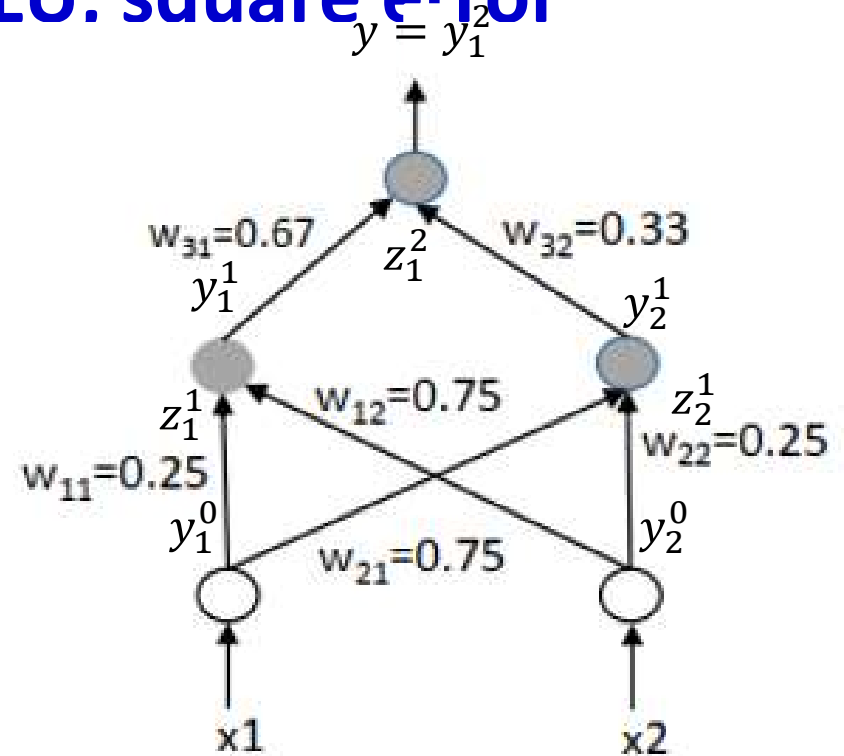
$$\text{div} = \frac{1}{2} (d - y)^2$$

$$\frac{\delta \text{div}}{\delta y} = (y - d) = 1.0$$

$$\frac{\delta \text{div}}{\delta z_1^2} = \frac{\delta \text{div}}{\delta y_1^2} * \frac{\delta y_1^2}{\delta z_1^2} = 1 * 1 = 1$$

$$\frac{\delta \text{div}}{\delta w_{31}} = \frac{\delta \text{div}}{\delta z_1^2} * \frac{\delta z_1^2}{\delta w_{31}} = 1 * y_1^1 = 1$$

$$w_{31} = w_{31} - \eta * \frac{\delta \text{div}}{\delta w_{31}} = 0.67 - 0.1 = 0.57$$



$$\frac{\delta \text{div}}{\delta y_1^1} = \frac{\delta \text{div}}{\delta z_1^2} * \frac{\delta z_1^2}{\delta y_1^1} = 1 * w_{31}$$

$$\frac{\delta \text{div}}{\delta z_1^1} = \frac{\delta \text{div}}{\delta y_1^1} * \frac{\delta y_1^1}{\delta z_1^1} = 0.67 * 1 = 0.67$$

$$\frac{\delta \text{div}}{\delta w_{12}} = \frac{\delta \text{div}}{\delta z_1^1} * \frac{\delta z_1^1}{\delta w_{12}} = 0.67 * y_2^0 = 0.67$$

$$w_{12} = w_{12} - \eta * \frac{\delta \text{div}}{\delta w_{12}} = 0.75 - 0.067 = 0.683$$

