

Data Analytics

EEE 4774 & 6777

Module 4 - Classification

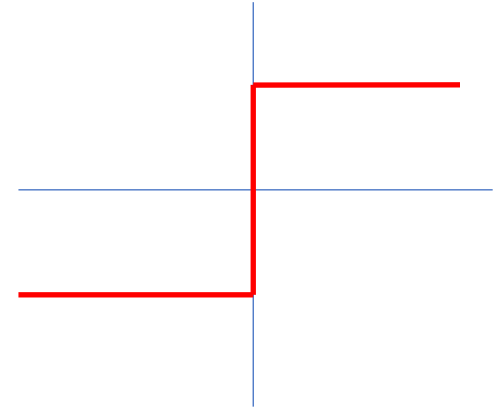
Perceptron, Multilayer Perceptron (MLP) / Artificial Neural Network

Spring 2022

Perceptron

- 2-class model, $t_n \in \{-1, +1\}$

$$y(\mathbf{x}) = f(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})) = \begin{cases} +1, & \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \geq 0 \\ -1, & \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) < 0 \end{cases}$$



Cost/Error function: (**Perceptron criterion**) $E_P(\mathbf{w}) = -\sum_{n \in \mathcal{M}} \mathbf{w}^T \boldsymbol{\phi}_n t_n$

\mathcal{M} : set of misclassified instances

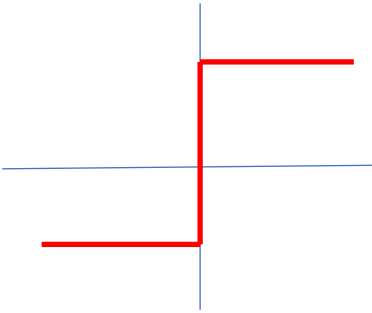
Stochastic gradient descent: $\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta \nabla E_P(\mathbf{w}) \cong \mathbf{w}^{(i)} + \eta \boldsymbol{\phi}_n t_n$

- When \mathbf{w} is multiplied by a constant, $y(\mathbf{x})$ does not change, hence w.l.o.g. $\eta = 1$

