# Loss function and Back propagation

Prof. (Dr.) Keyur Rana

---

- Loss function / Error function
  - Sum of Square error (quadratic error) – to minimize

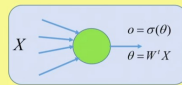$$Div(Y,d) = \frac{1}{2}\|Y - d\|^2 = \frac{1}{2}\sum_i (y_i - d_i)^2$$

---

## Cross entropy loss – Two class problem

$o \Rightarrow$ likelihood that $y$ is 1

$(1-o) \Rightarrow$ likelihood that $y$ is 0

Likelihood that is to be maximized $\Rightarrow o^y(1-o)^{(1-y)}$

Loglikelihood $\Rightarrow y \log o + (1-y)\log(1-o)$
to maximized

Diagram: $X$ → node, $o = \sigma(\theta)$, $\theta = W^t X$

---

## Cross entropy loss – Two class problem

$$\text{Minimize} \Rightarrow C = -\frac{1}{N}\sum_{\forall X}\left[ y \log o + (1-y)\log(1-o)\right]$$

$$\frac{\partial C}{\partial W_i} = -\frac{1}{N}\sum_{\forall X}\left[\frac{y}{\sigma(\theta)} - \frac{(1-y)}{1-\sigma(\theta)}\right]\frac{\partial\sigma(\theta)}{\partial W_i}$$

$$= -\frac{1}{N}\sum_{\forall X}\left[\frac{y}{\sigma(\theta)} - \frac{(1-y)}{1-\sigma(\theta)}\right]\frac{\partial\sigma(\theta)}{\partial\theta}\cdot\frac{\partial\theta}{\partial W_i} \quad \text{(Chain rule is applied)}$$
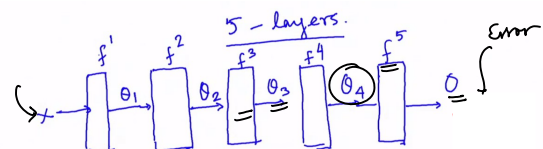
o is $\sigma(\theta)$ = Sigmoidal ($\theta$), here $\theta$ = W$^T$X

---

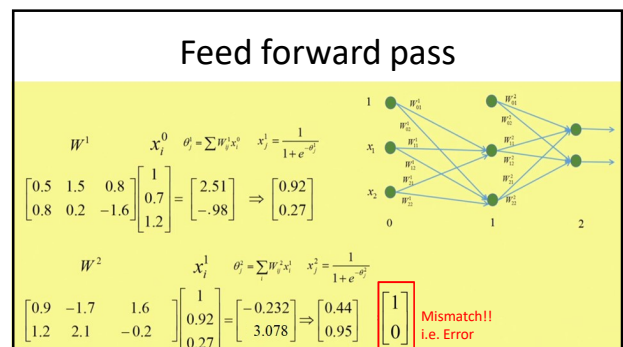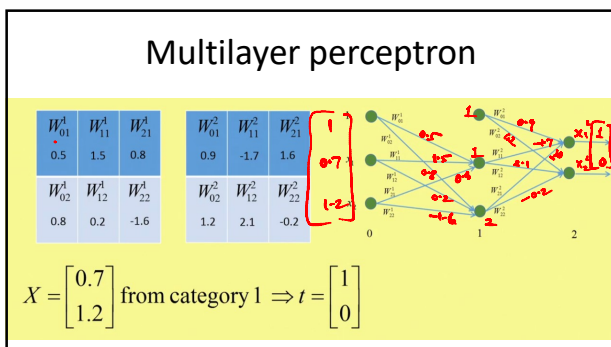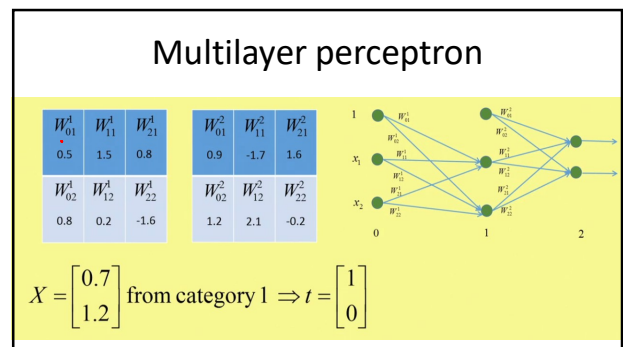## Cross entropy loss – Two class problem

