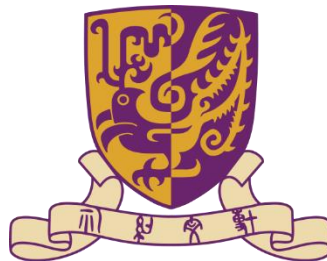


# Neural Network Fundamental

# Single Layer Perceptron

Piji Li

[pjli@se.cuhk.edu.hk](mailto:pjli@se.cuhk.edu.hk)



# Outline

---

- Introduction
- Structure
- Learning
- Tricks

# Outline

---

- **Introduction**
- Structure
- Learning
- Tricks

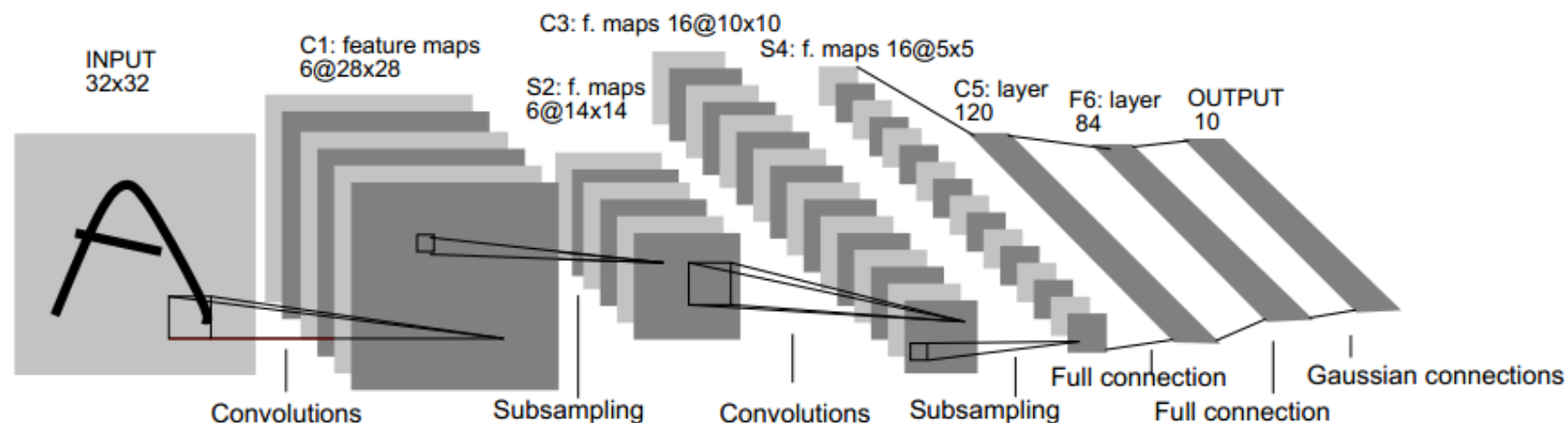
# Introduction

- **Deep Learning**

- Nowadays, deep learning techniques have been applied to fields like **computer vision, speech recognition, natural language processing and bioinformatics** where they have been shown to produce **state-of-the-art** results on various tasks.

- **Neural Network**

- Deep neural networks, convolutional deep neural networks, deep recurrent neural networks.



# Outline

---

- Introduction
- **Structure**
- Learning
- Tricks

# Structure-single layer perceptron

- Single hidden layer FF and BP
- Regression( $K = 1$ )
- Classification( $K$ -class)
- Linear combinations
- Activation Function  $\sigma(v)$

$$Z_m = \sigma(\alpha_{0m} + \alpha_m^T X), \quad m = 1, \dots, M$$

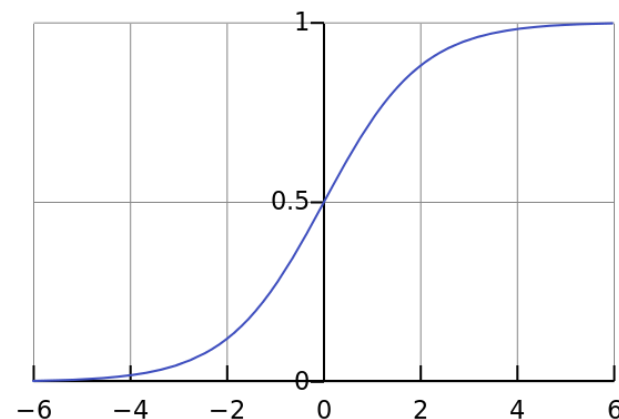
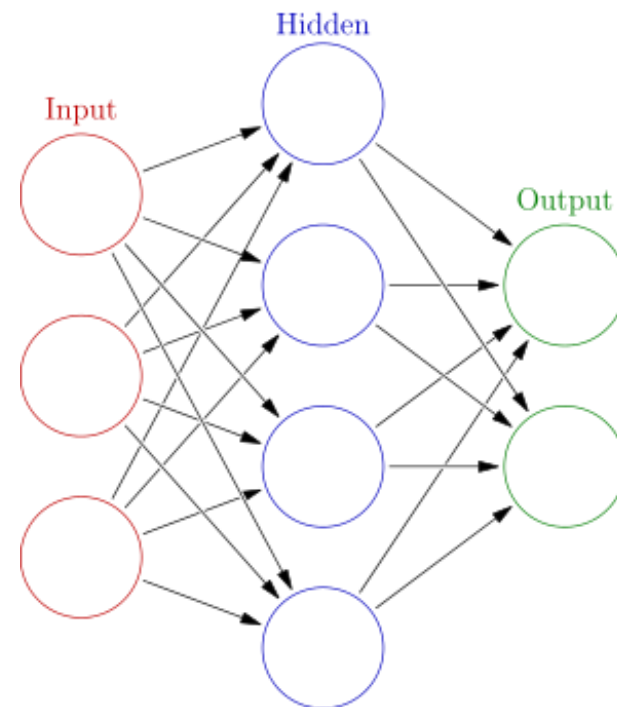
$$T_k = \beta_{0k} + \beta_k^T Z, \quad k = 1, \dots, K$$

$$f_k(X) = g_k(T), \quad k = 1, \dots, K$$

$$\sigma(v) = 1/(1 + e^{-v})$$

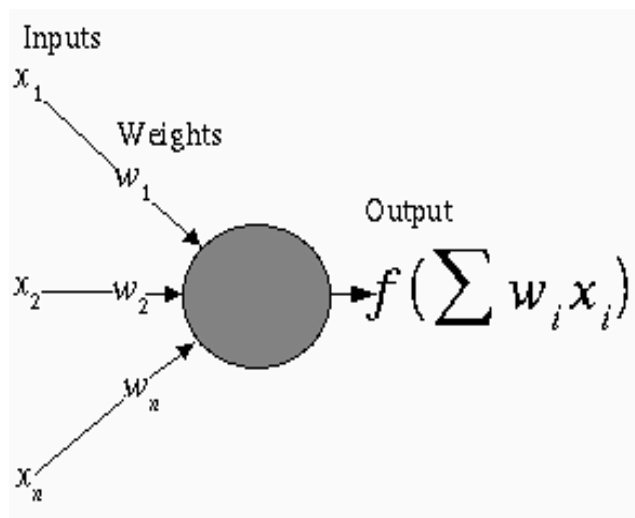
$$g_k(T) = T_k$$

$$g_k(T) = \frac{e^{T_k}}{\sum_{\ell=1}^K e^{T_\ell}} \quad \text{Softmax function}$$



# Non-linear Model

- Notice that if  $\sigma$  is the identity function, then the entire model collapses to a linear model in the inputs.
- Hence a neural network can be thought of as a nonlinear generalization of the linear model, both for regression and classification.
- By introducing the nonlinear transformation , it greatly enlarges the class of linear models.