Weight Updates with ReLU Hidden, sigmoid output and binary cross-

$$z_1^1 = 0.25 * 1 + 0.75 * 1 = 1.0$$

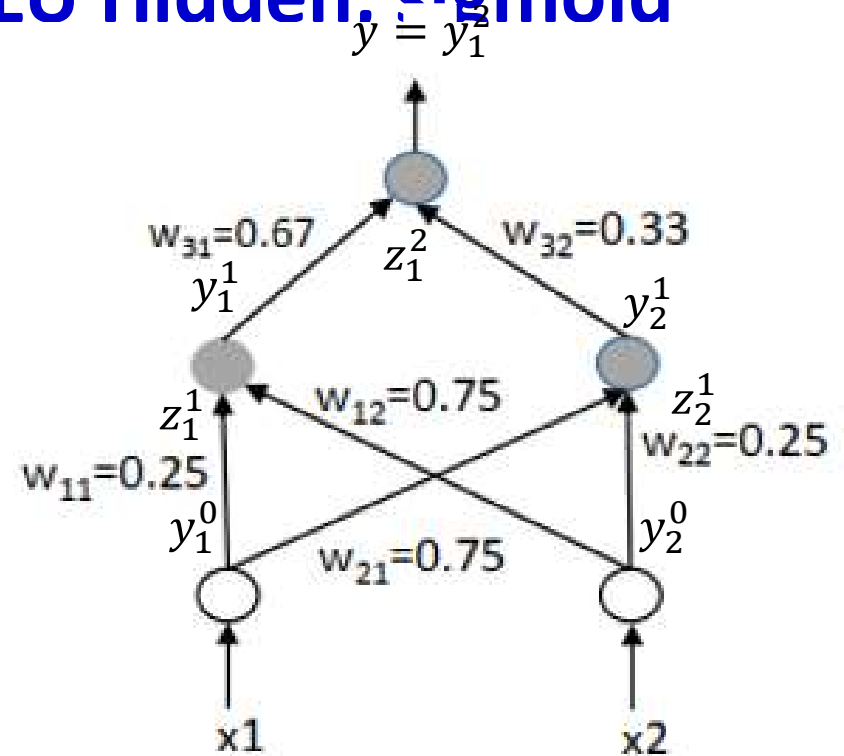
$$z_2^1 = 0.75 * 1 + 0.25 * 1 = 1.0$$

$$y_1^1 = \text{ReLU}(z_1^1) = 1 \quad y_2^1 = \text{ReLU}(z_2^1) = 1$$

$$z_1^2 = 0.67 * 1 + 0.33 * 1 = 1.0$$

$$y_1^2 = \text{sigmoid}(z_1^2) = e / (1 + e)$$

$$\text{div} = -d * \log(y) - (1 - d) * \log(1 - y)$$



Weight Updates with ReLU Hidden, sigmoid output and binary cross-

$$z_1^1 = 0.25 * 1 + 0.75 * 1 = 1.0$$

$$z_2^1 = 0.75 * 1 + 0.25 * 1 = 1.0$$

$$y_1^1 = \text{ReLU}(z_1^1) = 1 \quad y_2^1 = \text{ReLU}(z_2^1) = 1$$

$$z_1^2 = 0.67 * 1 + 0.33 * 1 = 1.0$$

$$y_1^2 = \text{sigmoid}(z_1^2) = e / (1 + e)$$

