

# BITS F464: Machine Learning

# 31

# Neural Network

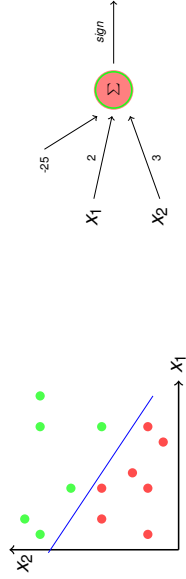


**Dr. Kamlesh Tiwari**  
 Assistant Professor, Department of CSIS,  
 BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA  
 April 09, 2021 **ONLINE** (Campus @ BITS-Pilani Jan-May 2021)

<http://kti.wari.in/ml>

## What about this arrangement?

With chosen *decision boundary*  $2x_1 + 3x_2 - 25 = 0$



- This illustration is called as **perceptron**
- Provides a graphical way to represent the linear boundary
- Values **3, 2, -25** are its parameters or weights

Given a data

"How to find appropriate parameters?" is an important **issue**

## Example

Consider the same data

$x_1$	$x_2$	$y$
1	9	green
10	9	green
4	7	green
4	5	red
5	3	red
8	9	green
4	2	red
2	5	red
7	1	red
2	10	green
8	5	green
1	2	red
8	2	red

$\eta = 0.01$

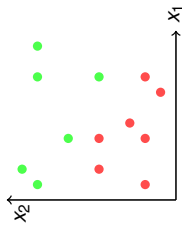
$w_1=0.500$ , $w_1=-0.500$ , $w_2=0.500$	err=7
$w_1=0.380$ , $w_1=-0.120$ , $w_2=0.100$	err=6
$w_1=0.300$ , $w_1=-0.180$ , $w_2=0.060$	err=5
$w_1=0.240$ , $w_1=-0.140$ , $w_2=0.140$	err=4
$w_1=0.180$ , $w_1=-0.200$ , $w_2=0.100$	err=5
$w_1=0.120$ , $w_1=-0.160$ , $w_2=0.190$	err=4
$w_1=0.060$ , $w_1=-0.060$ , $w_2=0.180$	err=5
$w_1=0.020$ , $w_1=-0.120$ , $w_2=0.140$	err=4
$w_1=0.000$ , $w_1=-0.140$ , $w_2=0.180$	err=4
$w_1=0.140$ , $w_1=-0.040$ , $w_2=0.180$	err=5
$w_1=0.200$ , $w_1=-0.100$ , $w_2=0.140$	err=3
$w_1=0.260$ , $w_1=-0.160$ , $w_2=0.100$	err=4
$w_1=0.320$ , $w_1=-0.120$ , $w_2=0.180$	err=3
$w_1=0.360$ , $w_1=-0.020$ , $w_2=0.180$	err=2
$w_1=0.420$ , $w_1=-0.060$ , $w_2=0.140$	err=2
$w_1=0.420$ , $w_1=-0.060$ , $w_2=0.240$	err=2
$w_1=0.900$ , $w_1=-0.020$ , $w_2=0.180$	err=1
$w_1=0.900$ , $w_1=-0.020$ , $w_2=0.240$	err=2
$w_1=0.920$ , $w_1=0.020$ , $w_2=0.220$	err=2
$w_1=0.960$ , $w_1=0.020$ , $w_2=0.220$	err=3
$w_1=0.980$ , $w_1=0.020$ , $w_2=0.200$	err=2
$w_1=1.000$ , $w_1=0.060$ , $w_2=0.180$	err=2
$w_1=1.040$ , $w_1=0.020$ , $w_2=0.180$	err=0

## Linear Classification

Consider Following data

$x_1$	$x_2$	$y$
1	9	green
10	9	green
4	7	green
4	5	red
5	3	red
8	9	green
4	2	red
2	5	red
7	1	red
2	10	green
8	5	green
1	2	red
8	2	red

Data is in 2D, so let us visualize



- Data looks **linearly separable**
- What is the decision boundary?