# Outline

---

- Introduction
- Structure
- **Learning**
- Tricks



# Parameters

$\{\alpha_{0m}, \alpha_m; m = 1, 2, \dots, M\}$   $M(p + 1)$  weights

$\{\beta_{0k}, \beta_k; k = 1, 2, \dots, K\}$   $K(M + 1)$  weights

**Loss functions:**

$$R(\theta) = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2$$

Sum-of-squared errors

$$R(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log f_k(x_i)$$

Cross-entropy

Minimize by GD

# Back-Propagation(BP)

1986



D.E. Rumelhart, G.E. Hinton, R.J. Williams  
**Learning representation by back-propagating errors.** *Nature*, 323 (1986), pp. 533–536

- ❑ Solved learning problem
- ❑ Biological system
- ❑ ...
- Hard to train (non-convex, tricks)
- Hard to do theoretical analysis
- Small training sets ...

# Back-Propagation(BP)

Here is back-propagation in detail for squared error loss:

$$\begin{aligned} R(\theta) &\equiv \sum_{i=1}^N R_i \\ &= \sum_{i=1}^N \sum_{k=1}^K (y_{ik} - f_k(x_i))^2 \end{aligned} \quad \begin{aligned} Z_m &= \sigma(\alpha_{0m} + \alpha_m^T X), \quad m = 1, \dots, M \\ T_k &= \beta_{0k} + \beta_k^T Z, \quad k = 1, \dots, K \\ f_k(X) &= g_k(T), \quad k = 1, \dots, K \end{aligned}$$

Take derivatives:

$$\begin{aligned} \frac{\partial R_i}{\partial \beta_{km}} &= -2(y_{ik} - f_k(x_i))g'_k(\beta_k^T z_i)z_{mi}, \\ \frac{\partial R_i}{\partial \alpha_{m\ell}} &= -\sum_{k=1}^K 2(y_{ik} - f_k(x_i))g'_k(\beta_k^T z_i)\beta_{km}\sigma'(\alpha_m^T x_i)x_{i\ell}. \end{aligned}$$

Given these derivatives, a gradient descent update method:

$$\begin{aligned} \beta_{km}^{(r+1)} &= \beta_{km}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}^{(r)}} \\ \alpha_{m\ell}^{(r+1)} &= \alpha_{m\ell}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \alpha_{m\ell}^{(r)}} \end{aligned}$$

# Back-Propagation(BP)

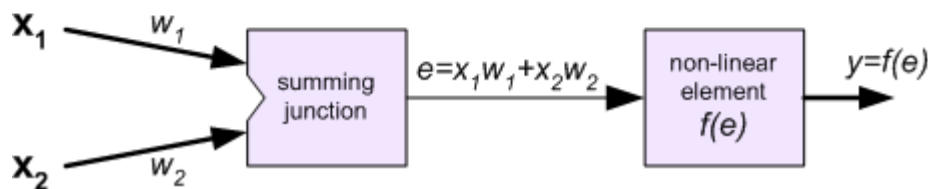
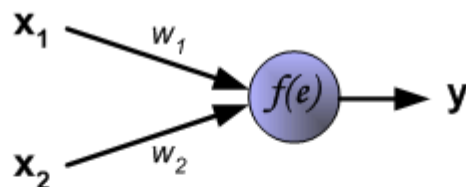
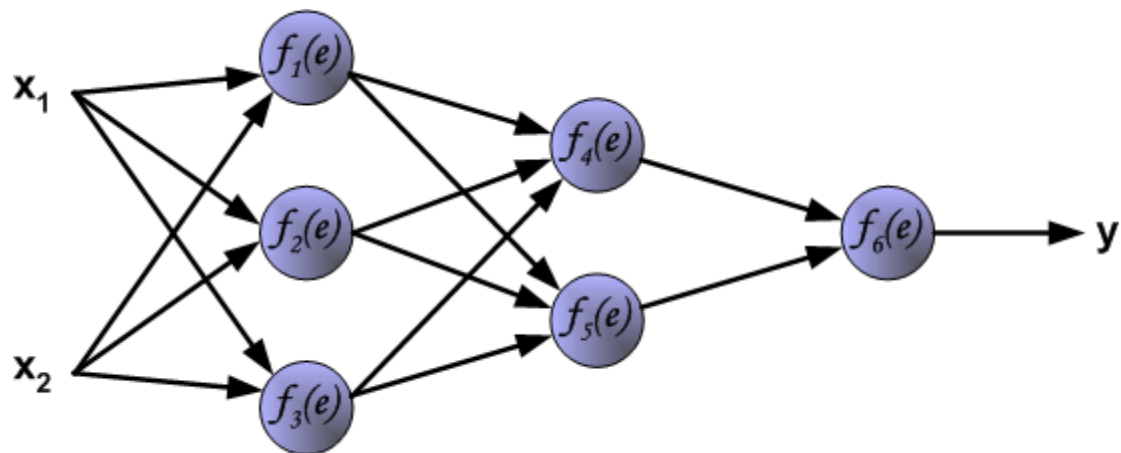
Rewriting:

$$\begin{aligned} \frac{\partial R_i}{\partial \beta_{km}} &= \delta_{ki} z_{mi} \\ \frac{\partial R_i}{\partial \alpha_{ml}} &= s_{mi} x_{il} \end{aligned} \quad \text{errors} \quad \text{and} \quad s_{mi} = \sigma'(\alpha_m^T x_i) \sum_{k=1}^K \beta_{km} \delta_{ki}$$

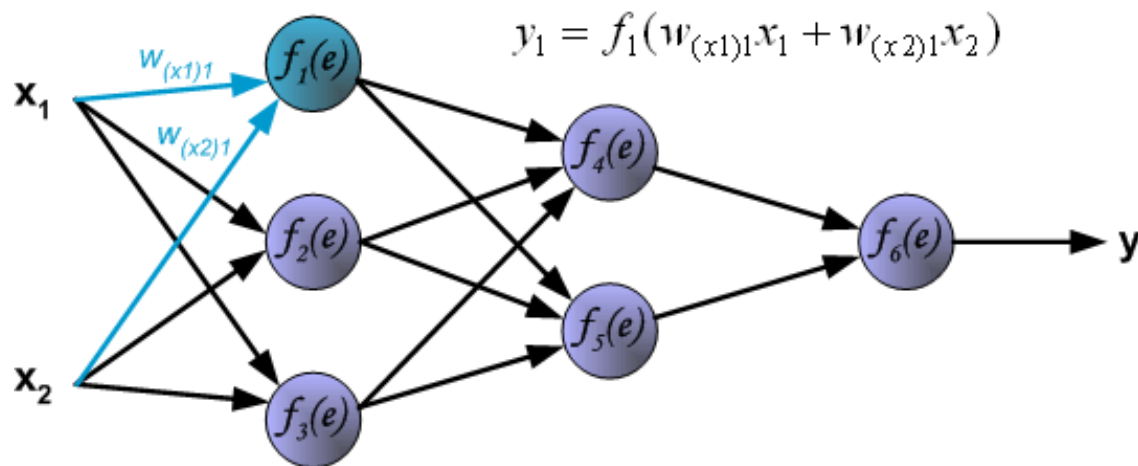
Two-pass algorithm:

- In the forward pass, the current weights are fixed and the predicted values are computed.
- In the backward pass, the output-layer errors are computed, and then back propagated to the hidden-layer errors.

