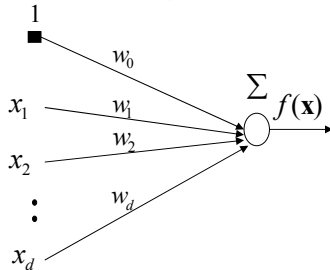**CS 2750 Machine Learning**
**Lecture 13**

# Multi-layer Neural Networks

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

---

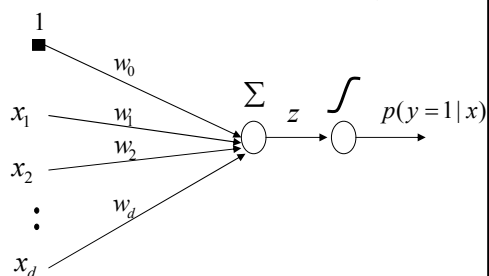# Limitations of basic linear units

**Linear regression**

$$f(\mathbf{x}) = w_0 + \sum_{j=1}^{d} w_j x_j$$
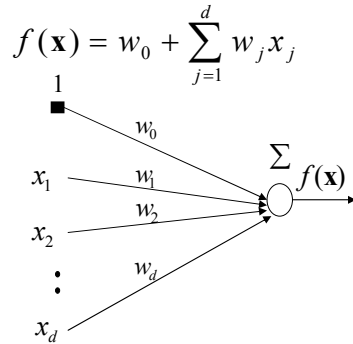


**Function linear in inputs !!**

**Logistic regression**

$$f(\mathbf{x}) = p(y=1 \mid \mathbf{x}, \mathbf{w}) = g(w_0 + \sum_{j=1}^{d} w_j x_j)$$



**Linear decision boundary!!**
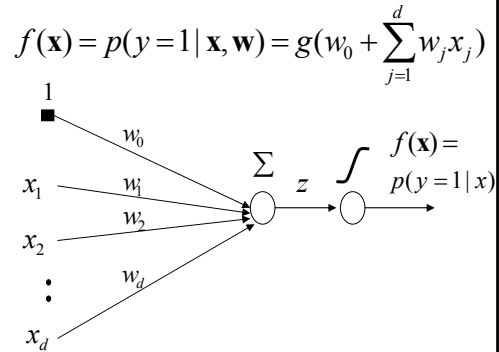
# Linear units

### Linear regression

$$f(\mathbf{x}) = w_0 + \sum_{j=1}^{d} w_j x_j$$



### Logistic regression

$$f(\mathbf{x}) = p(y = 1 \mid \mathbf{x}, \mathbf{w}) = g(w_0 + \sum_{j=1}^{d} w_j x_j)$$

**On-line gradient update:**

$$w_0 \leftarrow w_0 + \alpha(y - f(\mathbf{x}))$$
$$\vdots$$
$$w_j \leftarrow w_j + \alpha(y - f(\mathbf{x}))x_j$$

**The same**

**On-line gradient update:**

$$w_0 \leftarrow w_0 + \alpha(y - f(\mathbf{x}))$$
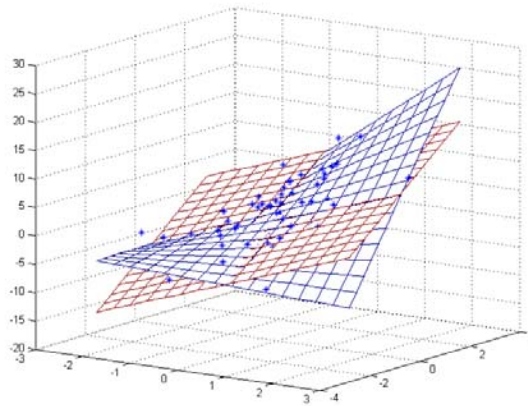$$\vdots$$
$$w_j \leftarrow w_j + \alpha(y - f(\mathbf{x}))x_j$$