Many Possibilities, such as  
 if  $(2x_1 + 3x_2 - 25 > 0)$  it is **green**  
 otherwise **red**

## Perceptron Training Rule

Different algorithms may converge to different acceptable hypotheses

### Algorithm 1: Perceptron training rule

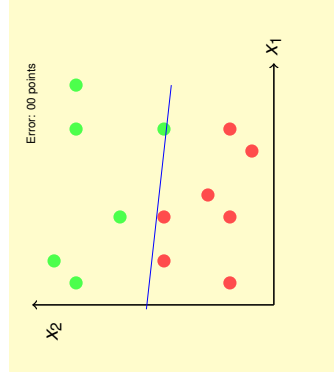
- 1 Begin with **random** weights  $w$
- 2 **repeat**
- 3     **for each misclassified example do**
- 4          $w_i = w_i + \eta(t - o)x_i$
- 5     **until all training examples are correctly classified;**
- 6 **return**  $w$

### Why would this strategy converge?

- Weight does not change when classification is correct
- If perceptron outputs -1 when target is +1: weight increases  $\uparrow$
- If perceptron outputs +1 when target is -1: weight decreases  $\downarrow$

Conversion with perceptron training rule is subject to linear separability of training example and appropriate  $\eta$

## Visual Interpretation

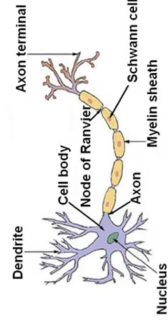


- Conversion is not gradual. (**Error is NOT reducing monotonically**)
- It is **difficult to decide when to stop** if data is not linearly separable

## Neural Network (NN)

NN is biologically motivated learning **model** that mimic human brain

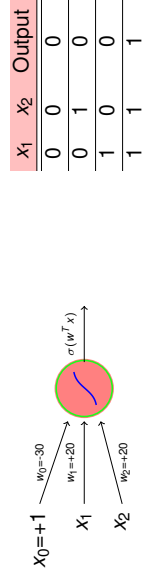
- Started by *W. McCulloch* study on working of neurons in 1943
- MADALINE (1959), an adaptive filter that eliminates echoes on phone lines was the first neural network
- Popularity of Neural Network diminished in 90's but, due to advances in **processing power** and availability of **large data** it again became state-of-the-art



- Cell, Axon, Synapses, Molecules, and Dendrites
- Humans have  $10^{11}$  neurons, each connected to  $10^4$  others, switches in  $10^{-3}$  sec

## An Example

Consider a perceptron with output 0/1 as below



This perceptron computes **logical AND**

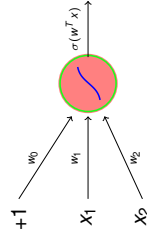
- $w_0 = -10$  gives **logical OR**
- $w_0 = 10, w_1 = -20$  with single input gives **logical NOT**
- XOR is not possible

## An Example

Design a **perceptron** for

$x_1$	$x_2$	Classification
0	0	0
0	1	0
1	0	1
1	1	0

Let us assume following



We have following four equations

- $w_0 + w_1 \times (0) + w_2 \times (0) < 0$
- $w_0 + w_1 \times (0) + w_2 \times (1) < 0$
- $w_0 + w_1 \times (1) + w_2 \times (0) \geq 0$
- $w_0 + w_1 \times (1) + w_2 \times (1) < 0$

By (1)  $w_0 < 0$  so let  $w_0 = -1$

By (2)  $w_0 + w_2 < 0$  so let  $w_2 = -1$

By (3)  $w_0 + w_1 \geq 0$  so let  $w_1 = 1.5$

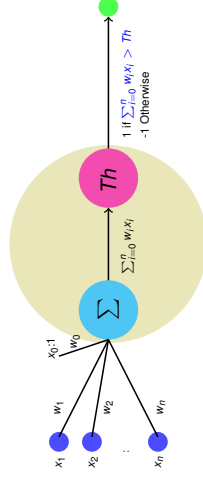
By (4)  $w_0 + w_1 + w_2 < 0$  that is valid

So  $(w_0, w_1, w_2) = (-1, 1.5, -1)$

Other possibilities are also there

## A Single Perceptron

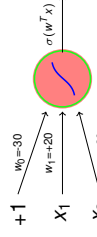
Perceptron representation



- A single perceptron can represent many boolean functions
- Any  $m$ -of- $n$  function (at least  $m$  of the  $n$  inputs must be true) can be represented by perceptron. OR ( $m=1$ ) and AND ( $m=n$ )