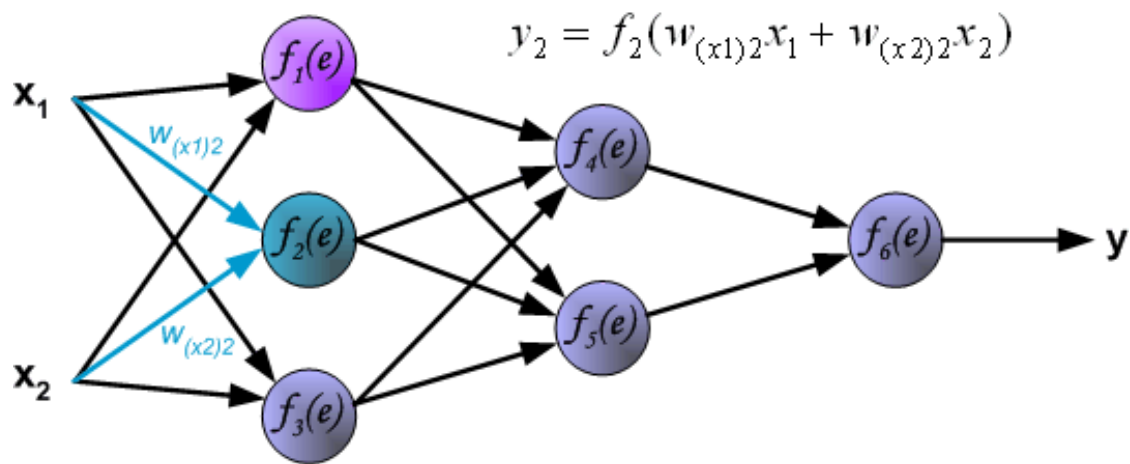
# Back-Propagation(BP)-Examples



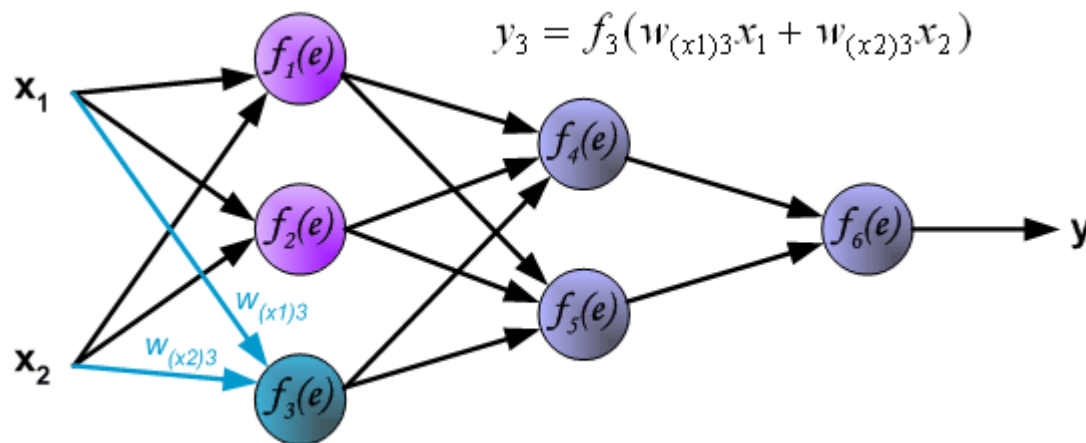
# Back-Propagation(BP)-Examples



# Back-Propagation(BP)-Examples

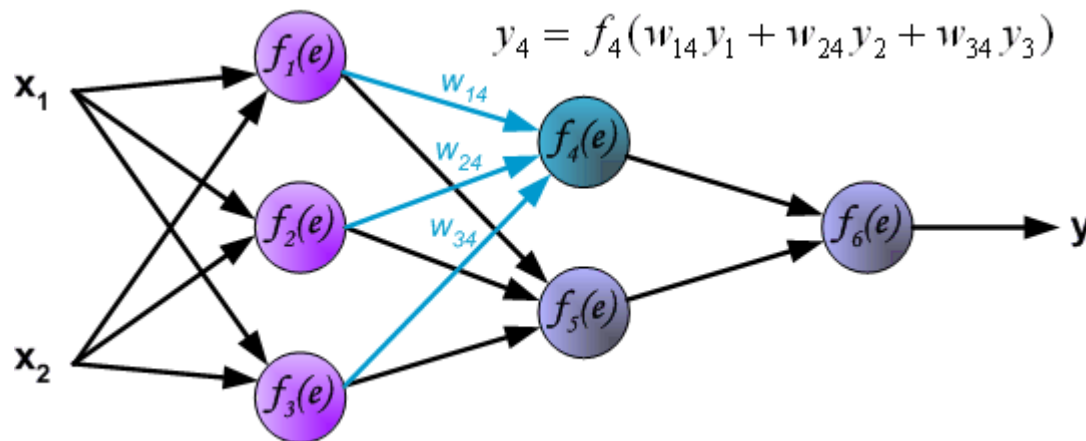


# Back-Propagation(BP)-Examples

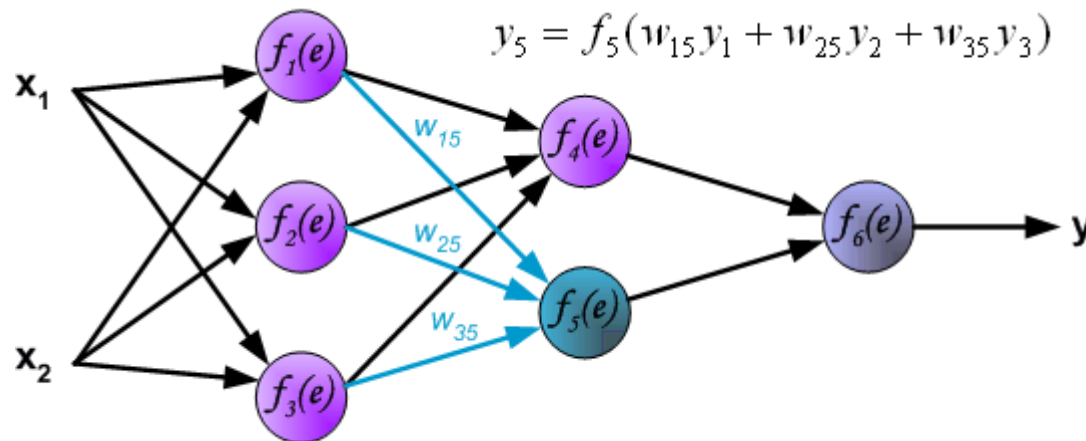




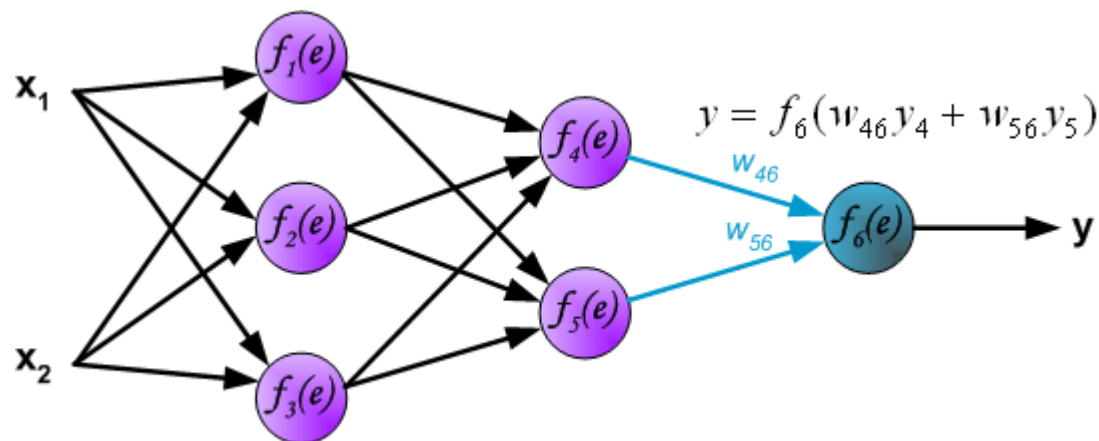
# Back-Propagation(BP)-Examples



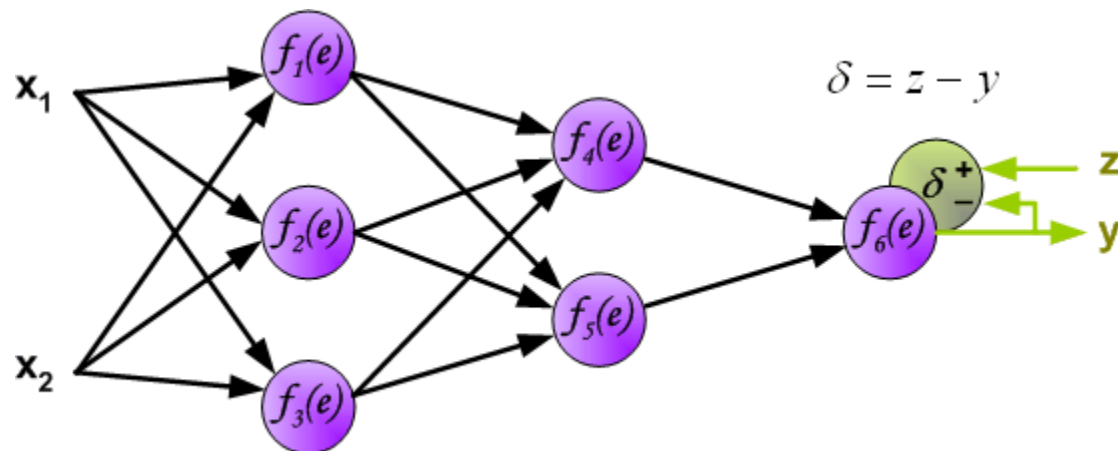
