

# Non-linear Hypothesis

## Introduction of Multilayer Perceptron

**Machine Learning (CS 306)**

**Instructor:** Dr. Moumita Roy

**Teaching Assistants:** Indrajit Kalita, Veronica Naosekpam

Email-ids:

moumita@iiitg.ac.in

veronica.naosekpam@iiitg.ac.in

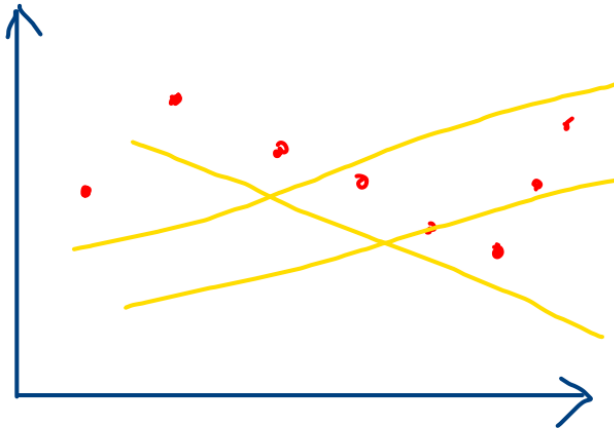
indrajit.kalita@iiitg.ac.in

**Mobile No:** +91-8420489325 (only for emergency quires)

**Reference:** <https://www.cse.iitm.ac.in/~miteshk/CS6910/Slides/Lecture2.pdf>

# Regression vs. Classification (considering linear hypothesis)

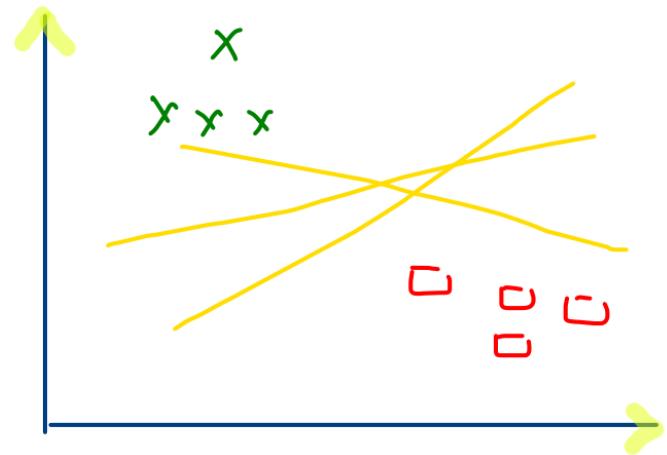
Regression



Linear Relation between independent and dependent variables

Try to find best fit one by minimizing cost function

classification



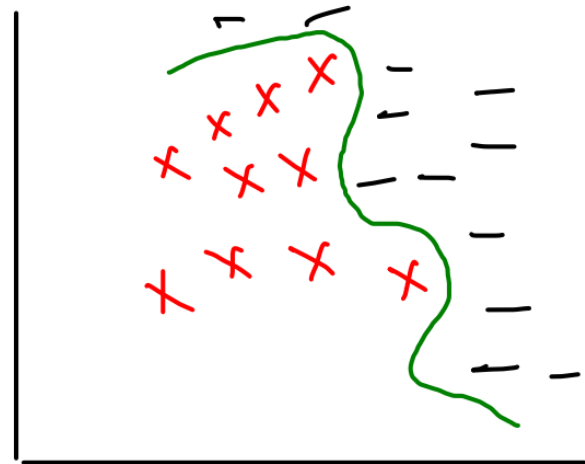
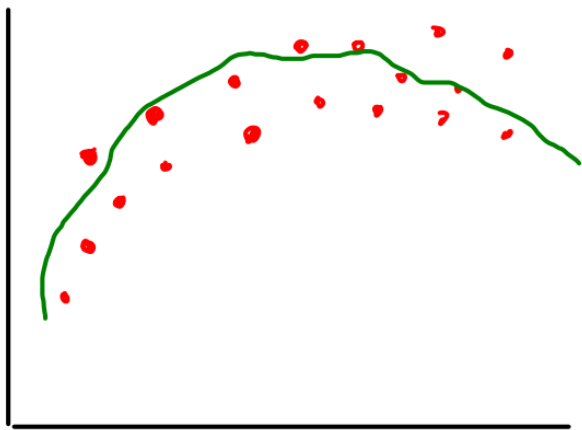
Data set is linearly separable