---

# Regression with the quadratic model.

**Limitation:** linear hyper-plane only
- a non-linear surface can be better
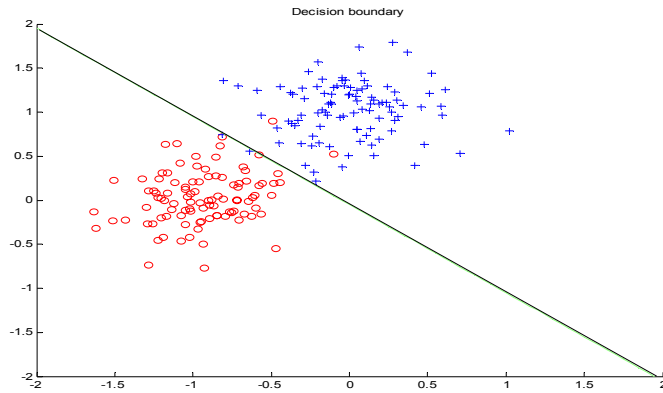
# Classification with the linear model.

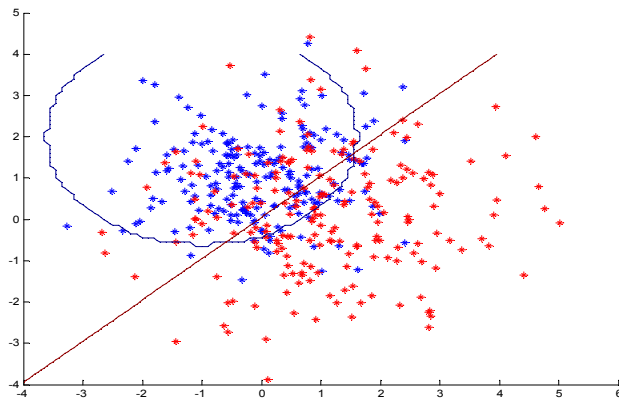**Logistic regression model defines a linear decision boundary**

- Example: 2 classes (blue and red points)



Decision boundary

# Linear decision boundary
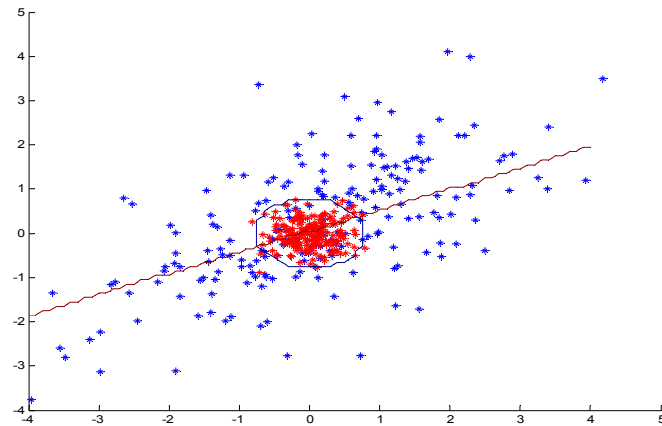
- logistic regression model is not optimal, but not that bad

# When logistic regression fails?
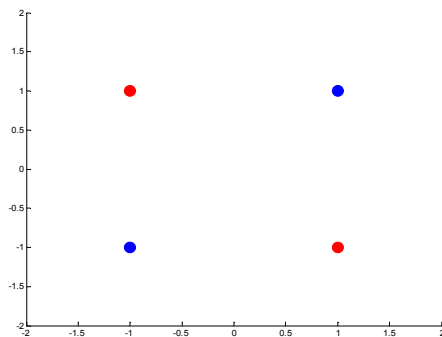
- Example in which the logistic regression model fails

# Limitations of linear units.

- Logistic regression does not work for **parity function**s
  - no linear decision boundary exists