# Back-Propagation(BP)-Examples



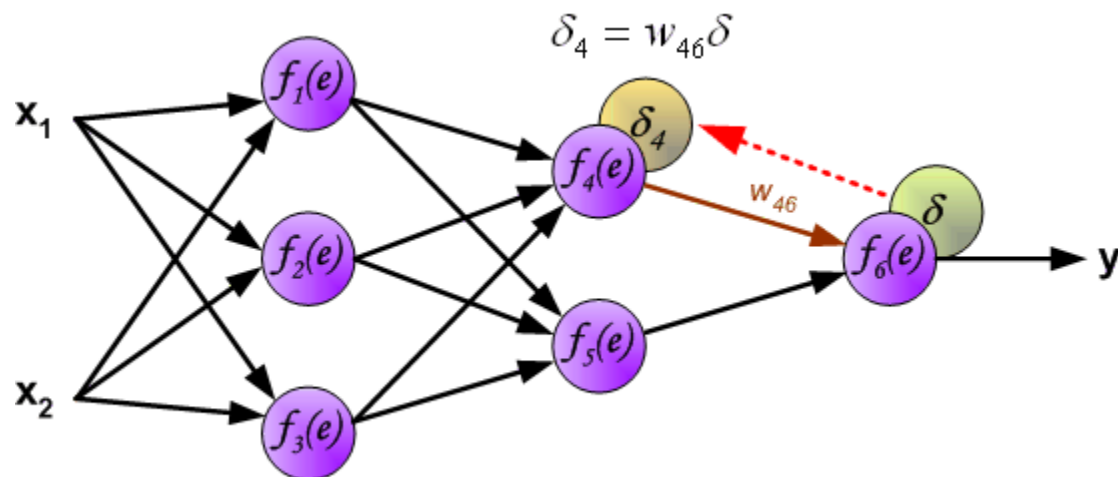
# Back-Propagation(BP)-Examples



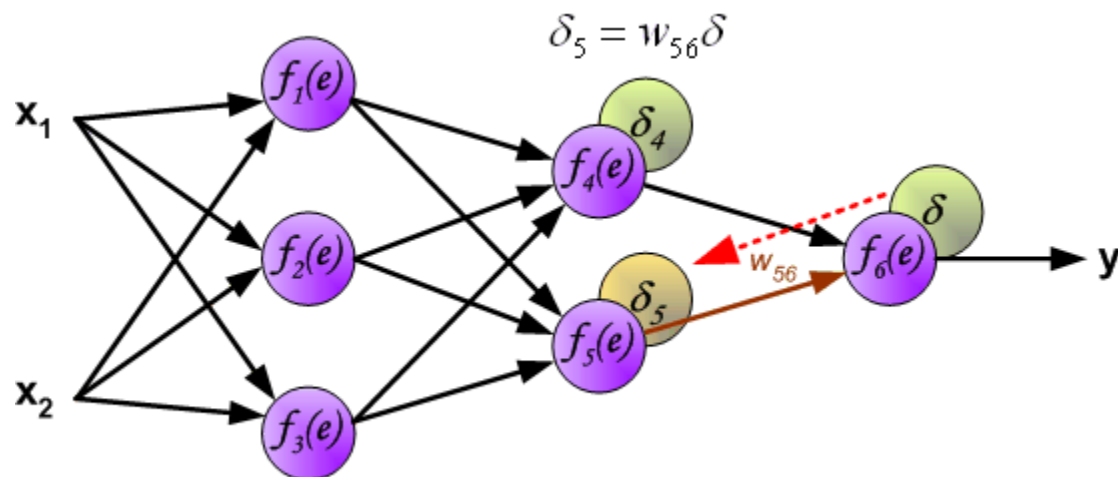
# Back-Propagation(BP)-Examples



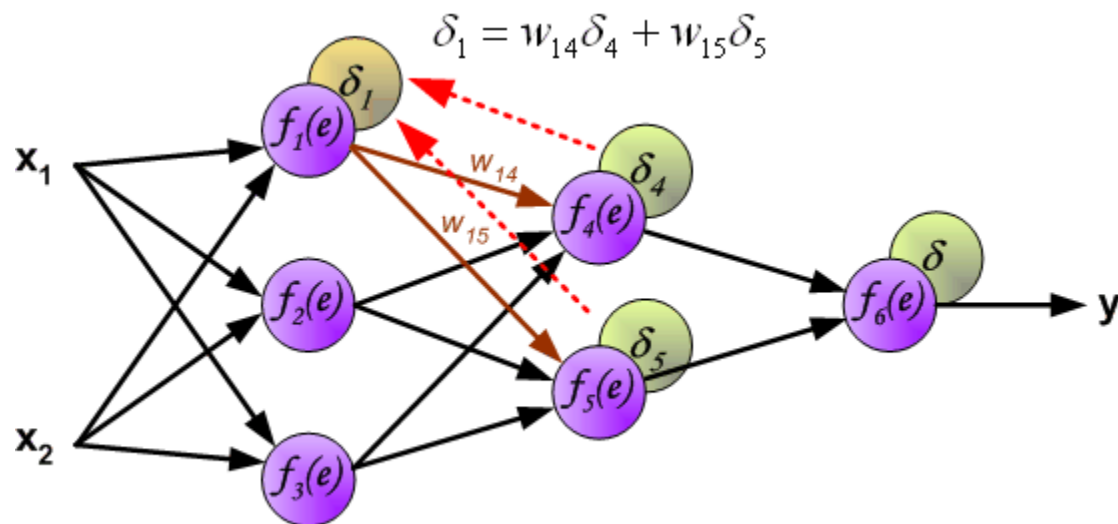
# Back-Propagation(BP)-Examples



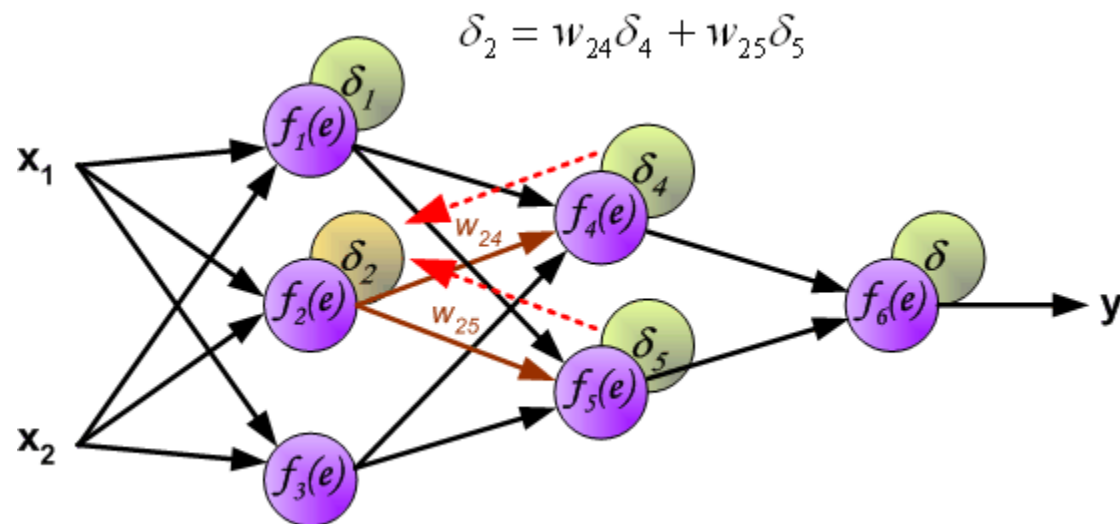
# Back-Propagation(BP)-Examples



# Back-Propagation(BP)-Examples

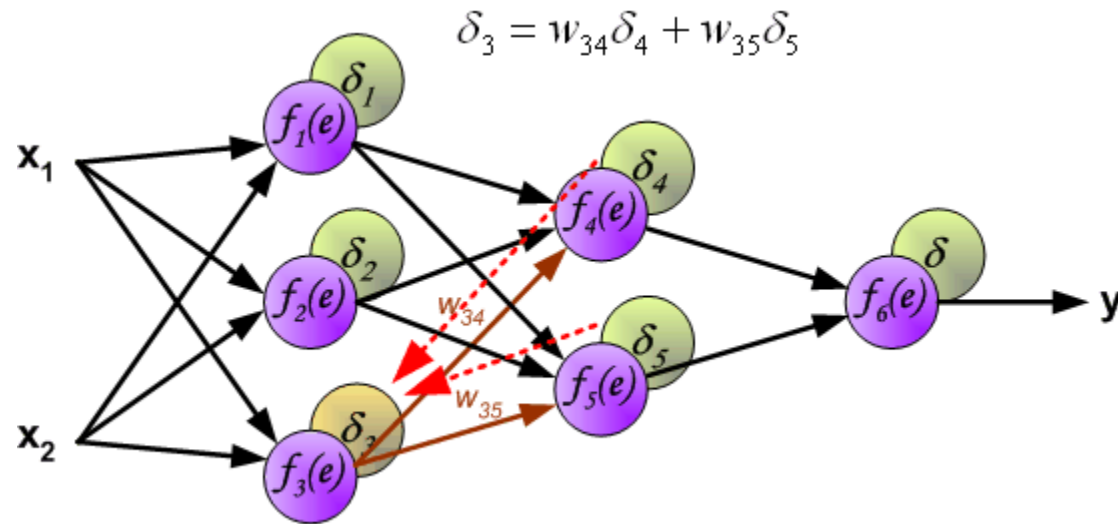