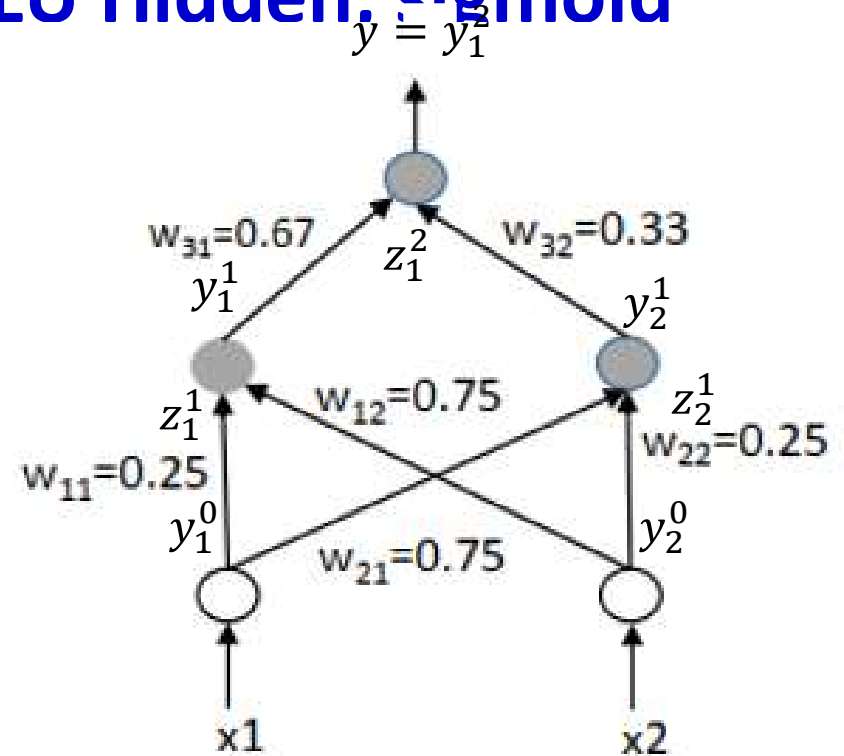
$$\text{div} = -d * \log(y) - (1 - d) * \log(1 - y)$$

$$\frac{\delta \text{div}}{\delta y} = -d/y + (1 - d)/(1 - y) = 1/(1 - y) = (1 + e)$$

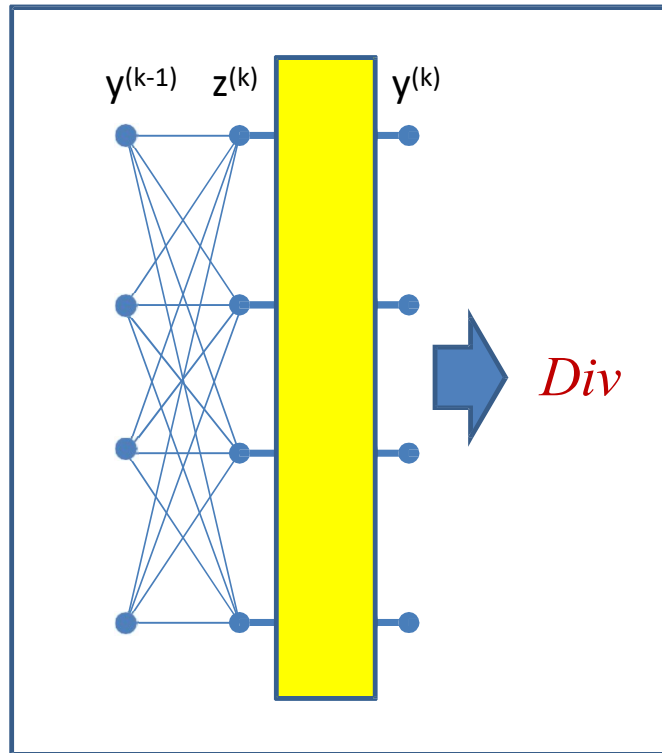
$$\frac{\delta \text{div}}{\delta z_1^2} = \frac{\delta \text{Div}}{\delta y_1^2} * \frac{\delta y_1^2}{\delta z_1^2} = (1 + e) * \text{sigmoid}(z_1^2) * (1 - \text{sigmoid}(z_1^2)) = e / (1 + e)$$

$$\frac{\delta \text{div}}{\delta w_{31}} = \frac{\delta \text{Div}}{\delta z_1^2} * \frac{\delta z_1^2}{\delta w_{31}} = e / (1 + e) * y_1^1 = e / (1 + e)$$

$$w_{31} = w_{31} - \eta * \frac{\delta \text{div}}{\delta w_{31}} = 0.67 - 0.1 * e / (1 + e) = 0.597$$