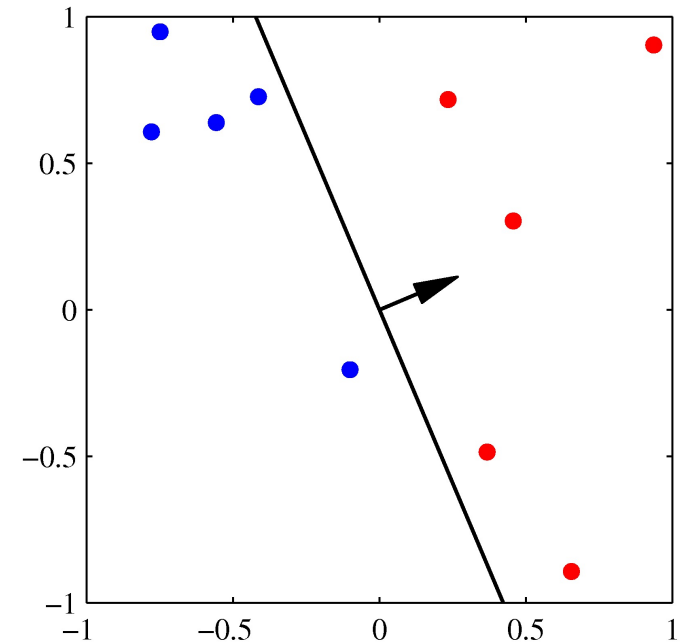
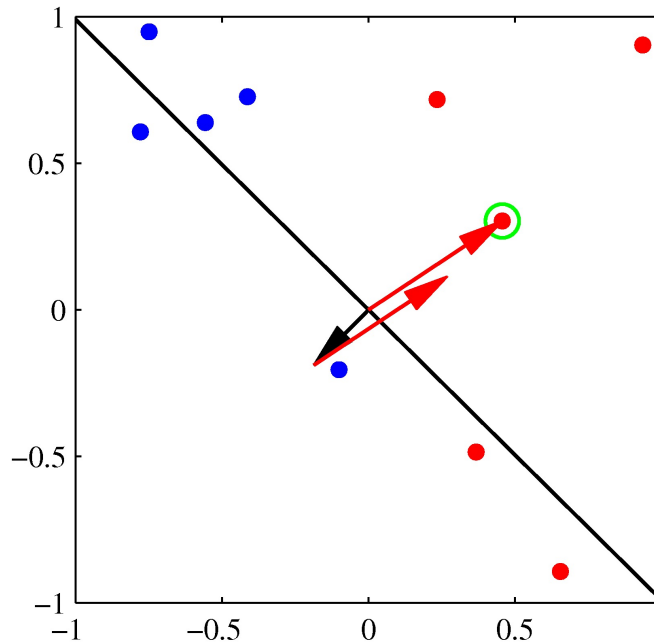
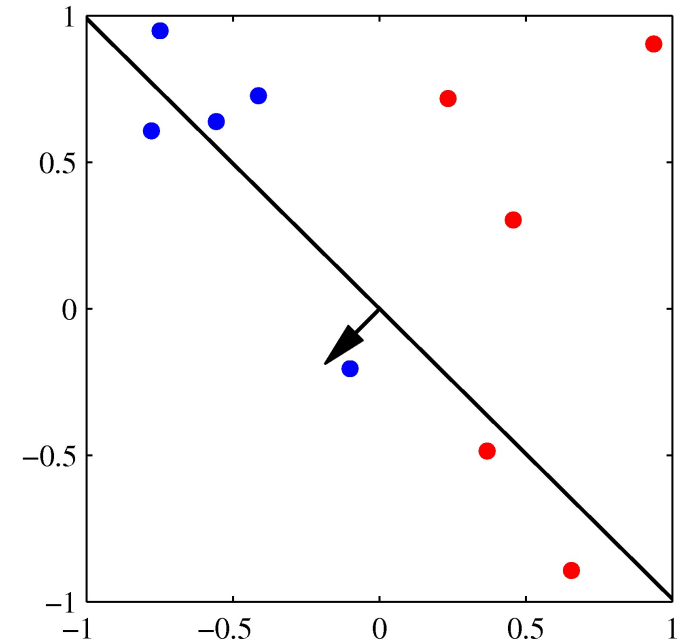
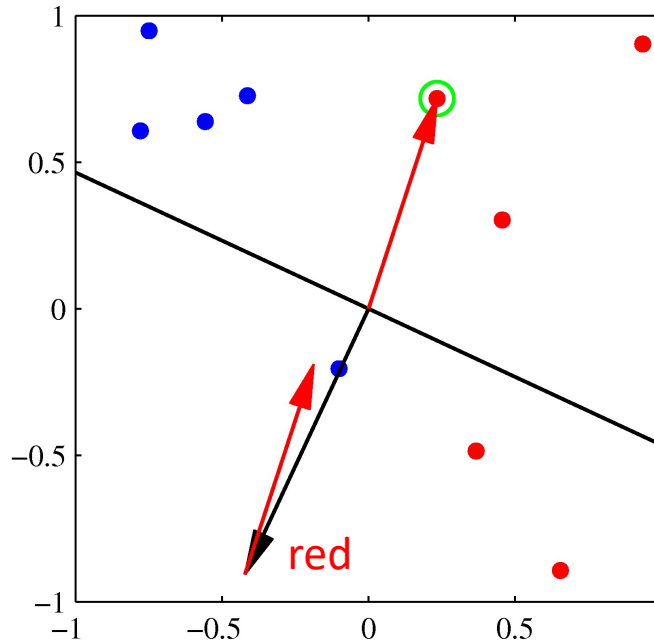
Perceptron procedure: For each instance x_n

- Compute $f(\mathbf{w}^T \phi(x_n))$

- Decide



- If incorrectly classified
 - For class 1, add $\phi(x_n)$ to the \mathbf{w} estimate
 - For class 2, subtract $\phi(x_n)$ from the \mathbf{w} estimate



- At each step error for the considered point decreases

$$-\mathbf{w}^{(i+1)T} \phi_n t_n = -\mathbf{w}^{(i)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(i)T} \phi_n t_n$$

- No guarantee that the total error decreases,
e.g., a previously correctly classified point may be misclassified after the update

Perceptron convergence theorem:

If training data is linearly separable, the perceptron algorithm is guaranteed to find an exact solution in a finite number of steps