Simplifying further…

$$\frac{\partial C}{\partial W_i} = -\frac{1}{N}\sum_{\forall X}\left[\frac{y}{\sigma(\theta)} - \frac{(1-y)}{1-\sigma(\theta)}\right]\frac{\partial\sigma(\theta)}{\partial\theta}\cdot\frac{\partial\theta}{\partial W_i}$$

$$= -\frac{1}{N}\sum_{\forall X}\left[\frac{y}{\sigma(\theta)} - \frac{(1-y)}{1-\sigma(\theta)}\right]\frac{\partial\sigma(\theta)}{\partial\theta}\cdot\frac{\partial\theta}{\partial W_i}$$

$$= -\frac{1}{N}\sum_{\forall X}\left[\frac{y-\sigma(\theta)}{\sigma(\theta)(1-\sigma(\theta))}\right]\sigma(\theta)(1-\sigma(\theta))x_i$$

$$= \frac{1}{N}\sum_{\forall X}x_i(\sigma(\theta) - y) \qquad = \frac{1}{N}\sum_{\forall X}x_i(o - y)$$

**Derivative** of Sigmodial function $\sigma(x)$ is **σ(x)(1- σ(x))**

$\theta$ = W$^T$X
σ(W$_i$X$_i$) = X$_i$

---

**Slide 1 (top-left, handwritten):**

5 - layers.

$f^1$ $f^2$ $f^3$ $f^4$ $f^5$

$x \rightarrow \theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \rightarrow O \quad \frac{\partial O}{\partial X}$

$O = f^5\left(f^4\left(f^3\left(f^2\left(f(x)\right)\right)\right)\right)$

Gradient of the Error

Local gradient

$\frac{\partial O}{\partial X} = \frac{\partial O}{\partial \theta_4} \cdot \frac{\partial \theta_4}{\partial \theta_3} \cdot \frac{\partial \theta_3}{\partial \theta_2} \cdot \frac{\partial \theta_2}{\partial \theta_1} \cdot \frac{\partial \theta_1}{\partial X}$

Gradients which are back propagated in to the input layers

**Slide 2 (top-right):**

# Multilayer perceptron

$1 \quad W^1_{01}$ $\quad W^2_{01}$

$W^1_{02}$ $\quad W^2_{02}$ $\quad X_1^2$

$x_1 \quad W^1_{11}$ $\quad W^2_{11}$ $\quad X_2^2$

$W^1_{12}$ $\quad W^2_{12}$

$W^1_{21}$ $\quad W^2_{21}$

$x_2 \quad W^1_{22}$ $\quad W^2_{22}$

$0 \quad\quad 1 \quad\quad 2$

**Slide 3 (middle-left):**

$W^k_{ij} \Rightarrow$ weight from node $i$ (from $k-1$ layer) to node $j$ at $k^{th}$ layer

# Multilayer perceptron

$\delta^k_i$ → local grad $i$ at $k^{th}$ layer

$x^k_i$ → gradient at node $i$ at $k^{th}$ layer

$1 \quad W^1_{01}$ $\quad W^2_{01}$

$W^1_{02}$ $\quad W_{01}^1 \quad W^2_{02}$ $\quad \delta^2_1 \quad x_1^2$

$x_1 \quad W^1_{11} \quad \delta^1_1 \quad W^2_{11}$

$W^1_{12} \quad \delta^1_1 \quad W^2_{12}$ $\quad X_2^2$

$W^1_{21} \quad W^2_{21}$ $\quad \delta^2_2 \quad x^2_2$

$x_2 \quad W^1_{22} \quad \delta^1_2 \quad W^2_{22}$

$0 \quad\quad 1 \quad\quad 2$

**Slide 4 (middle-right):**

# Multilayer perceptron

| $W^1_{01}$ | $W^1_{11}$ | $W^1_{21}$ | $W^2_{01}$ | $W^2_{11}$ | $W^2_{21}$ |
|---|---|---|---|---|---|
| 0.5 | 1.5 | 0.8 | 0.9 | -1.7 | 1.6 |
| $W^1_{02}$ | $W^1_{12}$ | $W^1_{22}$ | $W^2_{02}$ | $W^2_{12}$ | $W^2_{22}$ |
| 0.8 | 0.2 | -1.6 | 1.2 | 2.1 | -0.2 |