Example Special Case : Softmax



$$y_i^{(k)} = \frac{\exp(z_i^{(k)})}{\sum_j \exp(z_j^{(k)})}$$

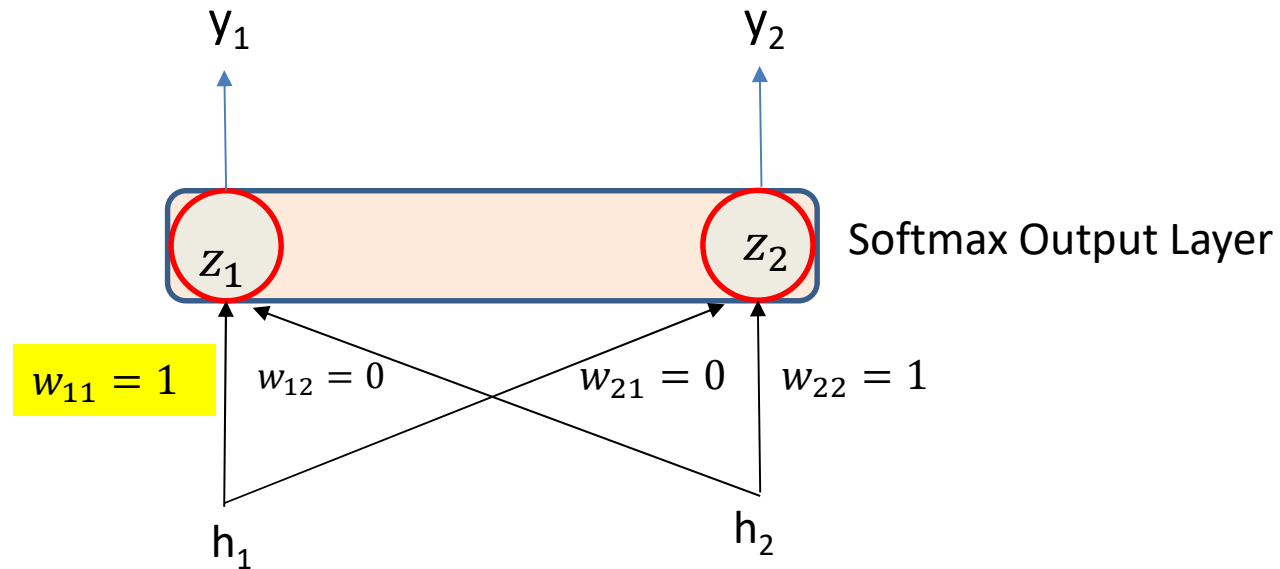
$$\frac{\partial Div}{\partial z_i^{(k)}} = \sum_j \frac{\partial Div}{\partial y_j^{(k)}} \frac{\partial y_j^{(k)}}{\partial z_i^{(k)}}$$

$$\frac{\partial y_j^{(k)}}{\partial z_i^{(k)}} = \begin{cases} y_i^{(k)} (1 - y_i^{(k)}) & \text{if } i = j \\ - y_i^{(k)} y_j^{(k)} & \text{if } i \neq j \end{cases}$$

$$\frac{\partial Div}{\partial z_i^{(k)}} = \sum_j \frac{\partial Div}{\partial y_j^{(k)}} y_i^{(k)} (\delta_{ij} - y_j^{(k)})$$

- For future reference
- δ_{ij} is the Kronecker delta: $\delta_{ij} = 1$ if $i = j$, 0 if $i \neq j$

Weight Updates with Softmax



- Training input: $(h_1, h_2) = (1, 0)$ target output: $(d_1, d_2) = (0, 1)$
 - *div* function: cross-entropy
 - Learning rate: 0.1
 - Bias = 0
 - What will be the value of w_{11} in next iteration?
-
- Input to softmax node $z_1 = w_{11} = 1$; $z_2 = 0$
 - $y_1 = e / (1+e)$
 - $y_2 = 1 / (1+e)$
 - $div = -d_1 \log(y_1) - d_2 \log(y_2) = -\log(y_2) = \log(1+e) = 1.313$