- Number of steps can be substantial
- Solution found depend on initialization (many solutions possible)
- Never converge if not linearly separable
- Not possible to distinguish between a nonseparable problem and slow convergence
- Does not generalize readily to multi-class
- Fixed basis functions, not adaptive to the input

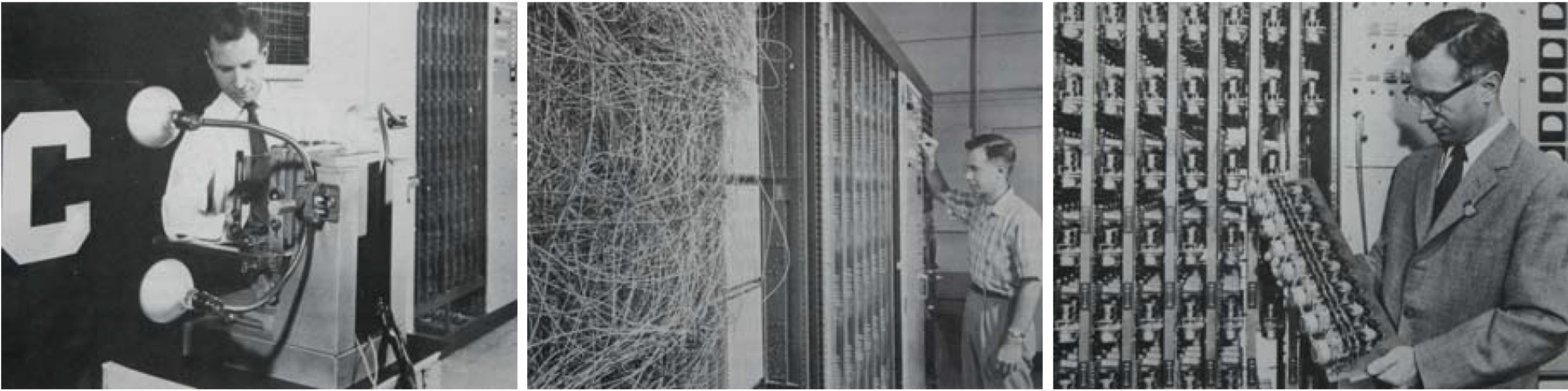


Figure 4.8 Illustration of the Mark 1 perceptron hardware. The photograph on the left shows how the inputs were obtained using a simple camera system in which an input scene, in this case a printed character, was illuminated by powerful lights, and an image focussed onto a 20×20 array of cadmium sulphide photocells, giving a primitive 400 pixel image. The perceptron also had a patch board, shown in the middle photograph, which allowed different configurations of input features to be tried. Often these were wired up at random to demonstrate the ability of the perceptron to learn without the need for precise wiring, in contrast to a modern digital computer. The photograph on the right shows one of the racks of adaptive weights. Each weight was implemented using a rotary variable resistor, also called a potentiometer, driven by an electric motor thereby allowing the value of the weight to be adjusted automatically by the learning algorithm.

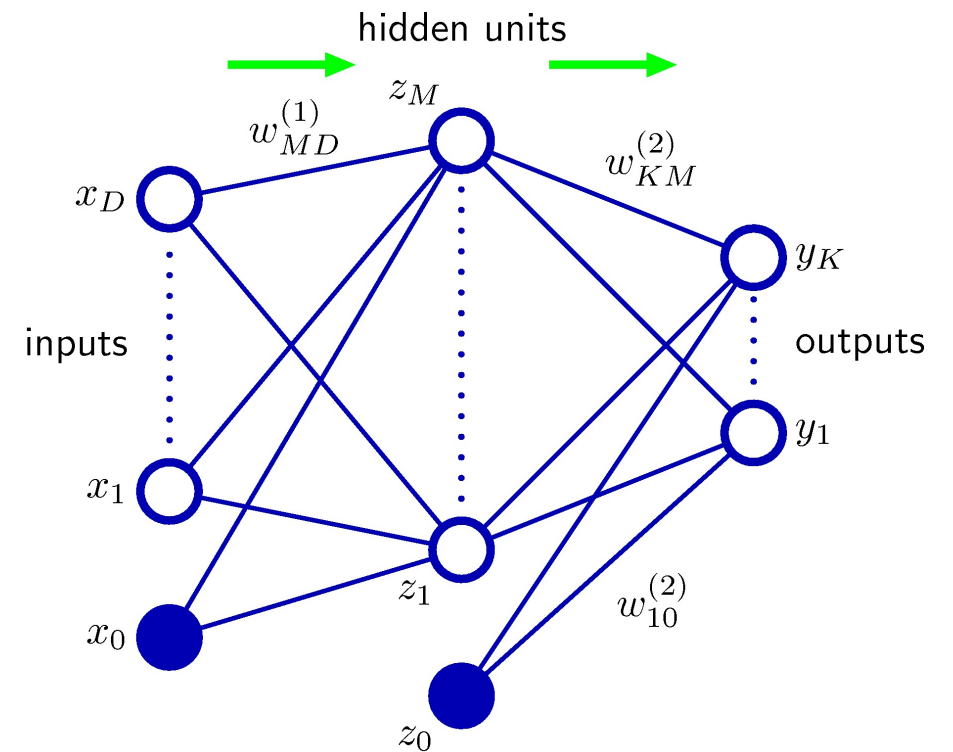
Artificial Neural Network (ANN) / Multilayer Perceptron (MLP)