**Two layer NN** can represent **any boolean function** (Consider SOP)

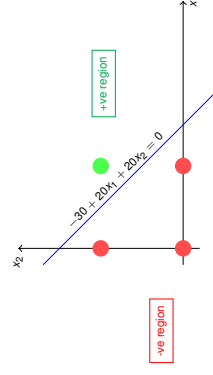
## Essentially it Represents A Decision Boundary



Provides **positive** classification if

$$-30 + 20x_1 + 20x_2 \geq 0$$

Represents a linear decision boundary

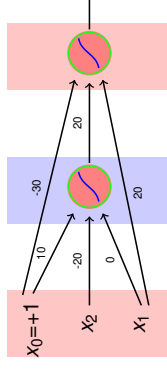


## An Example

Design a **neural network** for

$x_1$	$x_2$	Classification
0	0	0
0	1	0
1	0	1
1	1	0

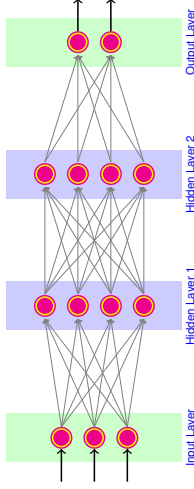
$x_1$	$x_2$	$(x_1) \text{ AND } (x_2)$
0	0	0
0	1	0
1	0	0
1	1	1



This arrangement is mostly avoided, as training is very challenging

## Neural Network

When neurons are interconnected in layers



- Number of layers may differ
- Nodes in each intermediate layers may also differ
- Multiple output neurons are used for different class
- **Two levels deep** NN can represent any boolean function

## Neural Network Applications

NN is appropriate for problems with the following characteristics:

- Instances are provided by many attribute-value pairs (more data)
- The target function output may be discrete-valued, real-valued, or a vector of several real or discrete valued attributes
- The training examples may contain errors
- Long training times are acceptable
- Fast evaluation of the target function may be required
- The ability of humans to understand the learned target function is not important

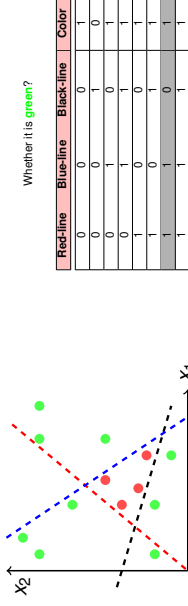
## Perceptron Training (delta rule)

**Algorithm 2:** Gradient Descent ( $D, \eta$ )

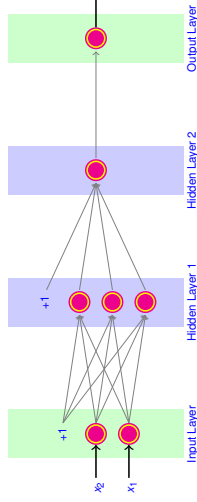
- 1 Initialize  $w_i$  with random weights
- 2 **repeat**
- 3     For each  $w_i$ , initialize  $\Delta w_i = 0$
- 4     **for each training example**  $d \in D$  **do**
- 5         Compute output  $o$  using model for  $d$  whose target is  $t$
- 6         For each  $w_i$ , update  $\Delta w_i = \Delta w_i + \eta(t - o)x_i$
- 7     For each  $w_i$ , set  $w_i = w_i + \Delta w_i$
- 8     **until** *termination condition is met*;
- 9     **return**  $w$

- A date item  $d \in D$ , is supposed to be multidimensional  $d = (x_1, x_2, \dots, x_n, t)$
- Algorithm converges toward the minimum error hypothesis.
- Linear programming can also be an approach

## More Example: Design NN for the following data



Whether it is green?



Note: Weights and activation in subsequent layers add power to the model in terms of non linearity.