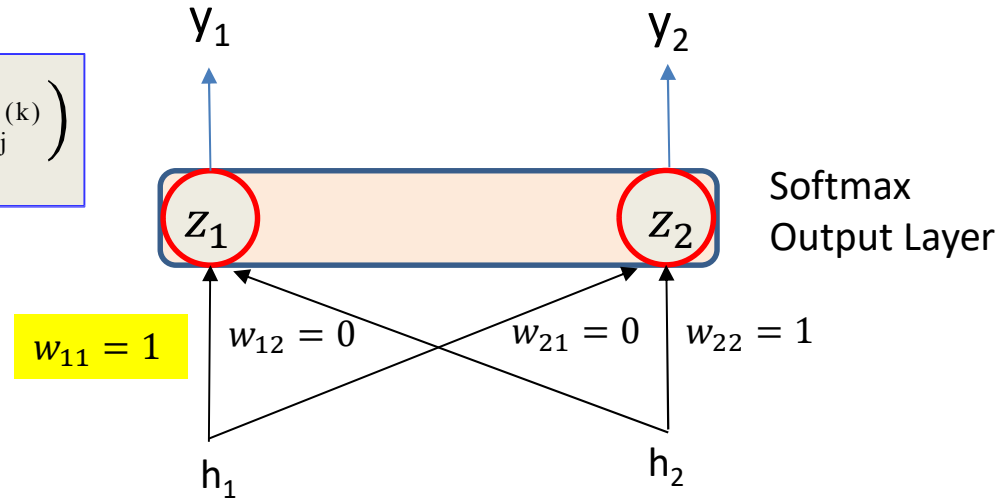
Weight Updates with Softmax ...

$$\frac{\partial Div}{\partial z_i^{(k)}} = \sum_j \frac{\partial Div}{\partial y_j^{(k)}} y_i^{(k)} (\delta_{ij} - y_j^{(k)})$$