Try to find linear DB by minimizing cost function

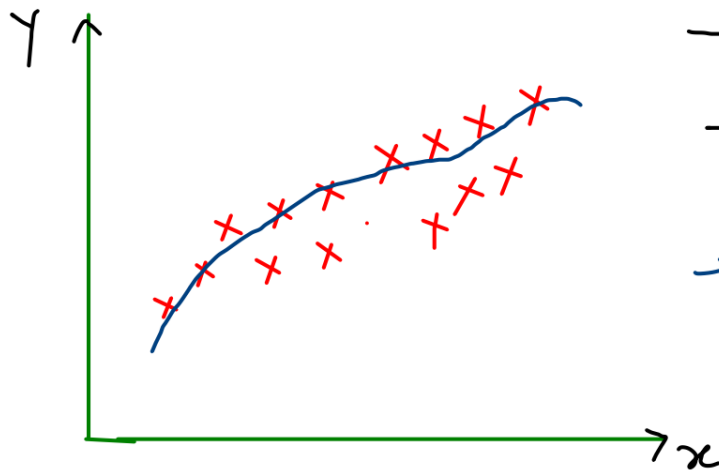
Consider binary classification

# Scenarios

## (Non-linear hypothesis)



# Polynomial Regression



$$\rightarrow h(x) = w_0 + w_1 x$$

$$\rightarrow h(x) = w_0 + w_1 x + w_2 x^2$$

$$\rightarrow h(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

If, we consider,  $h(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$

Can we use linear regression to  
Solve the problem?

Introduce more features

$$x_1 = x$$

$$x_2 = x^2$$

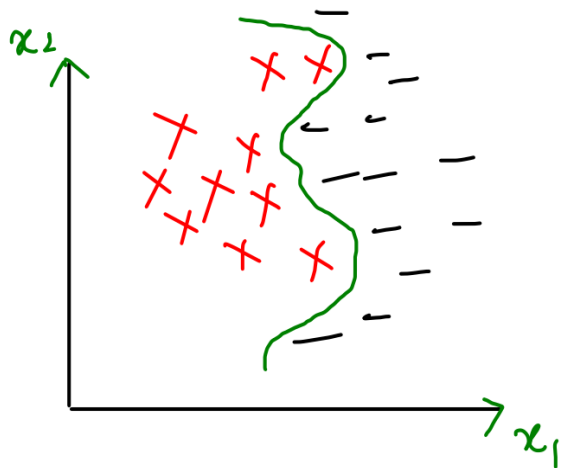
$$x_3 = x^3$$

Rewrite the hypothesis,

$$h(x) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

Can we solve now?!

# Non-linear classification



$$\rightarrow h(x) = g(w_0 + w_1 x_1 + w_2 x_2)$$

$$\rightarrow h(x) = g(w_0 + w_1 x_1^2 + w_2 x_2^2 + w_3 x_1 x_2 + w_4 x_1 x_2^2 + \dots)$$

Introduce more features

$$x'_1 = x_1^2$$

$$x'_2 = x_2^2$$

$$x'_3 = x_1 x_2$$

$\vdots$

so many ....

$$h(x) = g(w_0 + w_1 x'_1 + w_2 x'_2 + w_3 x'_3 + \dots)$$

Can we solve this problem??