## Perceptron Training (delta rule)

When data is not linearly-separable, error fluctuates with parameter update so, it becomes difficult to decide when to stop

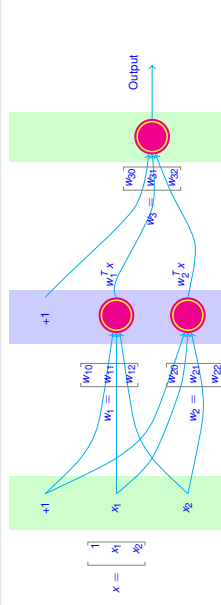
- **Delta rule** converges to a best-fit approximation of the target
- Uses **gradient descent**
- Consider unthresholded perceptron,  $\sigma(\vec{x}) = \vec{w} \cdot \vec{x}$
- Training error is defined as

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- Gradient would specify direction of steepest increase  $\nabla E(\vec{w}) = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$
- Weights can be learned as  $w_i = w_i - \eta \frac{\partial E}{\partial w_i}$
- It can be seen that  $\frac{\partial E}{\partial w_i} = \sum_{d \in D} d(t_d - o_d)(-x_{id})$

## Linear Activation is Not Much Interesting

NN with perceptrons have limited capability, even with many layers

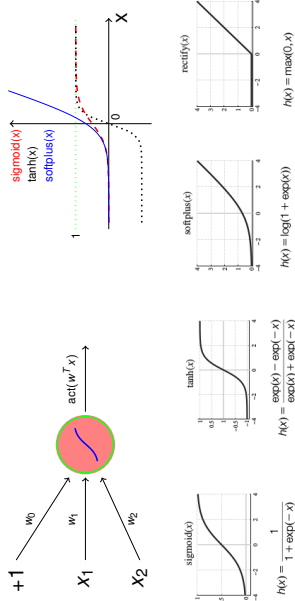


$$\begin{aligned} \text{Output} &= w_{30} \times 1 + w_{31} \times (w_1^T x) + w_{32} \times (w_2^T x) \\ &= w_{30} \times 1 + w_{31} \times (w_{10} \times 1 + w_{11} \times x_1 + w_{12} \times x_2) \\ &= w_{30} \times 1 + w_{32} \times (w_{20} \times 1 + w_{21} \times x_1 + w_{22} \times x_2) \\ &= (w_{30} + w_{31}w_{10} + w_{32}w_{20}) + (w_{31}w_{11} + w_{32}w_{21}) \times x_1 \\ &= w_0' + w_1' \times x_1 + w_2' \times x_2 \end{aligned}$$

**Expression of single perceptron**

## Neuron

Neuron uses nonlinear **activation functions** (sigmoid, tanh, ReLU, softplus etc.) at the place of thresholding



## Backpropagation (for 2 layers)

### Algorithm 3: Backpropagation ( $D, \eta, \eta_{in}, \eta_{out}, \eta_{hidden}$ )

- 1 Create the feedforward network with  $\eta_{in}, \eta_{out}, \eta_{hidden}$  layers
  - 2 Randomly initialize weights to small values  $\in [-0.05, +0.05]$
  - 3 **repeat**
  - 4     **for each**  $\langle \vec{x}, \vec{t} \rangle \in D$  **do**
  - 5          $O_u$  = get output from network  $\forall$  unit  $u$
  - 6          $\delta_k = o_k(1 - o_k)(t_k - o_k)$  for all **output unit**  $k$
  - 7          $\delta_h = o_h(1 - o_h) \sum_{k \in \text{outputs}} (W_{kh} \delta_k)$  for all **hidden unit**  $h$
  - 8          $W_{ij} = W_{ij} + \Delta W_{ij}$  where  $\Delta W_{ij} = \eta \delta_j x_{ij}$
  - 9     **until converge;**
- Recall error function is  $E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - O_{kd})^2$
- For a single training example  $E_d(\vec{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - O_k)^2$
- Weight  $w_{ji}$  is updated by adding  $\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$