**Solution:** a model of a non-linear decision boundary
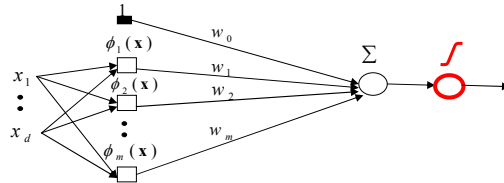
# Extensions of simple linear units

**Feature (basis) functions** model **nonlinearities**

**Linear regression**

$$f(\mathbf{x}) = w_0 + \sum_{j=1}^{m} w_j \phi_j(\mathbf{x})$$

**Logistic regression**

$$f(\mathbf{x}) = g\left(w_0 + \sum_{j=1}^{m} w_j \phi_j(\mathbf{x})\right)$$



**Important property:**
• The same problem as learning of the weights for linear units , the input has changed– but the weights are linear in the new input
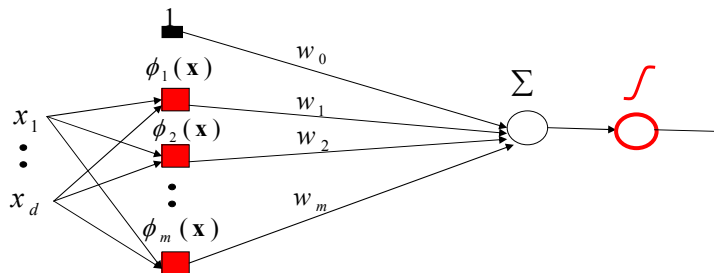**Problem:** too many weights to learn

---

# Multi-layer neural networks

• **Problems of extended linear units:**
  • fixed basis functions,
  • too many weights
• **One possible solution:**
  • Assume parametric feature (basis) functions
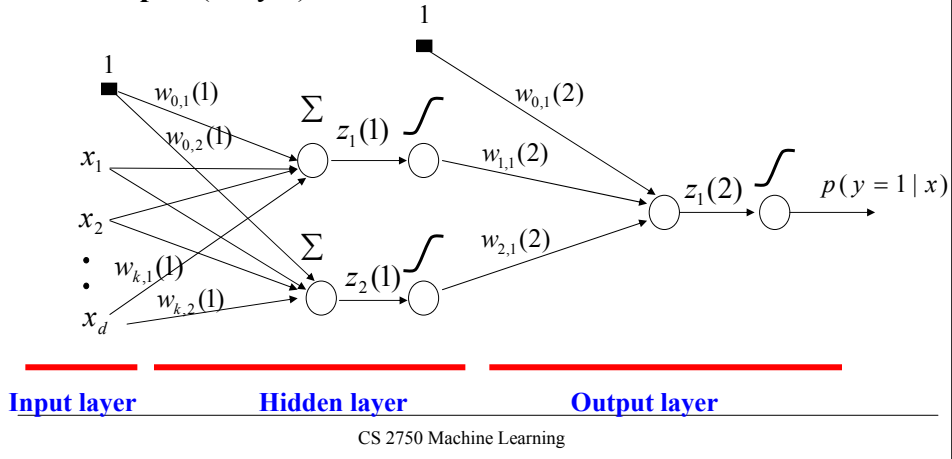  • Learn the parameters together with the remaining weights

# Multilayer neural network

Also called a **multilayer perceptron (MLP)**

Cascades multiple logistic regression units

**Example:** (2 layer) classifier with non-linear decision boundaries



**Input layer**            **Hidden layer**            **Output layer**

---

# Multilayer neural network

- Models **non-linearities through logistic regression units**
- Can be applied to both **regression and binary classification problems**