- Origins in attempts to find mathematical representations of information processing in biological systems
- Also known as *Multilayer Perceptron* / Can be also regarded as *Multilayer Logistic Regression*
- Finds basis functions $\phi(\mathbf{x})$ adaptive to the training data

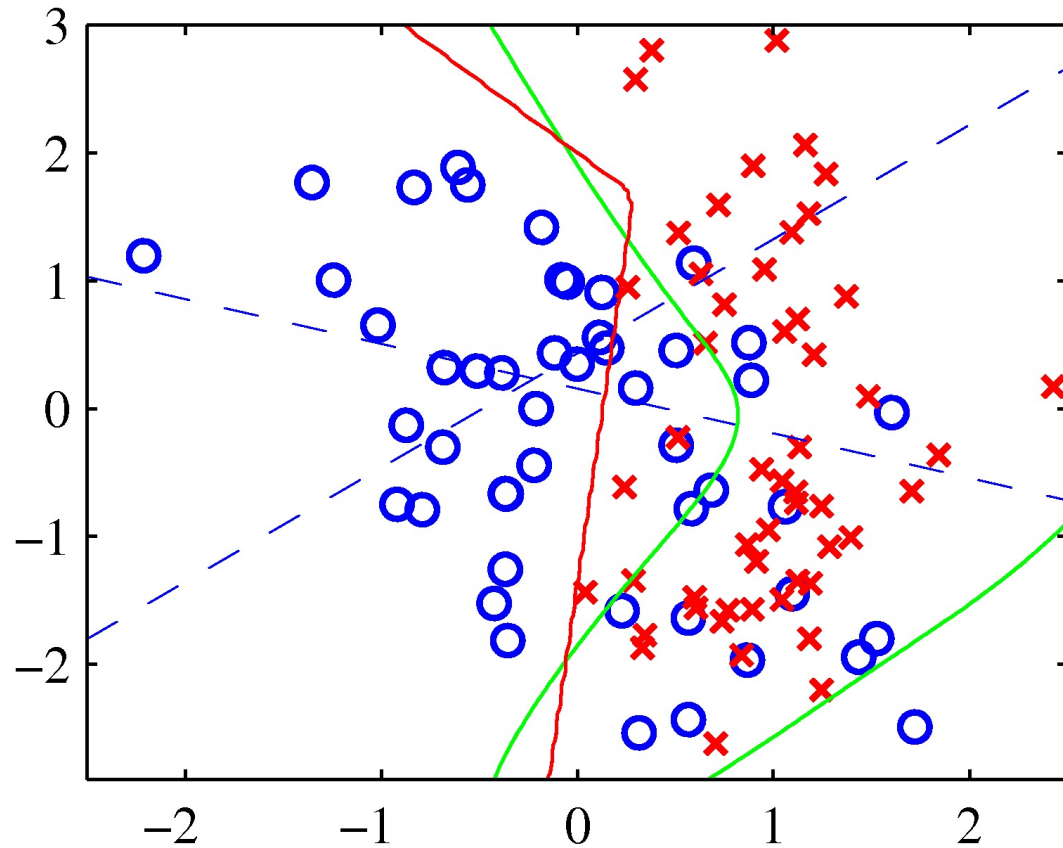
$$y_k(\mathbf{x}) = f(\mathbf{w}_{2k}^T \phi(\mathbf{x}) + b_{2k}) = \sigma(\mathbf{w}_{2k}^T \phi(\mathbf{x}) + b_{2k}), k = 1, \dots, K$$