$X = \begin{bmatrix} 0.7 \\ 1.2 \end{bmatrix}$ from category $1 \Rightarrow t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

**Slide 5 (bottom-left):**

# Multilayer perceptron

| $W^1_{01}$ | $W^1_{11}$ | $W^1_{21}$ | $W^2_{01}$ | $W^2_{11}$ | $W^2_{21}$ |
|---|---|---|---|---|---|
| 0.5 | 1.5 | 0.8 | 0.9 | -1.7 | 1.6 |
| $W^1_{02}$ | $W^1_{12}$ | $W^1_{22}$ | $W^2_{02}$ | $W^2_{12}$ | $W^2_{22}$ |
| 0.8 | 0.2 | -1.6 | 1.2 | 2.1 | -0.2 |

$X = \begin{bmatrix} 0.7 \\ 1.2 \end{bmatrix}$ from category $1 \Rightarrow t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

**Slide 6 (bottom-right):**

# Feed forward pass

$W^1 \qquad x^0_i \quad \theta^1_j = \sum W^1_{ij} x^0_i \quad x^1_j = \frac{1}{1+e^{-\theta^1_j}}$

$\begin{bmatrix} 0.5 & 1.5 & 0.8 \\ 0.8 & 0.2 & -1.6 \end{bmatrix} \begin{bmatrix} 1 \\ 0.7 \\ 1.2 \end{bmatrix} = \begin{bmatrix} 2.51 \\ -.98 \end{bmatrix} \Rightarrow \begin{bmatrix} 0.92 \\ 0.27 \end{bmatrix}$

$W^2 \qquad x^1_i \quad \theta^2_j = \sum W^2_{ij} x^1_i \quad x^2_j = \frac{1}{1+e^{-\theta^2_j}}$

$\begin{bmatrix} 0.9 & -1.7 & 1.6 \\ 1.2 & 2.1 & -0.2 \end{bmatrix} \begin{bmatrix} 1 \\ 0.92 \\ 0.27 \end{bmatrix} = \begin{bmatrix} -0.232 \\ 3.078 \end{bmatrix} \Rightarrow \begin{bmatrix} 0.44 \\ 0.95 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ Mismatch!! i.e. Error
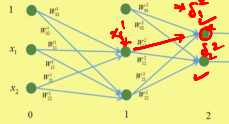
## Back propagation – Output layer

$$E = \frac{1}{2}\sum_{j=1}^{2}\left(x_j^2 - t_j\right)^2 \qquad x_j^2 = \frac{1}{1+e^{-\theta_j^2}} \qquad \theta_j^2 = \sum_{i=0}^{2}W_{ij}^2 x_i^1$$

$$\frac{\partial E}{\partial W_{ij}^2} = \frac{\partial E}{\partial x_j^2}\cdot\frac{\partial x_j^2}{\partial \theta_j^2}\cdot\frac{\partial \theta_j^2}{\partial W_{ij}^2} = (x_j^2 - t_j)x_j^2(1-x_j^2)x_i^1$$

We set $\boxed{\delta_j^2 = x_j^2(1-x_j^2)(x_j^2 - t_j)}$ $\Rightarrow \frac{\partial E}{\partial W_{ij}^2} = \delta_j^2 x_i^1$

$$\boxed{W_{ij}^2 \leftarrow W_{ij}^2 - \eta\,\frac{\partial E}{\partial W_{ij}^2}}$$