# Observations

- Introducing more features (mainly polynomial terms; complex non-linear hypothesis)
- These features help the model to use existing learning algorithm to handle the problem (linear regression/logistic regression)
- Such hand-crafted feature engineering is not desirable option to solve non-linear problems (with more features)

# ANN

- ANN is a very powerful and widely used model to learn a complex non-linear hypothesis
- ANN Representation helps us to automatically extract such features
- Solve XOR problem with ANN topics discussed so far (MP neuron, Perceptron, Sigmoid neuron)



# XOR Problem (using MP neuron)

$x_1$	$x_2$	$y$
0	0	1
0	1	0
1	0	0
1	1	1

$$y = 1, \quad \sum x_i > \theta$$

$$y = 0, \quad \sum x_i < \theta$$

$$\rightarrow 0 > \theta$$

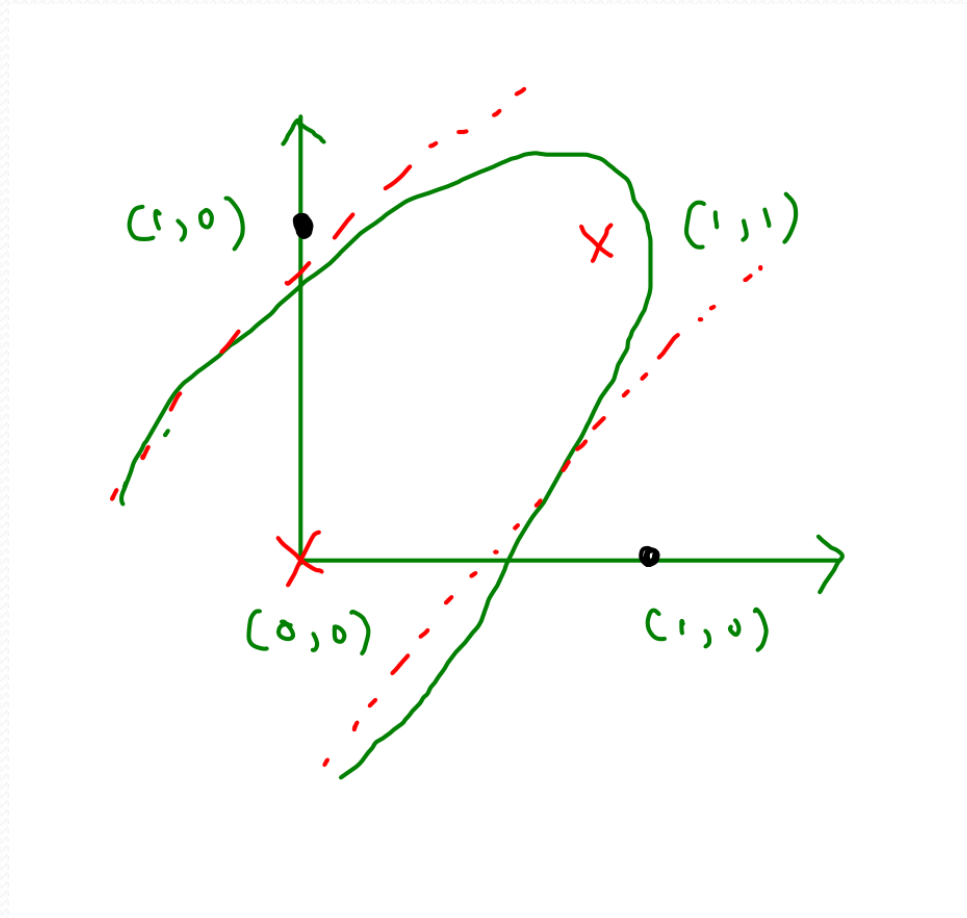