$$\phi_m(\mathbf{x}) = h(\mathbf{w}_{1m}^T \mathbf{x} + b_{1m}) = \sigma(\mathbf{w}_{1m}^T \mathbf{x} + b_{1m}), m = 1, \dots, M$$

tanh (hyperbolic tangent), another sigmoid function, ReLU, GELU, etc. can be used instead of logistic sigmoid function.



Neural Network for Classification



- Two hidden units with activation functions

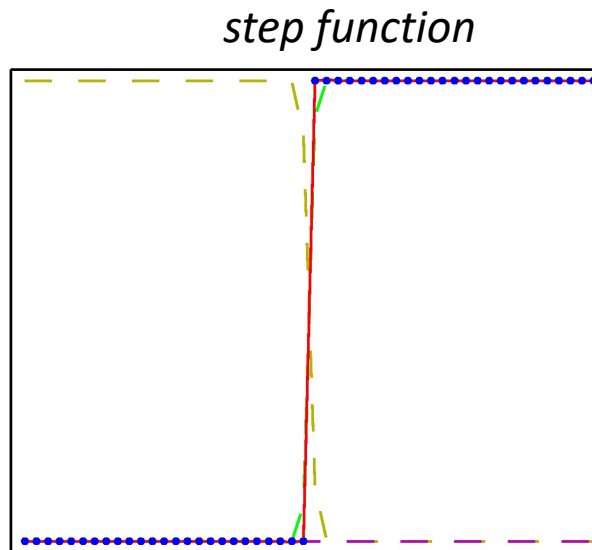
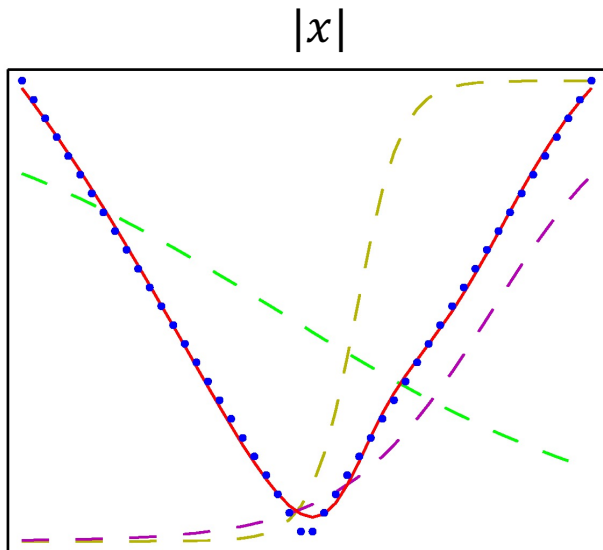
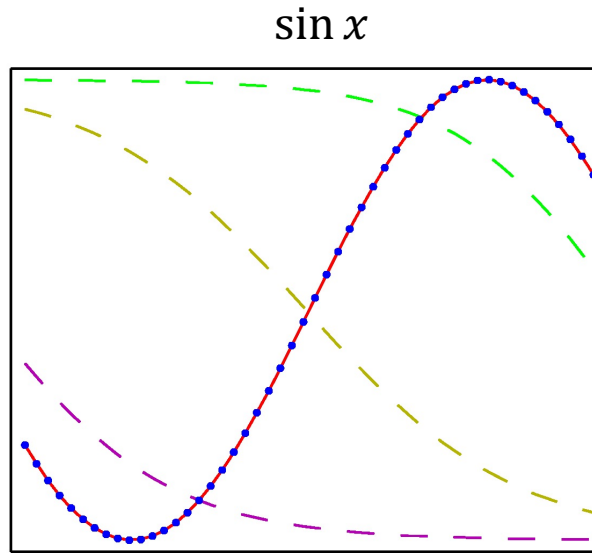
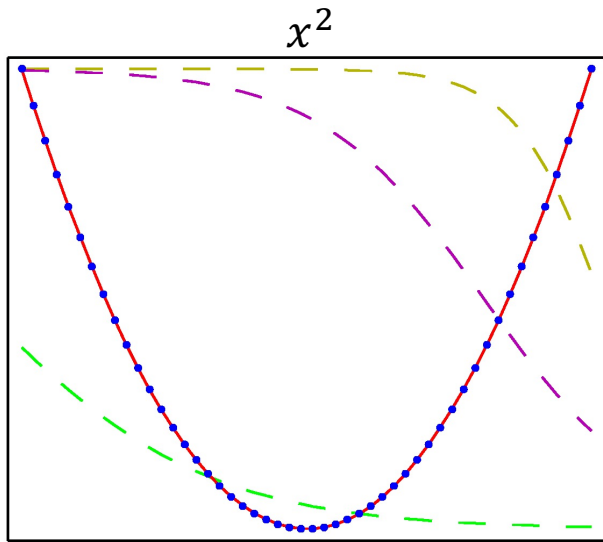
$$h(\mathbf{w}_{1i}^T \mathbf{x}) = \tanh(\mathbf{w}_{1i}^T \mathbf{x})$$