**Input layer**            **Hidden layer**            **Output layer**

**regression**

$$f(\mathbf{x}) = f(\mathbf{x}, \mathbf{w})$$

**classification**

$$f(\mathbf{x}) = p(y = 1 \mid \mathbf{x}, \mathbf{w})$$

**option**

# Multilayer neural network

- **Non-linearities are modeled using multiple hidden logistic regression units (organized in layers)**
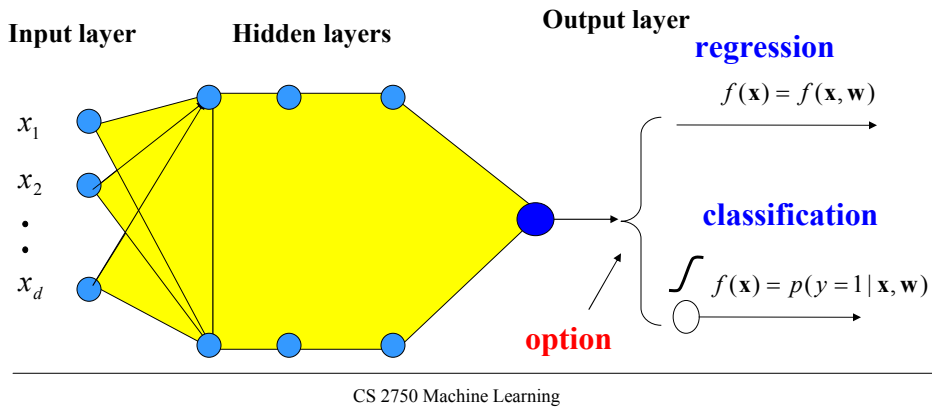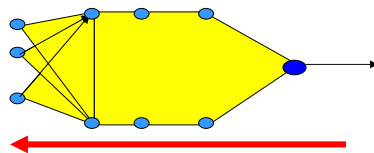- The output layer determines whether it is a **regression or a binary classification problem**

**Input layer**    **Hidden layers**    **Output layer**

**regression**

$$f(\mathbf{x}) = f(\mathbf{x}, \mathbf{w})$$

**classification**

$$f(\mathbf{x}) = p(y = 1 \mid \mathbf{x}, \mathbf{w})$$

$x_1$

$x_2$

$\vdots$

$x_d$

**option**

---

# Learning with MLP

- How to learn the parameters of the neural network?
- **Onl-line gradient descent algorithm**
  - Weight updates based on the error:   $J_{\text{online}}(D_i, \mathbf{w})$

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} J_{\text{online}}(D_i, \mathbf{w})$$

- We need to **compute gradients for weights in all units**
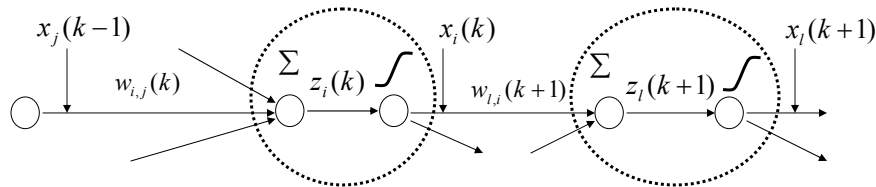- **Can be computed in one backward sweep through the net !!!**

- The process is called **back-propagation**

# Backpropagation

$x_i(k)$ - output of the unit i on level k

$z_i(k)$ - input to the sigmoid function on level k

$w_{i,j}(k)$ - weight between units j and i on levels (k-1) and k

$$z_i(k) = w_{i,0}(k) + \sum_j w_{i,j}(k)x_j(k-1)$$

$$x_i(k) = g(z_i(k))$$

---

# Backpropagation

**Update weight** $w_{i,j}(k)$ using a data point $D_u = <\mathbf{x}, y>$

$$w_{i,j}(k) \leftarrow w_{i,j}(k) - \alpha \frac{\partial}{\partial w_{i,j}(k)} J_{online}(D_u, \mathbf{w})$$

Let $\delta_i(k) = \frac{\partial}{\partial z_i(k)} J_{online}(D_u, \mathbf{w})$

Then: $\frac{\partial}{\partial w_{i,j}(k)} J_{online}(D_u, \mathbf{w}) = \frac{\partial J_{online}(D_u, \mathbf{w})}{\partial z_i(k)} \frac{\partial z_i(k)}{\partial w_{i,j}(k)} = \delta_i(k)x_j(k-1)$