$$\rightarrow 1 < \theta$$

$$\rightarrow 1 < \theta$$

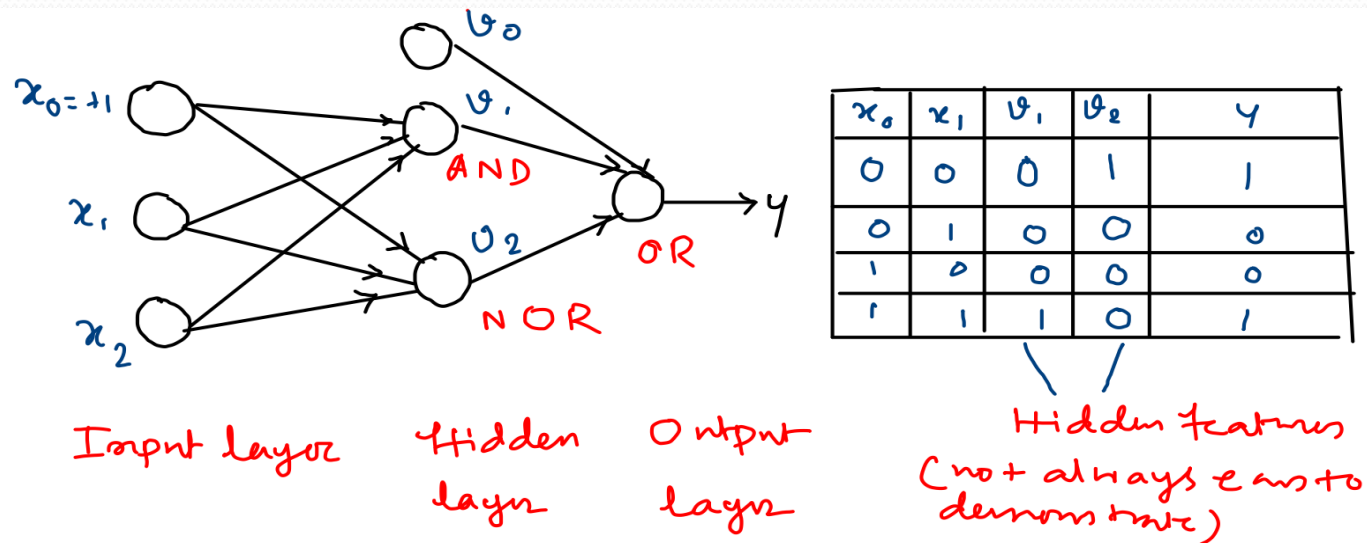
$$\rightarrow 2 > \theta$$

Can we find  $\theta$ ??

# Geometrical representation for XOP Problem



# Introduce multiple layers in ANN



XOR  $\rightarrow$  Simple problem (Solve by two lines)  
we can solve it by dividing it into  
two linear problem (AND and NOR,  
then OR)

# Multilayer Perceptron (MLP)

- It has one input layer, one output layer, and multiple hidden layers (feed-forward architecture).
- Number of input neurons= number of features+1
- Number of output neurons=number of classes
- Number of hidden layers and hidden neurons in each layer are fixed experimentally (application depended)
- We consider, number of hidden layers=1
- Activation function has been considered in each hidden neuron and output neuron.