---

## Back propagation – Output layer

$$\boxed{\delta_j^2 = x_j^2(1-x_j^2)(x_j^2 - t_j)}$$

$\delta_1^2 = x_1^2(1-x_1^2)(x_1^2 - t_1)$
$= 0.44*(1-0.44)*(0.44-1)$
$= -0.138$
$\Rightarrow \frac{\partial E}{\partial W_{11}^2} = \delta_1^2 x_1^1 = -0.126$

$\delta_2^2 = x_2^2(1-x_2^2)(x_2^2 - t_2)$
$= 0.95*(1-0.95)*(0.95-0)$
$= 0.045$
$\Rightarrow \frac{\partial E}{\partial W_{12}^2} = \delta_2^2 x_1^1 = 0.04$

Observe the updated weights

$W_{11}^2 \leftarrow W_{11}^2 + \eta * 0.126$   $W_{12}^2 \leftarrow W_{12}^2 - \eta * 0.04$

---

## Back propagation – Output layer

$\frac{\partial E}{\partial W_{21}^2} = \delta_1^2 x_2^1 = -0.037$   $\frac{\partial E}{\partial W_{22}^2} = \delta_2^2 x_2^1 = 0.012$
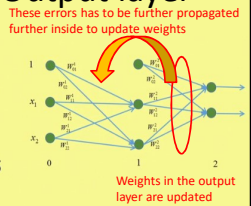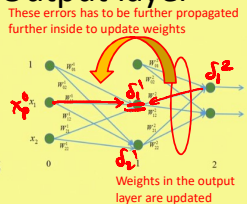
$\frac{\partial E}{\partial W_{01}^2} = \delta_1^2 x_0^1 = -0.138$   $\frac{\partial E}{\partial W_{02}^2} = \delta_2^2 x_0^1 = 0.045$

Error from output is propagated in the backward direction

Weights in the output layer are updated

---

## Back propagation – Output layer

These errors has to be further propagated further inside to update weights

$\frac{\partial E}{\partial W_{21}^2} = \delta_1^2 x_2^1 = -0.037$   $\frac{\partial E}{\partial W_{22}^2} = \delta_2^2 x_2^1 = 0.012$

$\frac{\partial E}{\partial W_{01}^2} = \delta_1^2 x_0^1 = -0.138$   $\frac{\partial E}{\partial W_{02}^2} = \delta_2^2 x_0^1 = 0.045$

Weights in the output layer are updated

---

## Back propagation – Output layer

These errors has to be further propagated further inside to update weights

$\frac{\partial E}{\partial W_{21}^2} = \delta_1^2 x_2^1 = -0.037$   $\frac{\partial E}{\partial W_{22}^2} = \delta_2^2 x_2^1 = 0.012$

$\frac{\partial E}{\partial W_{01}^2} = \delta_1^2 x_0^1 = -0.138$   $\frac{\partial E}{\partial W_{02}^2} = \delta_2^2 x_0^1 = 0.045$

Weights in the output layer are updated

$$\frac{\partial E}{\partial w_{11}} = \delta_1^1 x_1^0$$

---

## Back propagated Error term

$$\boxed{\delta_i^k = X_i^k(1-X_i^k)\sum_{j=1}^{M_{k+1}}\delta_j^{k+1}W_{ij}^{k+1}}$$

- $\delta_i^k$ : Back propagated Error Term at layer k for the $i^{th}$ node
- All the nodes to which $i^{th}$ node has fed the input, **accumulate** the error corresponding to the weights and Multiply with the local derivative *xi(1-xi)*

## Back propagation – Hidden layer

Local derivative of the output of the ith node

We set $\boxed{\delta_i^k = \boxed{X_i^k(1-X_i^k)}\sum_{j=1}^{M_{k+1}}\delta_j^{k+1}W_{ij}^{k+1}}$ $\Rightarrow$ $\delta_i^1 = x_i^1(1-x_i^1)\sum_{j=1}^{2}\delta_j^2 W_{ij}^2$ $\Rightarrow$ $\boxed{\dfrac{\partial E}{\partial W_{ij}^k} = \delta_j^k x_i^{k-1}}$