# Back-Propagation(BP)-Examples

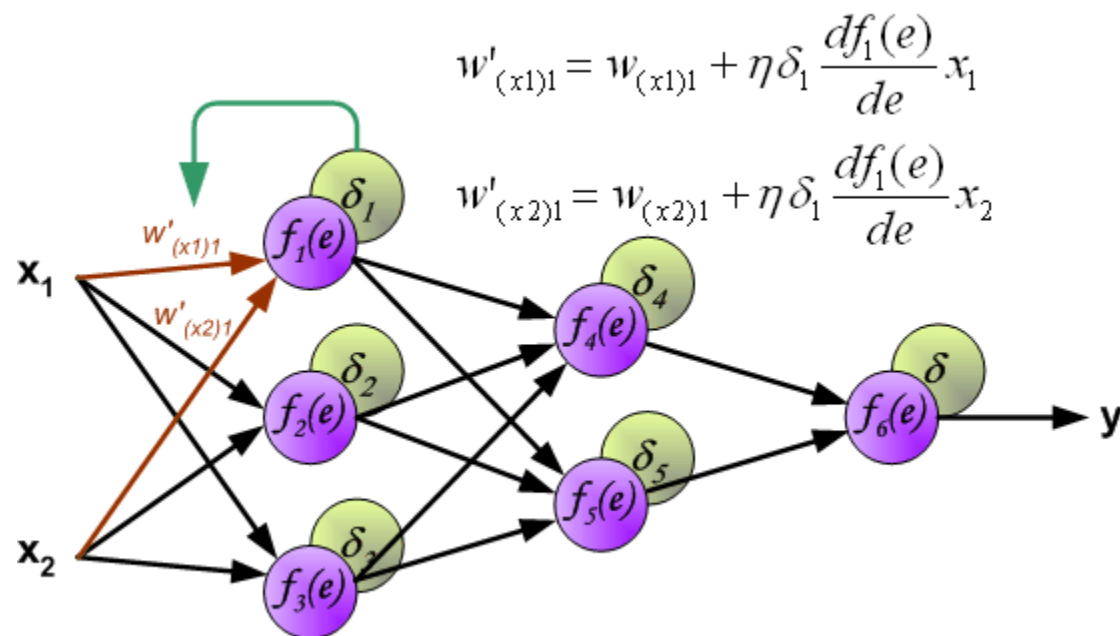




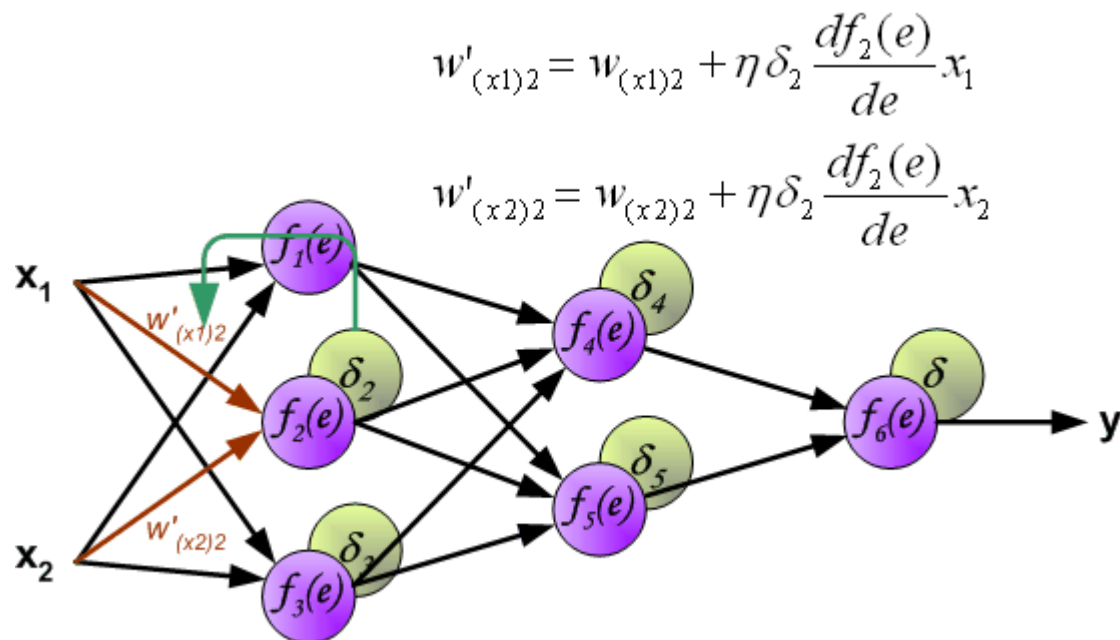
# Back-Propagation(BP)-Examples



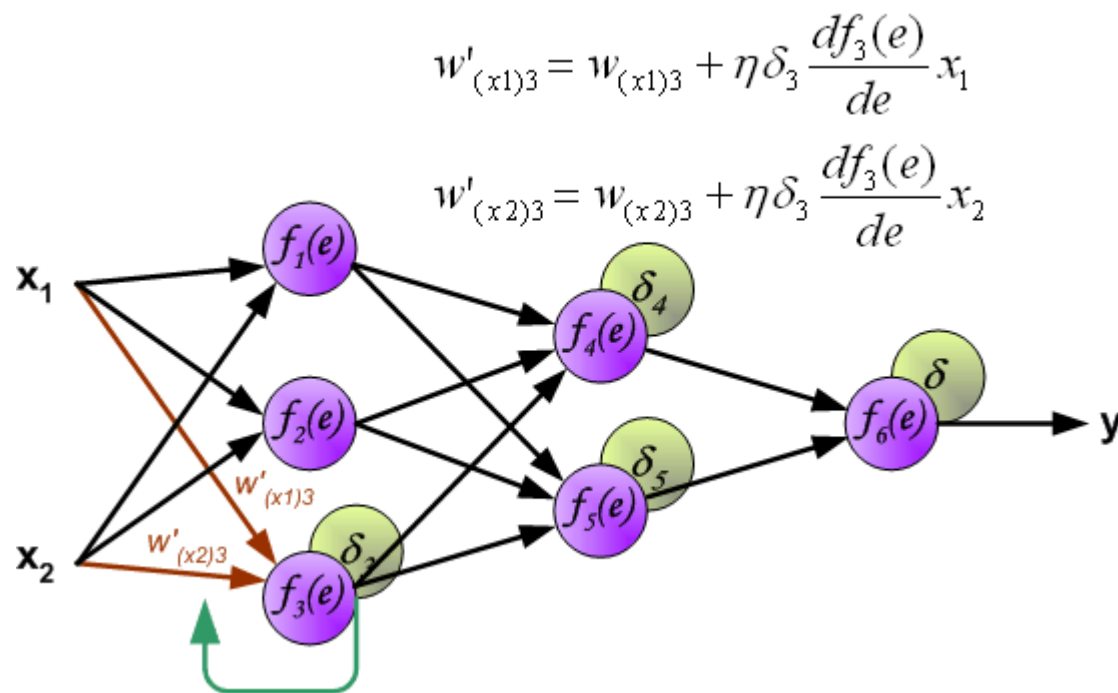
# Back-Propagation(BP)-Examples



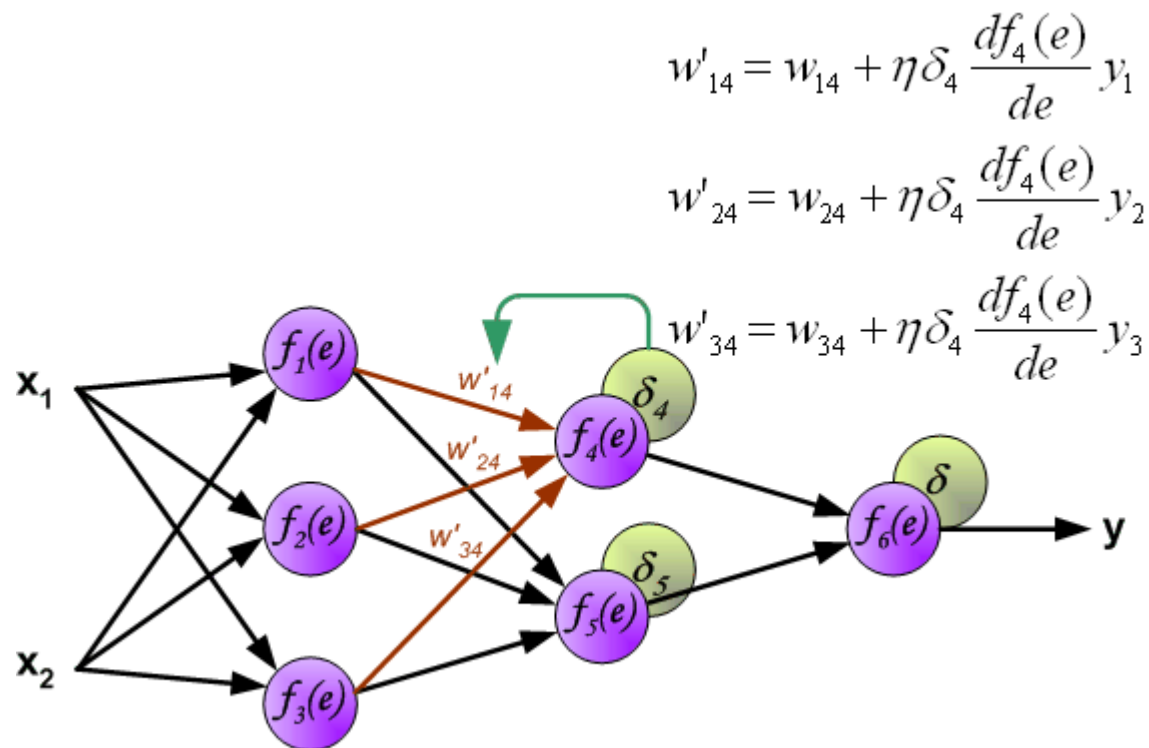
# Back-Propagation(BP)-Examples



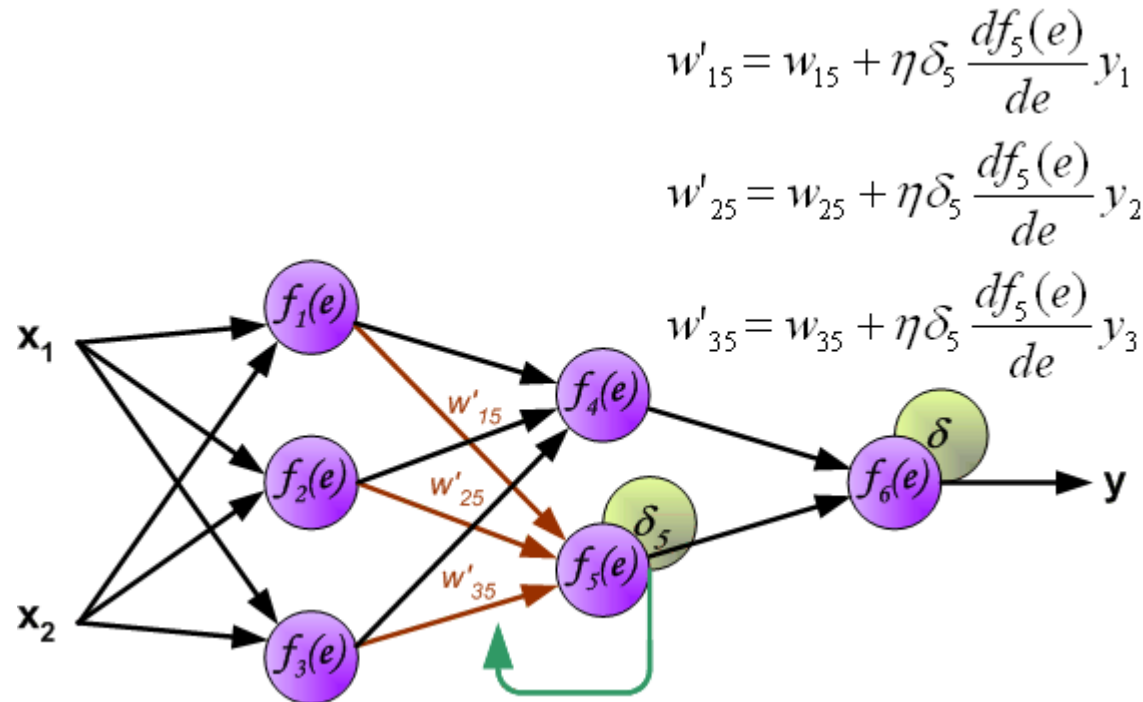