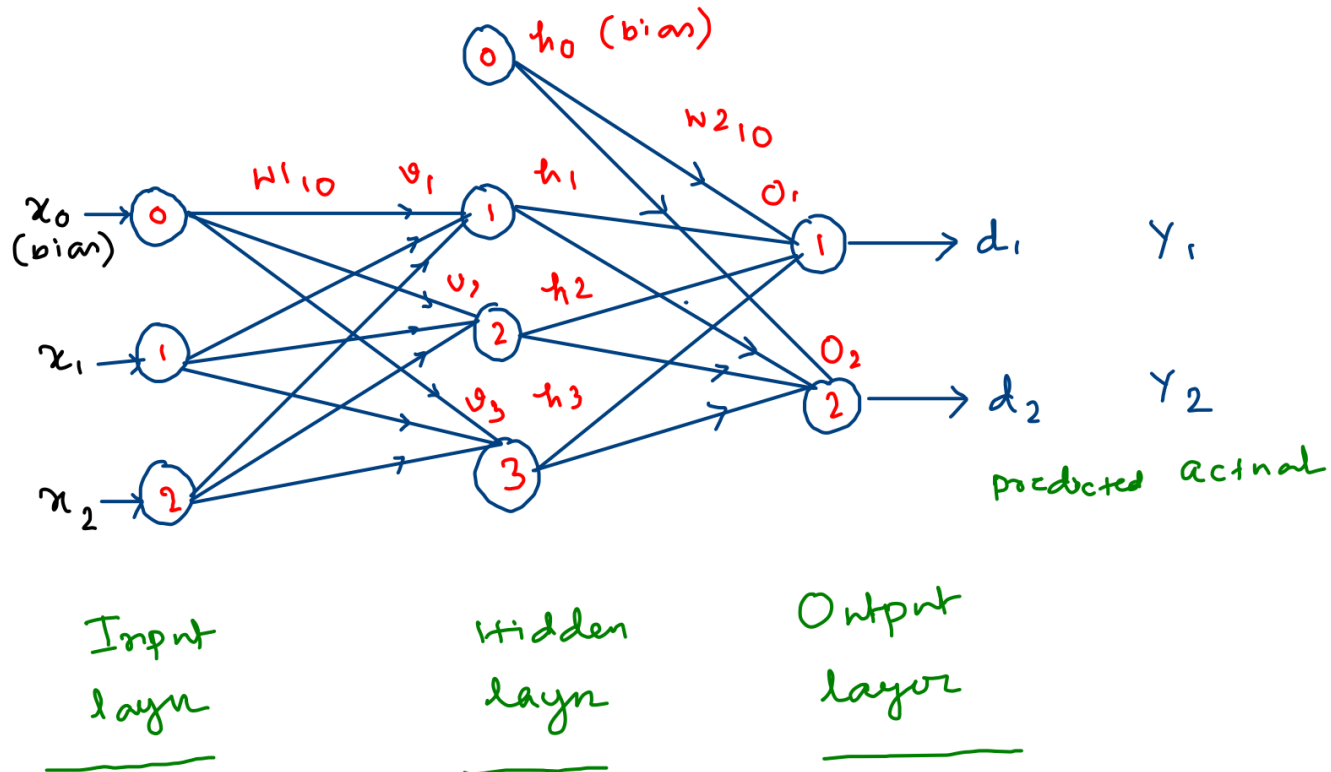
# Neural representation for non-linear hypothesis

- MLPs can be trained to implement any given nonlinear input-output mapping (hypothesis).
- Here, multiple hidden neurons (with non-linear activation function) are helpful to extract hidden features automatically.
- MLP is basically transformation of the original input space into a new one, where the classification task becomes linearly separable.
- In this regard, Cover's theorem may be referred:  
A complex pattern-classification problem, cast in a high-dimensional space nonlinearly, is more likely to be linearly separable than in a low-dimensional space.

# Training Process of MLP

- Design a MLP to solve a non-linearly separable problem (2 features and two classes)
  - Data preprocessing (one hot encoding)
  - Architecture
  - Training algorithm (?)
  - Cost function: MSE
  - Activation function: Sigmoid/logistic function

# Architecture



Start training (using incremental GD)

Input to hidden layer

Random initialized the  
weights

$$v_1 = w_{1,0} x_0 + w_{1,1} x_1 + w_{1,2} x_2$$

$$v_2 = \text{same} \dots$$

$$v_3 = \text{same} \dots$$

$$h_1 = g(v_1) = \frac{1}{1 + \exp(-v_1)}$$

$$h_2 = \text{same} \dots$$

$$h_3 = \text{same} \dots$$



Hidden to output

$$O_1 = W_{2,0} h_0 + W_{2,1} h_1 + W_{2,2} h_2 + W_{2,3} h_3$$

$$O_2 = \text{Same} \dots$$

$$d_1 = g(O_1)$$

$$d_2 = g(O_2)$$

Training GD

$$E_{\text{tr}} = \frac{1}{2} \sum_{j=1}^2 (y_j - d_j)^2$$

Error in each output neuron.