Propagated error

$\delta_1^1 = x_1^1(1-x_1^1)\left[\delta_1^2 * W_{11}^2 + \delta_2^2 W_{12}^2\right]$
$= 0.92*(1-0.92)[(-0.137)*(-1.7)+0.045*2.1]$
$= 0.024$

$\delta_2^1 = x_2^1(1-x_2^1)\left[\delta_1^2 * W_{21}^2 + \delta_2^2 W_{22}^2\right]$
$= 0.27*(1-0.27)[(-0.137)*0.8+0.045*(-0.2)]$
$= -0.02$

## Back propagation – Hidden layer

$\dfrac{\partial E}{\partial W_{11}^1} = \delta_1^1 * x_1^0 = 0.024*0.7 = 0.017$

$\dfrac{\partial E}{\partial W_{12}^1} = \delta_2^1 * x_1^0 = -0.02*0.7 = -0.014$

$\dfrac{\partial E}{\partial W_{21}^1} = \delta_1^1 * x_2^0 = 0.024*1.2 = 0.0288$

$\dfrac{\partial E}{\partial W_{22}^1} = \delta_2^1 * x_2^0 = -0.02*1.2 = -0.024$

$\dfrac{\partial E}{\partial W_{01}^1} = \delta_1^1 * x_0^0 = 0.024*1 = 0.024$

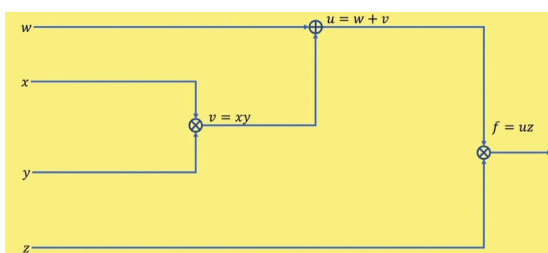$\dfrac{\partial E}{\partial W_{02}^1} = \delta_2^1 * x_0^0 = -0.02*1 = -0.02$



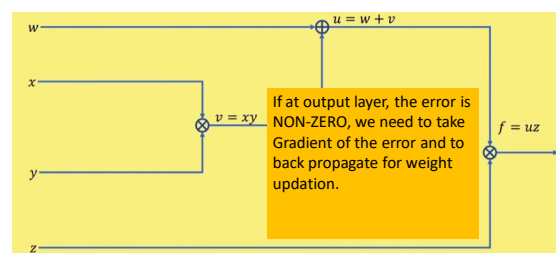$$W_{ij}^1 \leftarrow W_{ij}^1 - \eta \frac{\partial E}{\partial W_{ij}^1}$$

## Back propagation

- Applicable to any number of layers
- Highly parallelism possible
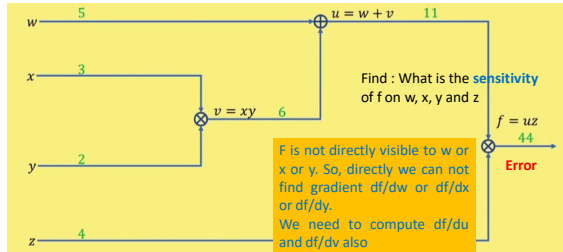
## Back propagation learning at the Node level
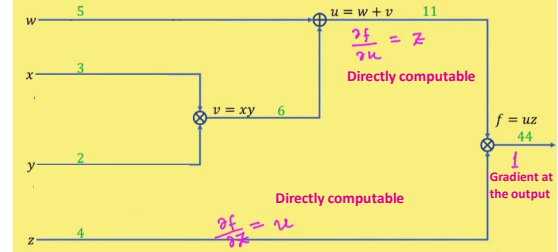
## Back propagation



$w$
$x$
$u = w + v$
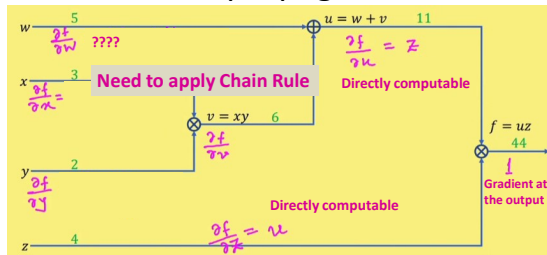$v = xy$
$y$
$f = uz$
$z$

## Back propagation



$w$
$x$
$u = w + v$
$v = xy$

If at output layer, the error is NON-ZERO, we need to take Gradient of the error and to back propagate for weight updation.