S.t. $\delta_i(k)$ is computed from $x_i(k)$ and the next layer $\delta_l(k+1)$

$$\delta_i(k) = \left[ \sum_l \delta_l(k+1)w_{l,i}(k+1) \right] x_i(k)(1-x_i(k))$$

**Last unit** (is the same as for the regular linear units):

$$\delta_i(K) = -(y - f(\mathbf{x}, \mathbf{w}))$$

It is the same for the classification with the log-likelihood measure of fit and linear regression with least-squares error!!!

# Learning with MLP

- **Online gradient descent algorithm**
  - Weight update for example $D_u$:

$$w_{i,j}(k) \leftarrow w_{i,j}(k) - \alpha \frac{\partial}{\partial w_{i,j}(k)} J_{online}(D_u, \mathbf{w})$$

$$\frac{\partial}{\partial w_{i,j}(k)} J_{online}(D_u, \mathbf{w}) = \frac{\partial J_{online}(D_u, \mathbf{w})}{\partial z_i(k)} \frac{\partial z_i(k)}{\partial w_{i,j}(k)} = \delta_i(k) x_j(k-1)$$

$$\boxed{w_{i,j}(k) \leftarrow w_{i,j}(k) - \alpha \delta_i(k) x_j(k-1)}$$

$x_j(k-1)$  - j-th output of the (k-1) layer

$\delta_i(k)$  - derivative computed via backpropagation

$\alpha$  - a learning rate

# Online gradient descent algorithm for MLP

**Online-gradient-descent** (*D, number of iterations*)

   **Initialize** all weights  $w_{i,j}(k)$

   **for** *i*=1:1: *number of iterations*

      **do**      **select** a data point $D_u = <\boldsymbol{x}, y>$ from $D$

                  **set  learning rate**  $\alpha$

                  **compute** outputs  $x_j(k)$  for each unit

                  **compute** derivatives $\delta_i(k)$ via **backpropagation**

                  **update** all weights (in parallel)

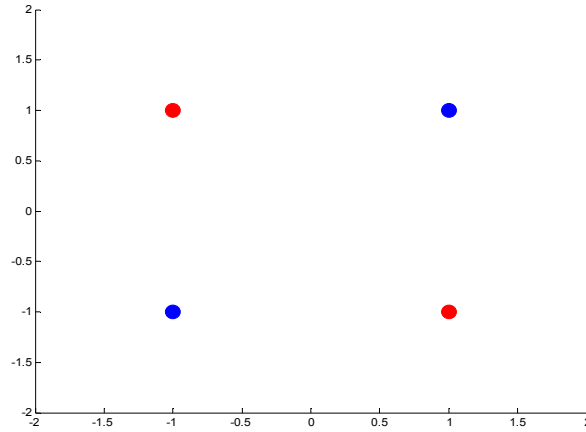$$w_{i,j}(k) \leftarrow w_{i,j}(k) - \alpha \delta_i(k) x_j(k-1)$$

   **end for**

   **return** weights **w**

# Xor Example.

• linear decision boundary does not exist

# Xor example. Linear unit

# Xor example.
## Neural network with 2 hidden units

# Xor example.
## Neural network with 10 hidden units

# MLP in practice

- **Optical character recognition** – digits 20x20
  - Automatic sorting of mails
  - 5 layer network with multiple output functions

**10 outputs (0,1,…9)**

. . .

**20x20 = 400  inputs**

| layer | Neurons | Weights |
|-------|---------|---------|
| 5 | 10 | 3000 |
| 4 | 300 | 1200 |
| 3 | 1200 | 50000 |
| 2 | 784 | 3136 |
| 1 | 3136 | 78400 |