$$w_{11}(t+1) = ?$$

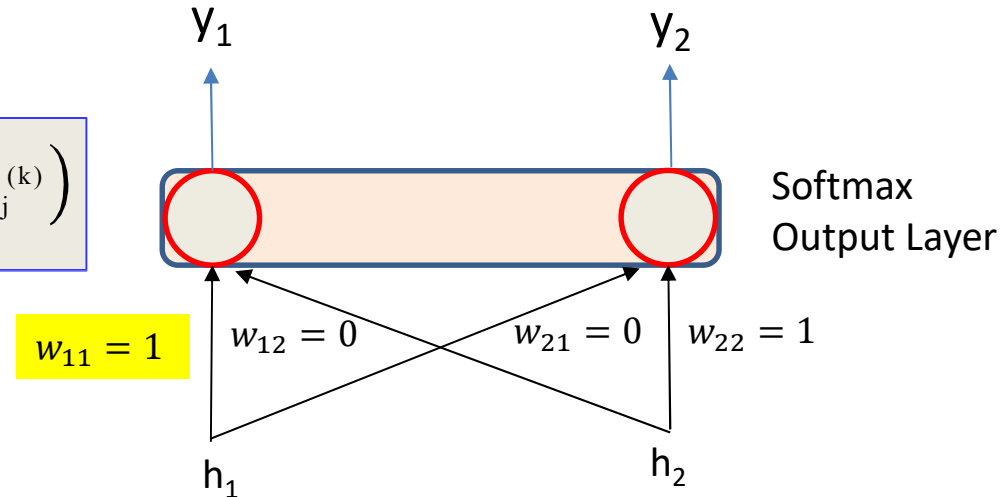
Training input: $(h_1, h_2) = (1, 0)$

Target output: $(d_1, d_2) = (0, 1)$



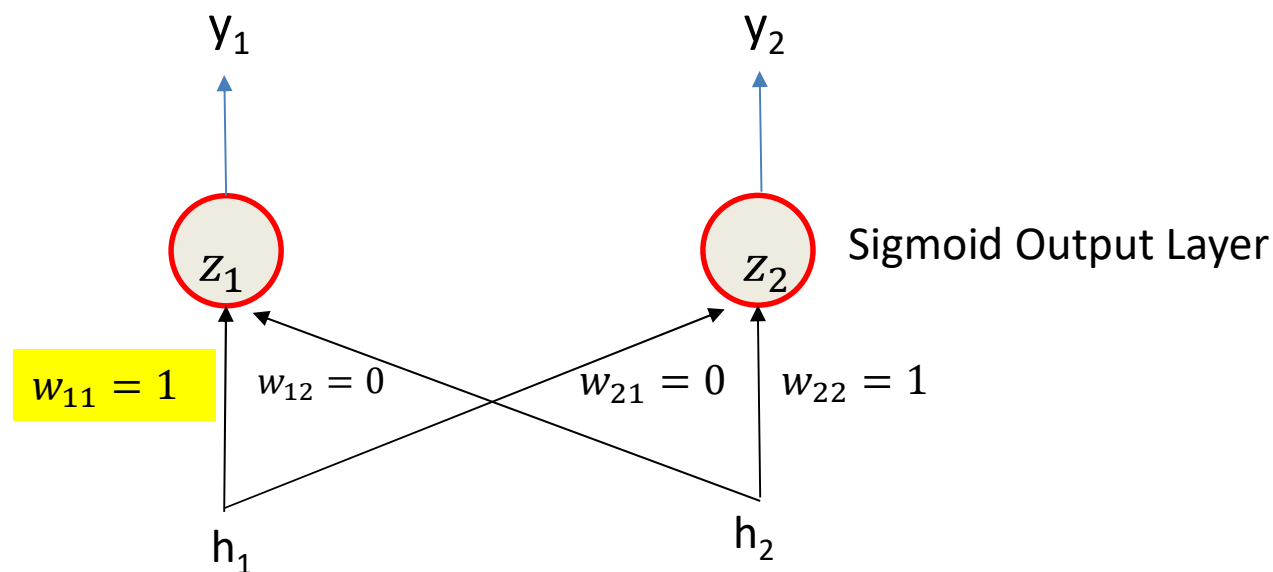
Weight Updates with Softmax ...

$$\frac{\partial Div}{\partial z_i^{(k)}} = \sum_j \frac{\partial Div}{\partial y_j^{(k)}} y_i^{(k)} (\delta_{ij} - y_j^{(k)})$$



- Change in $w_{11} = \delta w_{11} = -0.1 * ddiv/dw_{11}$
 $= -0.1 * ddiv/dz_1 * dz_1/dw_{11} = -0.1 * ddiv/dz_1 * h_1$
- $ddiv/dy_1 = -d_1/y_1 = 0$ $ddiv/dy_2 = -d_2/y_2 = -1/y_2$
- $ddiv/dz_1 = ddiv/dy_1 * y_1(1-y_1) + ddiv/dy_2 * y_1(-y_2) = -y_1 y_2 / y_2 = -y_1$
- Input to softmax node $z_1 = w_{11} = 1$; $z_2 = 0$; $y_1 = e / (1+e)$
- $\delta w_{11} = -0.1 * e/(1+e) = -0.0731$