$y$
$f = uz$
$z$

## Slide 27 (top-left)

# Back propagation

$w$ — 5
$u = w + v$    11
$x$ — 3
$v = xy$    6
$y$ — 2
$f = uz$    44
$z$ — 4

Find : What is the **sensitivity** of f on w, x, y and z

F is not directly visible to w or x or y. So, directly we can not find gradient df/dw or df/dx or df/dy.
We need to compute df/du and df/dv also

**Error**

## Slide 27 (top-right)

# Back propagation

$w$ — 5
$u = w + v$    11
$\frac{\partial f}{\partial u} = z$
**Directly computable**
$x$ — 3
$v = xy$    6
$f = uz$    44
$y$ — 2
$z$ — 4

**Directly computable**

$\frac{\partial f}{\partial z} = u$

Gradient at the output

27

## Slide 28 (middle-left)

# Back propagation

$w$ — 5   $\frac{\partial f}{\partial w}$ ????
$u = w + v$    11
$\frac{\partial f}{\partial u} = z$
**Need to apply Chain Rule**
**Directly computable**
$x$ — 3   $\frac{\partial f}{\partial x} =$
$v = xy$    6
$\frac{\partial f}{\partial v}$
$f = uz$    44
$y$ — 2   $\frac{\partial f}{\partial y}$
**Directly computable**
Gradient at the output
$z$ — 4   $\frac{\partial f}{\partial z} = u$

28

## Slide 29 (middle-right)

**Back propagated gradient**

# Back propagation

$w$ — 5   $\frac{\partial f}{\partial w} = \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial w}$
$u = w + v$    11
$\frac{\partial f}{\partial u} = z$
$x$ — 3   $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial v} \cdot \frac{\partial v}{\partial x}$
**Local gradient**
$v = xy$    6   $\frac{\partial f}{\partial v} = \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial v}$
**Local gradient**
$f = uz$    44
$y$ — 2   $\frac{\partial f}{\partial y} = \frac{\partial f}{\partial v} \cdot \frac{\partial v}{\partial y}$
**Back propagated gradient**
Gradient at the output
$z$ — 4   $\frac{\partial f}{\partial z} = u$

29

## Slide 30 (bottom-left)

# Back propagation

```
# Set Input
  w=5; x=3; y=2; z=4

# Forward Pass      # Backward Pass

  v = x*y             dfdu = z
  u = w+v             dfdz = u
  f = u+z             dfdw = 1*dfdu   # dudw = 1
                      dfdv = 1*dfdu   # dudv = 1
                      dfdx = y*dfdv   # dvdx = y
                      dfdy = x*dfdv   # dvdy = x
```
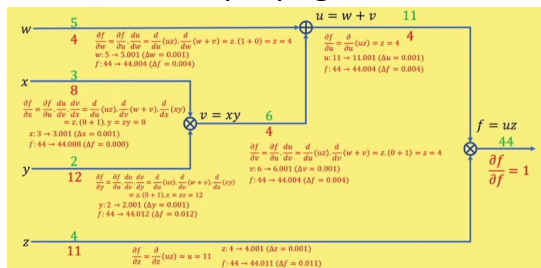
30