# Back-Propagation(BP)-Examples



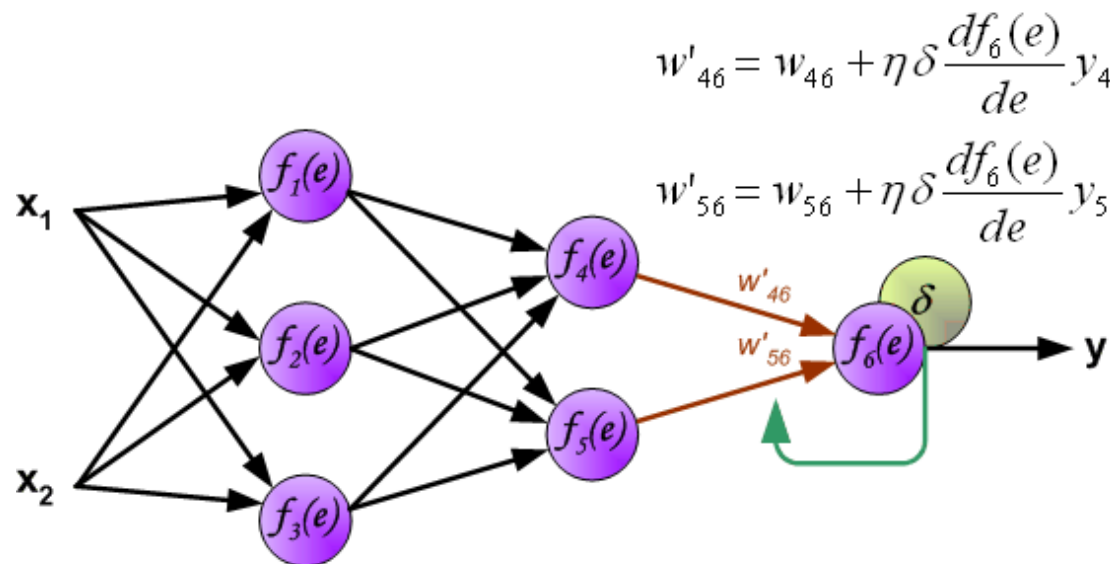
# Back-Propagation(BP)-Examples



# Back-Propagation(BP)-Examples



# Back-Propagation(BP)-Examples



# Outline

---

- Introduction
- Structure
- Learning
- **Tricks**



# Tricks-Initial values

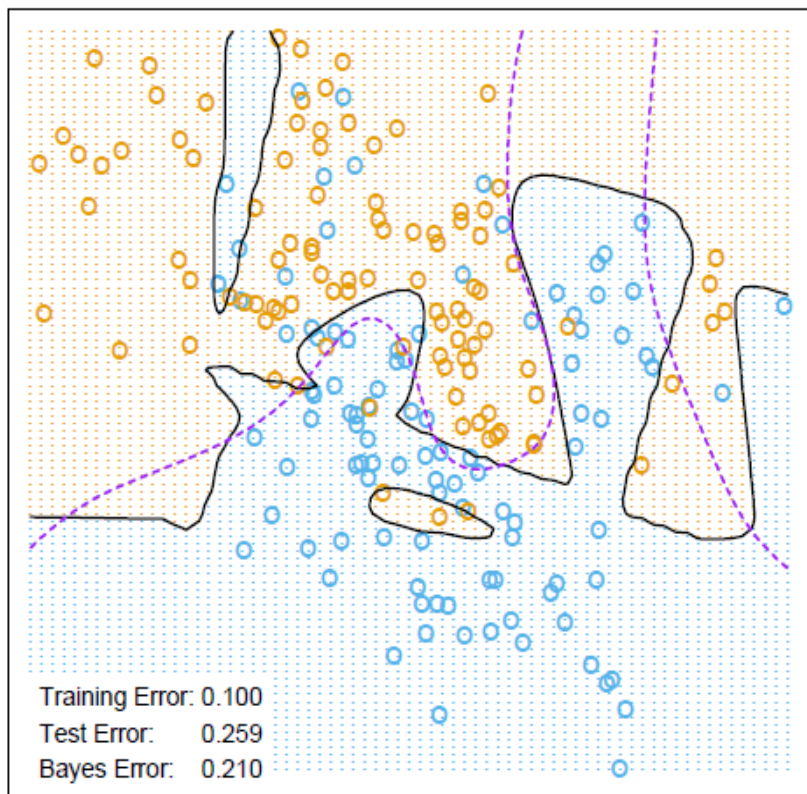
- The weights are typically initialized to small random values chosen from a zero-mean Gaussian with a standard deviation of about 0.01.