## Value of $\frac{\partial E_d}{\partial \text{net}_j}$ for output units

- $\frac{\partial E_d}{\partial \text{net}_j} = \frac{\partial E_d}{\partial o_j} \times \frac{\partial o_j}{\partial \text{net}_j}$
  - $\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2 = \frac{\partial}{\partial o_j} \frac{1}{2} (t_j - o_j)^2 = -(t_j - o_j)$
  - Note that  $o_j = \sigma(\text{net}_j)$  therefore  $\frac{\partial o_j}{\partial \text{net}_j}$  is derivative of sigmoid
- $$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \frac{1}{1 + e^{-x}} = \frac{(-1)(1 + e^{-x})^{-2} \frac{d}{dx} (1 + e^{-x})}{(1 + e^{-x})^2} \\ &= \frac{(-1)(1 + e^{-x})^{-2} (0 - e^{-x})}{1 + e^{-x} + 1 - 1} \\ &= \frac{1}{1 + e^{-x}} \times \frac{1 - 1}{1 + e^{-x}} = \sigma(x)(1 - \sigma(x)) \end{aligned}$$
- As a result  $\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial \sigma(\text{net}_j)}{\partial \text{net}_j} = \sigma(\text{net}_j)(1 - \sigma(\text{net}_j)) = o_j(1 - o_j)$
  - $\frac{\partial E_d}{\partial \text{net}_j} = -(t_j - o_j) o_j (1 - o_j)$
  - Term  $(t_j - o_j) o_j (1 - o_j)$  is treated as  $\delta_j$

Therefore,  $\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial \text{net}_j} \times x_{ij} = \eta (t_j - o_j) o_j (1 - o_j) x_{ij}$

## Multilayer Networks and Backpropagation

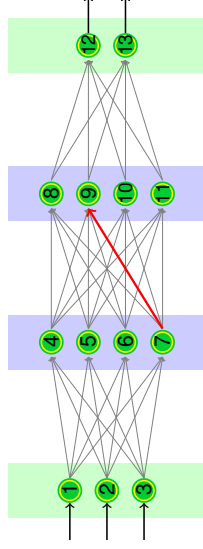
- Single perceptron can only express linear decision surface
- We need units whose output is a nonlinear function of its inputs AND is also differentiable (Use Neuron not Perceptron)

$$o(\vec{x}) = \sigma(\vec{w} \cdot \vec{x})$$

- where  $\sigma(y) = \frac{1}{1 + e^{-y}}$
- **Backpropagation** algorithm learns the weights for a fixed set of units and interconnections
  - It employs **gradient descent** to minimize the error between the network output values and the target values for these outputs
  - Let Error function is redefined as

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - O_{kd})^2$$

## Conventions Over The Network



- $x_{ij}$   $i$ th input to unit  $j$  ( $x_{44}$  is highlighted)
- $w_{ij}$  weight associated with  $i$ th input to unit  $j$
- $\text{net}_j$  be  $\sum_i w_{ij} x_{ij}$  the weighted sum of input for unit  $j$
- $o_j$  output computed by unit  $j$ . Consider it as  $\sigma(\text{net}_j)$
- $t_j$  target output for unit  $j$
- outputs** set of units in final layer ( $\{12, 13\}$  in our case)
- Downstream( $j$ )** units whose immediate input is the output of unit  $j$

We are interested in  $\frac{\partial E_d}{\partial \text{net}_j}$  it is  $\frac{\partial E_d}{\partial \text{net}_j} \times \frac{\partial \text{net}_j}{\partial w_{ij}}$  and therefore,  $\frac{\partial E_d}{\partial \text{net}_j} \times x_{ij}$

## Value of $\frac{\partial E_d}{\partial \text{net}_j}$ for hidden units

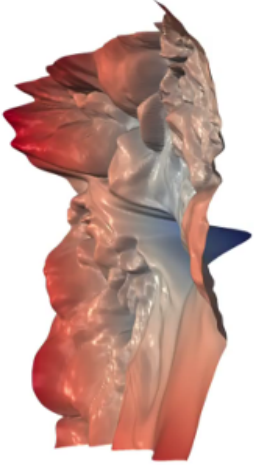
$$\begin{aligned} \frac{\partial E_d}{\partial \text{net}_j} &= \sum_{k \in \text{Downstream}(j)} \frac{\partial E_d}{\partial \text{net}_k} \times \frac{\partial \text{net}_k}{\partial \text{net}_j} \\ &= \sum_{k \in \text{Downstream}(j)} -\delta_k \times \frac{\partial \text{net}_k}{\partial \text{net}_j} \\ &= \sum_{k \in \text{Downstream}(j)} -\delta_k \times \frac{\partial \text{net}_k}{\partial o_j} \times \frac{\partial o_j}{\partial \text{net}_j} \\ &= \sum_{k \in \text{Downstream}(j)} -\delta_k \times w_{kj} \times o_j (1 - o_j) \end{aligned}$$