$$e_1 = y_1 - d_1$$

$$e_2 = y_2 - d_2$$

$$w_{ji} := w_{ji} + \alpha (y_j - d_j) d_j (1 - d_j) x_i$$

Update the weight between hidden to output:

$$w_{2,0} = w_{2,0} + \alpha (y_1 - d_1) d_1 (1 - d_1) x_0$$

⋮

Same for others

Update weight between input to hidden

---

$$w_{10} = w_{10} + \alpha \underbrace{(\text{?} - h_1)}_{\text{error (not known)}} h_1 (1 - h_1) x_0$$

**Solution:** back-propagating errors for  
Output neurons to hidden neurons

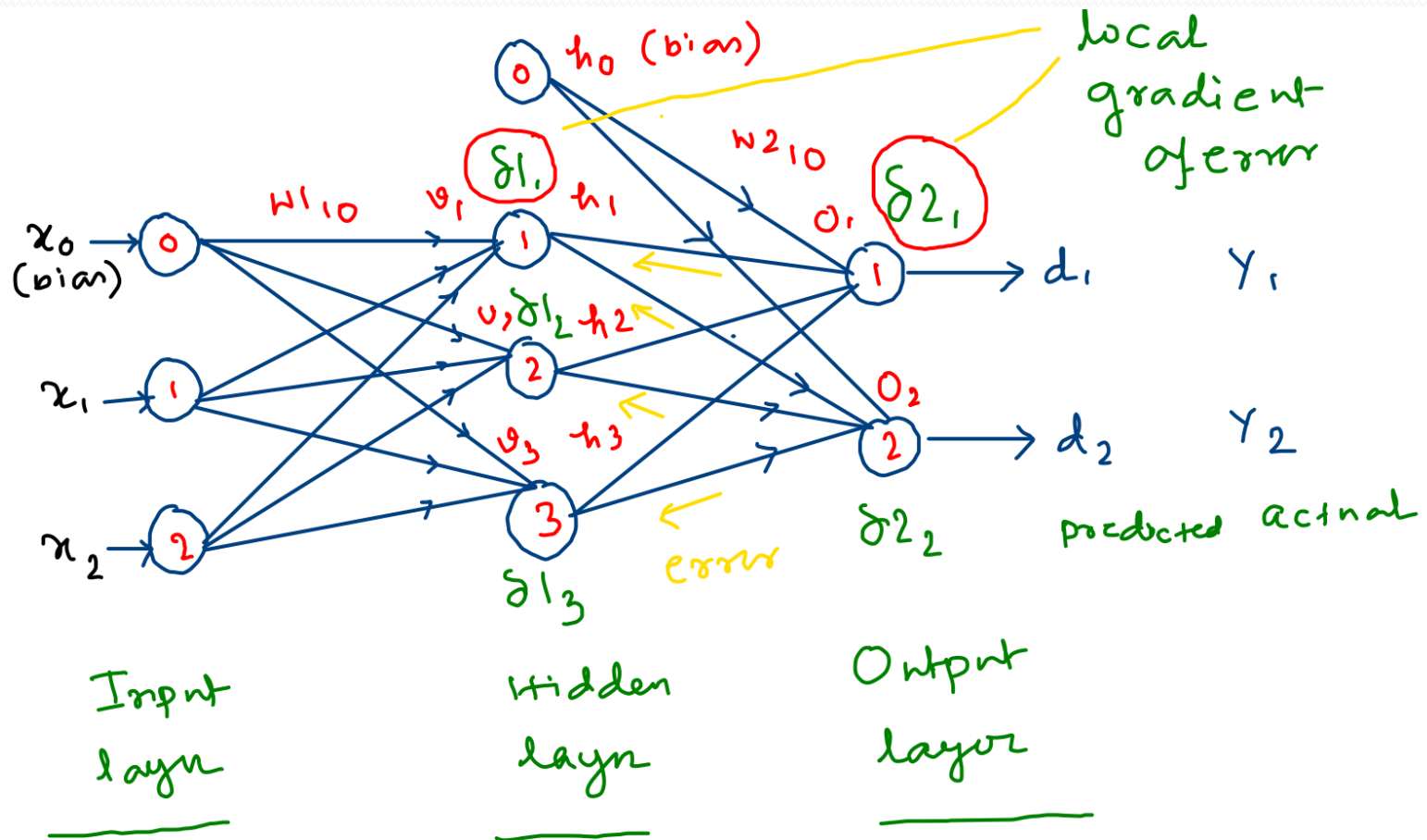
Back-propagation learning algorithm

---

Step 1: Forward pass

Step 2: Backward pass

# Back-propagation algorithm



# Back-propagation algorithm

Forward Pass (update for one sample)

Input to hidden layer

$$v_1 = w_{1,0} x_0 + w_{1,1} x_1 + w_{1,2} x_2$$

$$v_2 = \text{same} \dots$$

$$v_3 = \text{same} \dots$$

$$h_1 = g(v_1) = \frac{1}{1 + \exp(-v_1)}$$

$$h_2 = \text{same} \dots$$

$$h_3 = \text{same} \dots$$

Hidden to output

$$O_1 = w_{2,0} h_0 + w_{2,1} h_1 + w_{2,2} h_2 + w_{2,3} h_3$$

$$O_2 = \text{Same} \dots$$

$$d_1 = g(O_1)$$

$$d_2 = g(O_2)$$

## Backward pass

Error calculation for the pattern

$$E_{\text{tr}} = \frac{1}{2} \sum_{j=1}^2 (d_j - y_j)^2$$

weight updation in GD

$$w_{ji} := w_{ji} + \alpha \underbrace{(y_j - d_j)}_{\text{error}} (1 - d_j) d_j x_i$$