Weight Updates with Sigmoid Output



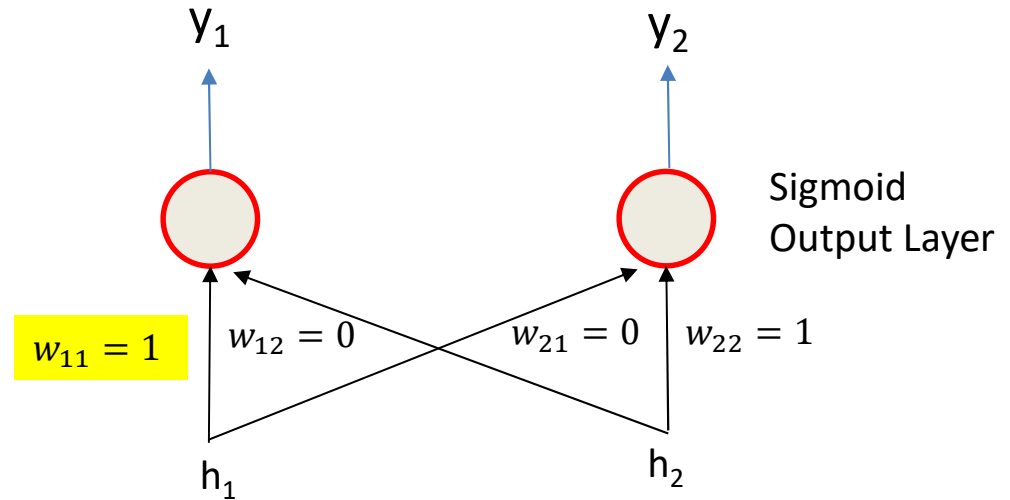
- Training input: $(h_1, h_2) = (1, 0)$ target output: $(d_1, d_2) = (0, 1)$
 - Activation: sigmoid, Bias=0
 - *div* function: cross-entropy
 - Learning rate: 0.1
 - What will be the value of w_{11} in next iteration?
-
- Input to output node $z_1 = w_{11} = 1$; $z_2 = 0$
 - $y_1 = 1 / (1 + e^{-1}) = e / (1 + e)$ $y_2 = 1 / (1 + e^0) = 1/2$
 - $div = -d_1 \log(y_1) - d_2 \log(y_2) = -\log(y_2) = \log 2$

Weight Updates with Sigmoid Outputs ...

$$w_{11}(t+1)=?$$

Input $(h_1, h_2)=(1.0, 0.0)$

target output: $(d_1, d_2)=(0,1)$



Weight Updates with Sigmoid Outputs ...

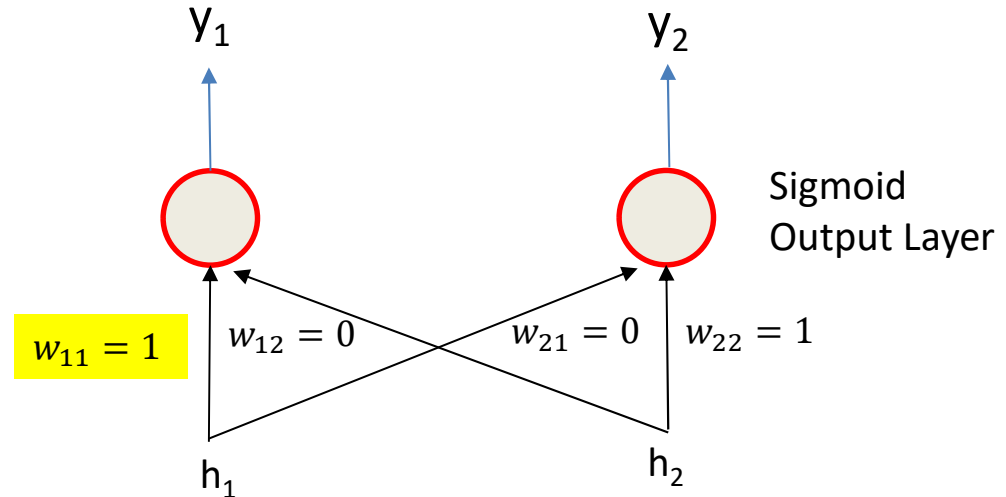
Input $(h_1, h_2) = (1.0, 0.0)$
 target output: $(d_1, d_2) = (0, 1)$

$w_{11}(t+1) = ?$

$div = -d_1 \log(y_1) - d_2 \log(y_2)$

$$\frac{\delta div}{\delta y_1} = -\frac{d_1}{y_1}$$

$$\frac{\delta y_1}{\delta z_1} = \frac{1}{1 + e^{-1}} * \left(1 - \frac{1}{1 + e^{-1}}\right) = \frac{e}{(1 + e)^2}$$



$$\begin{aligned} w_{11}(t + 1) &= w_{11}(t) - \eta \frac{\delta div}{\delta w_{11}} \\ &= w_{11}(t) - 0.1 * \frac{\delta div}{\delta y_1} \frac{\delta y_1}{\delta z_1} \frac{\delta z_1}{\delta w_{11}} \\ &= w_{11}(t) - 0.1 * 0 * \frac{e}{(1+e)^2} * h_1 \\ &= 1.0 \end{aligned}$$