- $\delta_j$  being  $-\frac{\partial E_d}{\partial \text{net}_j} = o_j(1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k \times w_{kj}$

Therefore,  $\Delta w_{ji} = \eta \delta_j x_{ij} = \eta (o_j(1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k \times w_{kj}) x_{ij}$

Note:  $\text{net}_k = w_{k1}x_{k1} + w_{k2}x_{k2} + \dots + w_{kj}x_{kj} + \dots$  as  $o_j$  is input to  $k$ th unit so  $o_j = x_{kj}$ ; so  $\frac{\partial \text{net}_k}{\partial o_j} = w_{kj}$

## Loss Landscape <sup>1</sup>



Backpropagation over multilayer networks converge to a local minimum, **NOT necessarily to the global minima**

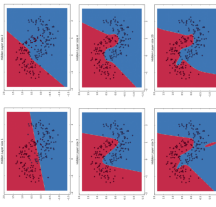
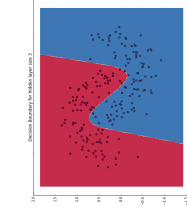
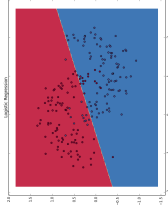
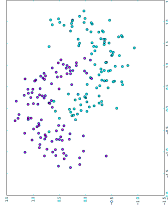
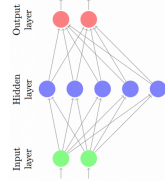
<sup>1</sup>Visualizing the Loss Landscape of Neural Nets, Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, Tom Goldstein, 2018

## Backpropagation

- Result of Backpropagation over multilayer networks is only guaranteed to converge toward some local minimum and not necessarily to the global minimum error
- No methods are known to predict with certainty when local minima will cause difficulties
- Suggested to use momentum, true gradient descent or multiple networks (initialized with different random weights)
- Any boolean function can precisely be represented by some network having only **two** layers of units (Cybenko 1989; Hornik et al. 1989)
- Every bounded continuous function can be approximated with arbitrarily small error (under a finite norm) by a network with two layers of units
- Any function can be approximated to arbitrary accuracy by a network with three layers of units (Cybenko 1988)

## Coding Example <sup>2</sup>

Consider two class data  
 Logistic Regression  
 Neural Network  
 Decision Boundary  
 Bigger hidden layer



<sup>2</sup><https://github.com/dennybritz/nn-from-scratch/blob/master/nn-from-scratch.ipynb>

## Backpropagation

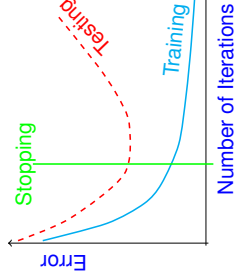
- Adding Momentum:** weight update during  $n^{\text{th}}$  iteration depend partially on the update that occurred during the  $(n - 1)^{\text{th}}$  iteration
- Learning in arbitrary acyclic network:** for feed-forward networks of arbitrary depth,  $\delta_r$  value for a unit in hidden layer is determined as

$$\Delta W_{ji}(n) = \eta \delta_j^r X_{ji} + \alpha \Delta W_{ji}(n - 1)$$

$$\delta_r = o_r(1 - o_r) \sum_{s \in \text{Downstream}(r)} W_{sr} \times \delta_s$$

## Generalization, Overfitting, and Stopping Criterion

Continue training until the error on the training examples falls below some predetermined threshold could be a poor strategy



- Weight decay or use of validation set ( $k$ -fold ?) is suggested
- Input or output encoding can be used

## An Example Code [1/4] <sup>3</sup>

```
import matplotlib.pyplot as plt
import numpy as np
import sklearn.datasets
import sklearn.linear_model
import matplotlib

matplotlib.rcParams['figure.figsize'] = (10.0, 8.0)
np.random.seed(0)
datasets.make_moons(200, noise=0.20)
X, y = datasets.load_data()
plt.scatter(X[:, 0], X[:, 1], s=40, c=y, cmap=plt.cm.Spectral)

def plot_decision_boundary(pred_func):
    x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
    y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
    xx = yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
    Z = pred_func(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)
    plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral)
    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Spectral)
    plt.title("Logistic Regression")
    plt.show()
```

<sup>3</sup><https://github.com/dennybritz/nn-from-scratch/blob/master/nn-from-scratch.ipynb>