error in each output node

$$e_1 = (y_1 - d_1)$$

$$e_2 = (y_2 - d_2)$$

local gradient of error in each output node

---

$$\delta_{2,1} = d_1(1-d_1)(y_1-d_1)$$

$$\delta_{2,2} = d_2(1-d_2)(y_2-d_2)$$

local gradient of error in each hidden node

---

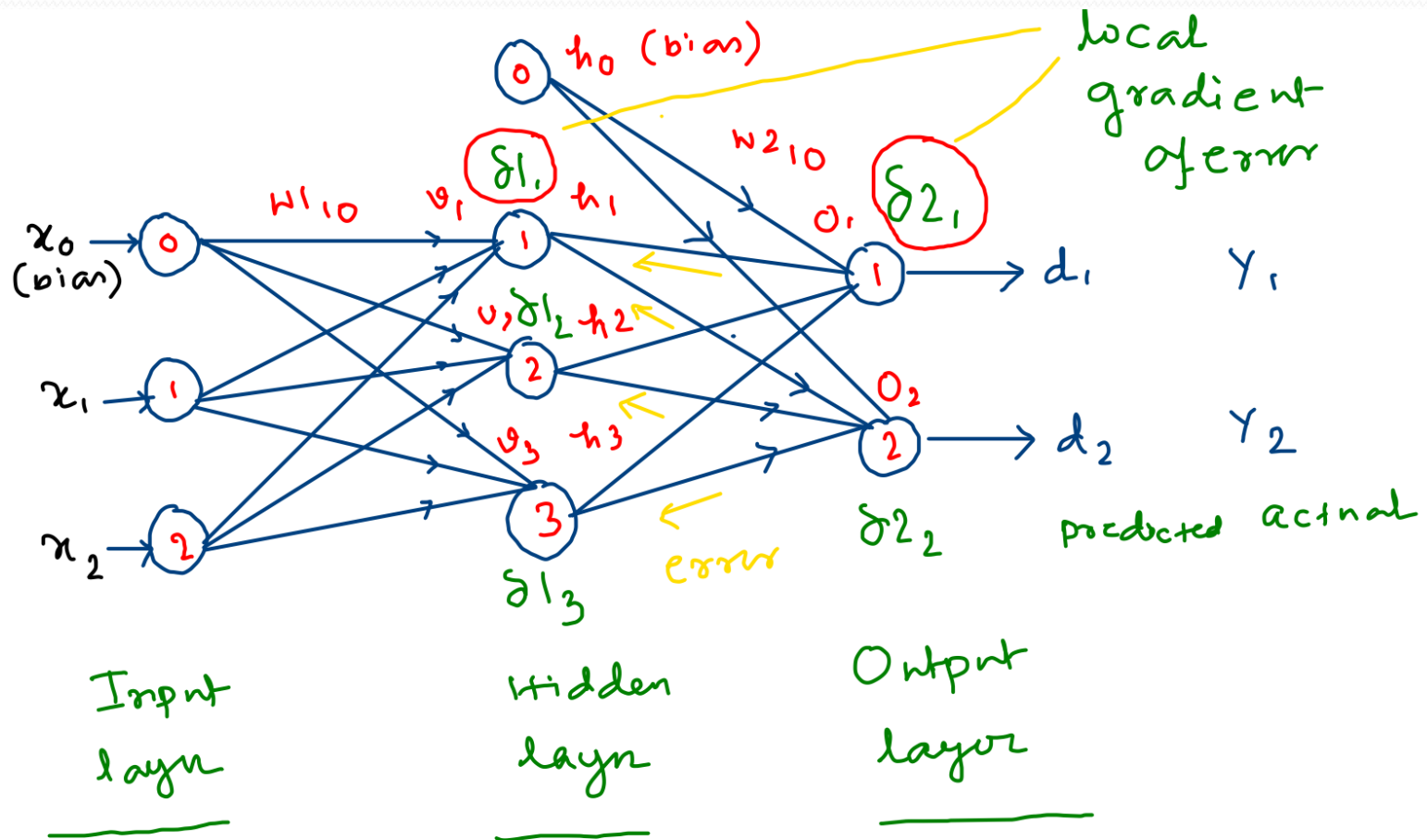
$$\delta_{1,1} = (\delta_{2,1} * w_{2,1,1} + \delta_{2,2} * w_{2,2,1}) h_1(1-h_1)$$

$$\delta_{1,2} = (\delta_{2,1} * w_{2,1,2} + \delta_{2,2} * w_{2,2,2}) h_2(1-h_2)$$

$$\delta_{1,3} = ?$$



# Back-propagation algorithm



## Update of weights

$$w_{ji} := w_{ji} + \alpha (y_j - d_j) d_j (1 - d_j) x_i$$

hidden to output

$$w_{2,0} := w_{2,0} + \alpha * \delta_{2,1} * h_0$$

$$w_{2,1} := w_{2,1} + \alpha * \delta_{2,1} * h_1$$

⋮

Input to hidden

$$w_{10} := w_{10} + \alpha * \delta_1 * x_0$$

$$w_{11} := w_{11} + \alpha * \delta_1 * x_1$$

⋮

Stopping criteria

$$\text{Epoch 1: } E_{itr} = \sum_{i=1}^{\text{no. of samples}} \sum_{j=1}^{\text{no. of classes}} (y_j - d_j)^2$$

Update weights for each sample

$$\text{Epoch 2: } E_{itr+1} = \sum_{i=1}^m \sum_{j=1}^2 (y_j - d_j)^2$$

Until

$$|E_{itr-1} - E_{itr}| > \text{small value or}$$

no. of epochs < large value