L2 Regularization

$$\text{Loss}(x,y)=3x^2+6xy+y^2$$

What is the value of (x_{opt}, y_{opt}) that minimizes the loss, if L2 regularization constant 0.1 is applied?

Now, for regularized case, optimal θ_R^* is given by

$$\theta_R^* \approx (H + 0.1 I)^{-1} H \theta^*$$

where H is the Hessian of Loss(x,y) and θ^* is the optimal value in the unregularized case.

L2 Regularization

$\text{Loss}(x,y)=3x^2+6xy+y^2$. What is the value of $(x_{\text{opt}}, y_{\text{opt}})$ that minimizes the loss, if L2 regularization constant 0.1 is applied?

$d \text{ Loss} / dx = 6x + 6y = 0$ and $d \text{ Loss} / dy = 6x+2y = 0$ at $(x_{\text{opt}}, y_{\text{opt}})$. Solving these two equations, we get, in unregularized case, $x_{\text{opt}} = y_{\text{opt}} = 0$.

Now, for regularized case, optimal θ_R^* is given by

where H is the Hessian of $\text{Loss}(x,y)$ and θ^* is the optimal value in the unregularized case. $\theta_R^* \approx (H + 0.1 I)^{-1} H \theta^*$

Since in unregularized case, $x_{\text{opt}}=y_{\text{opt}}=0$, θ_R^* is also a zero vector.

L1 Regularization

$$\text{Loss}(x, y) = 3x^2 - 6x + 3 + y^2$$

What is the value of $(x_{\text{opt}}, y_{\text{opt}})$ that minimizes the loss, if L1 regularization constant 0.1 is applied?

So, in regularized case $(\theta_R^*)_i \approx \begin{cases} \max \left\{ \theta^* - \frac{\alpha}{H_{ii}}, 0 \right\} & \text{if } \theta^* \geq 0 \\ \min \left\{ \theta^* + \frac{\alpha}{H_{ii}}, 0 \right\} & \text{if } \theta^* < 0 \end{cases}$ $x_{\text{opt}} = 1/6 = 0.1667$ and $y_{\text{opt}} = \max(0 - 0.1/2, 0) = 0$.

L1 Regularization

$\text{Loss}(x, y) = 3x^2 - 6x + 3 + y^2$. What is the value of $(\underline{x}_{opt}, \underline{y}_{opt})$ that minimizes the loss, if L1 regularization constant 0.1 is applied?

$H_{11} = d^2 \text{Loss}(x, y) / dx^2 = 6$, $H_{12} = H_{21} = d^2 \text{Loss}(x, y) / dx dy = 0$ and $H_{22} = d^2 \text{Loss}(x, y) / dy^2 = 2$. So, H is diagonal and +ve definite. Hence

$$(\theta_R^*)_i \approx \begin{cases} \max\left\{\theta^* - \frac{\alpha}{H_{ii}}, 0\right\} & \text{if } \theta^* \geq 0 \\ \min\left\{\theta^* + \frac{\alpha}{H_{ii}}, 0\right\} & \text{if } \theta^* < 0 \end{cases}$$

$d \text{Loss}(x, y) / dx = 6x - 6 = 0$, and $d \text{Loss}(x, y) / dy = 2y = 0$ at unregularized (x_{opt}, y_{opt}) . Solving these we get, in unregularized case, $x_{opt} = 1$ and $y_{opt} = 0$.

So, in regularized case, $x_{opt} = \max(1 - 0.1/6, 0) = 0.9833$ and $y_{opt} = \max(0 - 0.1/2, 0) = 0$.