## An Example Code [2/4] <sup>4</sup>

```
##### Implementing NN #####
num_examples = len(X) # training set size
nn_input_dim = 2 # input layer dimensionality
nn_output_dim = 2 # output layer dimensionality

epsilon = 0.01 # learning rate for gradient descent
reg_lambda = 0.01 # regularization strength

def calculate_loss(model):
    W1, W2, b2 = model['W1'], model['b1'], model['W2'], model['b2']
    z1 = X.dot(W1) + b1
    a1 = np.tanh(z1)
    z2 = a1.dot(W2) + b2
    exp_scores = np.exp(z2)
    probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
    correct_logprobs = -np.log(probs[range(num_examples), y])
    data_loss = np.sum(correct_logprobs)
    d = 1.0 / num_examples
    return 1./num_examples * data_loss

def predict(model, x):
    W1, W2, b2 = model['W1'], model['b1'], model['W2'], model['b2']
    z1 = x.dot(W1) + b1
    a1 = np.tanh(z1)
    z2 = a1.dot(W2) + b2
    exp_scores = np.exp(z2)
    probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
    return np.argmax(probs, axis=1)
```

<sup>4</sup><https://github.com/dennybritz/nn-from-scratch/blob/master/nn-from-scratch.ipynb>   Machine Learning (BITS F4-6) M W F (10-11 AM) online@BITS-Pilani Lecture-31(April 09, 2021) 31/34

## An Example Code [4/4] <sup>6</sup>

```
# Build a model with a 3-dimensional hidden layer
model = build_model(3, print_loss=True)

# Plot the decision boundary
plot_decision_boundary(lambda x: predict(model, x))
plt.title("Decision Boundary for hidden layer size 3")

plt.show()

##### Varying the number of hidden layers #####
plt.figure(figsize=(16, 32))
hidden_layer_dimensions = [1, 2, 3, 4, 5, 20, 50]
for i, nn_hdim in enumerate(hidden_layer_dimensions):
    plt.subplot(5, 2, i+1)
    plt.title("%d-Hidden Layer size %d" % nn_hdim)
    plot_decision_boundary(lambda x: predict(model, x))
plt.show()
```

<sup>6</sup><https://github.com/dennybritz/nn-from-scratch/blob/master/nn-from-scratch.ipynb>   Machine Learning (BITS F4-6) M W F (10-11 AM) online@BITS-Pilani Lecture-31(April 09, 2021) 33/34

## An Example Code [3/4] <sup>5</sup>

```
def build_model(nn_hdim, num_passes=20000, print_loss=False):
    np.random.seed(0)
    W1 = np.random.randn(nn_input_dim, nn_hdim) / np.sqrt(nn_input_dim)
    W2 = np.zeros((1, nn_hdim)), np.random.randn(nn_output_dim) / np.sqrt(nn_hdim)
    b2 = np.zeros((1, nn_output_dim))
    model = {'W1': W1, 'W2': W2, 'b2': b2}
    for i in range(0, num_passes):
        z1 = X.dot(W1) + b1
        a1 = np.tanh(z1)
        z2 = a1.dot(W2) + b2
        exp_scores = np.exp(z2)
        probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
        delta3 = range(num_examples), y) == 1
        dW2 = (a1.T).dot(delta3)
        dB2 = np.sum(delta3, axis=0, keepdims=True)
        dW1 = np.dot(X.T, delta2)
        dB1 = np.sum(delta2, axis=0)
        dW2 += reg_lambda * W2
        dW1 += reg_lambda * W1
        W1 += -epsilon * dW1
        W2 += -epsilon * dW2
        b2 += -epsilon * dB2
    model = {'W1': W1, 'b1': b1, 'W2': W2, 'b2': b2}
    return model
```

<sup>5</sup><https://github.com/dennybritz/nn-from-scratch/blob/master/nn-from-scratch.ipynb>   Machine Learning (BITS F4-6) M W F (10-11 AM) online@BITS-Pilani Lecture-31(April 09, 2021) 32/34

## Thank You!

Thank you very much for your attention!

Machine Learning (BITS F4-6) M W F (10-11 AM) online@BITS-Pilani Lecture-31(April 09, 2021) 34/34