- Single output with activation function

$$f(\mathbf{w}_2^T \mathbf{h}(\mathbf{x})) = \sigma(\mathbf{w}_2^T \mathbf{h}(\mathbf{x}))$$

- Blue dashed lines show $h = 0$ for hidden units
- Red line shows $f = 0.5$ decision boundary
- Green line shows the optimum decision boundary

Neural Network for Regression



- $N = 50$ data points
- 2-layer network having 3 hidden units with \tanh activation function
- Single output unit with linear activation function
- Dashed curves show hidden units
- Red curves show output functions

Artificial Neural Network

An ANN is typically defined by three types of parameters:

1. f_j : The activation function that converts a neuron's weighted input to its output activation.
2. The interconnection pattern between the different layers of neurons
3. w_j, b_j : The weights of the interconnections and offset, which are updated in the learning process.

Supervised Learning:

- a) Feed the input x_i , and obtain the output y_i
- b) Compare the output y_i with the ground truth t_i , resulting in the error $y_i - t_i$
- c) Adjust the weights w_j to minimize a function of the error, e.g.,
 $t_i \log y_i + (1 - t_i) \log(1 - y_i)$ for classification or $(y_i - t_i)^2$ for regression

Learning ANN coefficients

- Mean squared error (MSE):
 - MAE, Huber loss, etc.
- $$E(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N [t_i - y_i(\mathbf{w}, \mathbf{x}_i)]^2$$
- **Gradient Descent** based on Backpropagation
 - Efficient technique for computing the **gradient** of error function $E(\mathbf{w})$.
 - **Local message passing** scheme in which information is sent alternately forward and backward through the network.
 - Given input \mathbf{x} , the **information flows forward to produce prediction y and scalar cost $E(\mathbf{w})$** .
 - In computing the derivatives, **errors are propagated backwards** through the network using the chain rule of differentiation, hence the name.
 - Finally, **an optimization algorithm** such as Stochastic Gradient Descent (SGD), Adam, RMSProp, etc. **updates the weights using the computed gradient**.

Backpropagation Algorithm

$$E(t, f_L(W_L f_{L-1}(W_{L-1} \dots f_2(W_2 f_1(W_1 x))))))$$

$$\frac{dE}{da_L} \circ \frac{da_L}{dz_L} \frac{dz_L}{da_{L-1}} \circ \frac{da_{L-1}}{dz_{L-1}} \frac{dz_{L-1}}{da_{L-2}} \dots \frac{da_1}{dz_1} \frac{dz_1}{dx}$$

$$\frac{dE}{da_L} \circ f'_L \circ W_L \circ f'_{L-1} \circ W_{L-1} \dots f'_1 \circ W_1$$

$$\nabla_x E = W_1^T f'_1 \dots \circ W_{L-1}^T f'_{L-1} \circ W_L^T f'_L \circ \frac{dE}{da_L}$$

δ_1 δ_{L-1} : error at level L-1

$$\nabla_{W_k} E = \delta_k a_{k-1}^T$$

$$\delta_k = f'_k \circ W_{k+1}^T \delta_{k+1}$$

Computing δ_k in terms δ_{k+1} avoids duplicate operations.