## Slide 31 (bottom-right)

# Back propagation

```
# Set Input
  w=5; x=3; y=2; z=4

# Forward Pass      # Backward Pass      sensitivity of f
                                          on w, x, y and z
  v = x*y             dfdu = z
  u = w+v             dfdz = u
  f = u+z             dfdw = 1*dfdu   # dudw = 1
                      dfdv = 1*dfdu   # dudv = 1
                      dfdx = y*dfdv   # dvdx = y
                      dfdy = x*dfdv   # dvdy = x
```
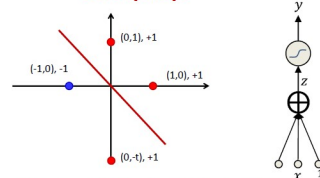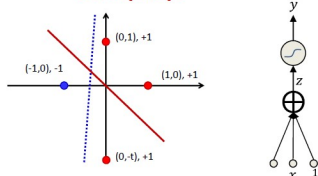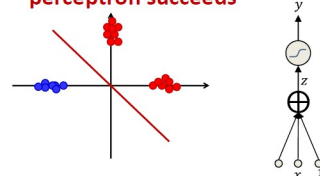
31

## Back propagation



$$u = w + v$$

$$v = xy$$

$$f = uz$$

$$\frac{\partial f}{\partial f} = 1$$

---

## Backprop



(0,1), +1
(-1,0), -1
(1,0), +1
(0,-t), +1

- Local optimum solution found by backprop
- Does not separate the points *even though the points are linearly separable!*
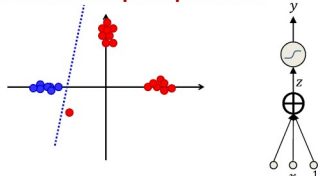
21

---

## Backprop



(0,1), +1
(-1,0), -1
(1,0), +1
(0,-t), +1

- Solution found by backprop
- Does not separate the points *even though the points are linearly separable!*
- Compare to the perceptron: *Backpropagation fails to separate where the perceptron succeeds*

22

---

## Backprop fails to separate where perceptron succeeds



- Brady, Raghavan, Slawny, '89
- *Several* linearly separable training examples
- Simple setup: **both backprop and perceptron algorithms find solutions**
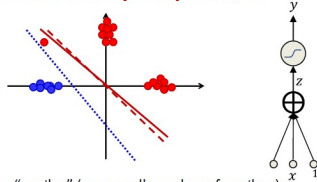
23

---

## A more complex problem



- Adding a "spoiler" (or a small number of spoilers)
  - Perceptron finds the linear separator,

24

---

## A more complex problem



- Adding a "spoiler" (or a small number of spoilers)
  - Perceptron finds the linear separator,
  - Backprop does not find a separator
    - A single additional input does not change the loss function significantly
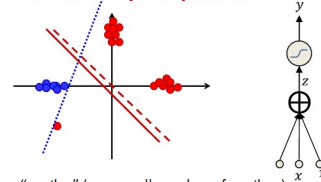      - **Assuming weights are constrained to be bounded**

25

## A more complex problem



- Adding a "spoiler" (or a small number of spoilers)
  - Perceptron finds the linear separator,
  - For bounded *w*, backprop does not find a separator
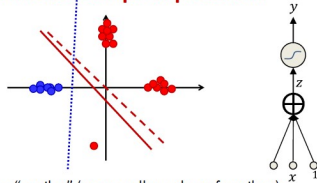    - A single additional input does not change the loss function significantly

26

## A more complex problem



- Adding a "spoiler" (or a small number of spoilers)
  - Perceptron finds the linear separator,
  - For bounded *w*, backprop does not find a separator
    - A single additional input does not change the loss function significantly

27

## A more complex problem



- Adding a "spoiler" (or a small number of spoilers)
  - Perceptron finds the linear separator,
  - For bounded *w*, backprop does not find a separator
    - A single additional input does not change the loss function significantly

28

## So what is happening here?

- The perceptron may change greatly upon adding just a single new training instance
  - But it fits the training data well
  - The perceptron rule has *low bias*
    - Makes no errors if possible
  - But high variance
    - Swings wildly in response to small changes to input
- Backprop is minimally changed by new training instances
  - Prefers consistency over perfection
  - It is a *low-variance* estimator, at the potential cost of bias

29

## Backpropagation: Finding the separator

- Backpropagation will often not find a separating solution *even though the solution is within the class of functions learnable by the network*

- This is because the separating solution is not a feasible optimum for the loss function

- One resulting benefit is that a backprop-trained neural network classifier has lower variance than an optimal classifier for the training data

32