- Note that if the weights are near zero, then the operative part of the sigmoid is roughly linear, and hence the neural network collapses into an approximately linear model.
- Use of exact zero weights leads to zero derivatives and perfect symmetry, and the algorithm never moves.
- Starting instead with large weights often leads to poor solutions.

# Tricks-Regularization

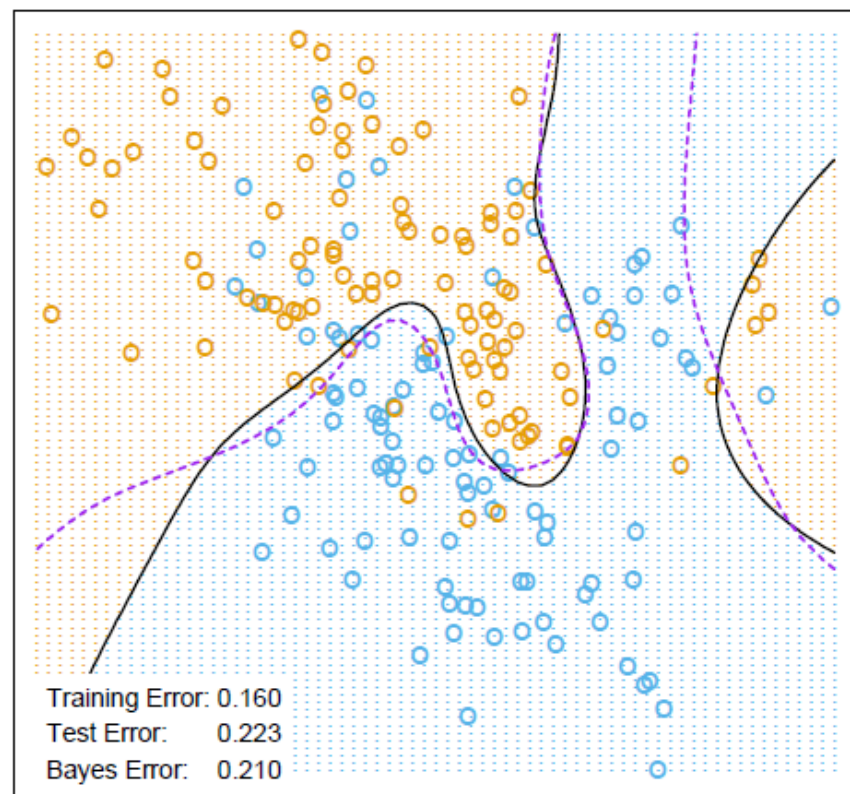
Add a penalty to the error function

$$R(\theta) + \lambda J(\theta)$$

Neural Network - 10 Units, No Weight Decay

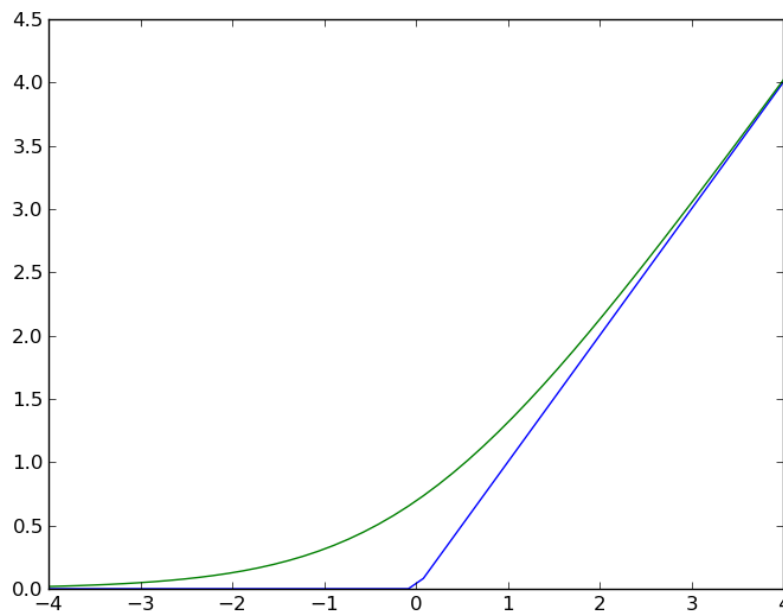


Neural Network - 10 Units, Weight Decay=0.02



# Tricks-Tuning

- Learning Rate
- Mini-batch size for batch learning
- Activation function
  - Rectifier function  $f(x) = \max(0, x)$



---

Thanks a lot!