# Data Analytics
# EEE 4774 & 6777

## Module 4 - Classification

Perceptron, Multilayer Perceptron (MLP) / Artificial Neural Network

Spring 2022

# Perceptron

$\nabla_{\vec{w}} E_P(w) = \begin{bmatrix} \frac{\partial E_P}{\partial w_1} \\ \vdots \\ \frac{\partial E_P}{\partial w_d} \end{bmatrix} = \begin{bmatrix} -\sum_{n \in M} \phi_n(1) t_n \\ \vdots \\ -\sum_{n \in M} \phi_n(d) t_n \end{bmatrix} = -\sum_{n \in M} \vec{\phi}_n t_n$

$(x_n, t_n)$

$w^T \phi_n = \sum_{i=1}^{d} w_i \phi_n(i)$

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

true label (ground truth)

$x \in \mathbb{R}^d \quad or \quad \phi(x) \in \mathbb{R}^d$

$w \in \mathbb{R}^d$

$\frac{\partial}{\partial w_i} w^T \phi_n = \phi_n(i)$

$w^T x$

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

predicted label

parameter vector

# misclassified instances

$\phi(x_n) = \phi_n$

$y(x_n) = y_n$

Loss / Cost/Error function: (**Perceptron criterion**) $\quad \boxed{E_P(w) = -\sum_{n \in M} w^T \phi_n t_n}$

set of misclassified instances

$\mathcal{M}$: set of misclassified instances

$\nabla_w E_P(w) = -\sum_{n \in M} \phi_n t_n \stackrel{\text{iteration } i}{\cong} -\phi_i t_i \quad, i \in \mathcal{M}$

Correctly classified: $y_n t_n = 1 > 0$

Misclassified : $y_n t_n = -1 < 0$

step size

**Stochastic gradient descent:**

$w^{(i+1)} = w^{(i)} - \eta \nabla E_P(w) \cong \boxed{w^{(i)} + \eta \phi_n t_n}$

Gradient descent

Stochastic Gradient descent

feature vector i.e. i data point

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

learning rate

**Perceptron procedure:** For each instance $x_n$

$f_n = \begin{cases} +1, \text{ red} \\ -1, \text{ blue} \end{cases}$

- Initialize $w$ → e.g. $w^{(0)} = \begin{pmatrix} -0.4 \\ -0.8 \end{pmatrix}$

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

- Decide

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

threshold

$w^T x < 0$ e.g. blue $w^{(0)}T x = -10$

$w^T x = 0$

$w^T x > 0$ red

$w^{(1)}$

$w^{(0)}$

red

e.g. $w^{(1)}T x = -5$

$w^T x < 0$ blue

$w^T x > 0$ red

$w^{(1)}$

$w^{(2)}$

$w^{(3)}$

- At each step error for the considered point decreases

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

*(handwritten annotations: $w^{(i)} + \phi_n t_n$ ; error at iter. i for nth pt. ; $> 0$ ; error at iter. i+1 for nth pt.)*

- 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 *(handwritten: linearly)*
- Does not generalize readily to multi-class
- Fixed basis functions, not adaptive to the input

Handwritten annotations on the figure:
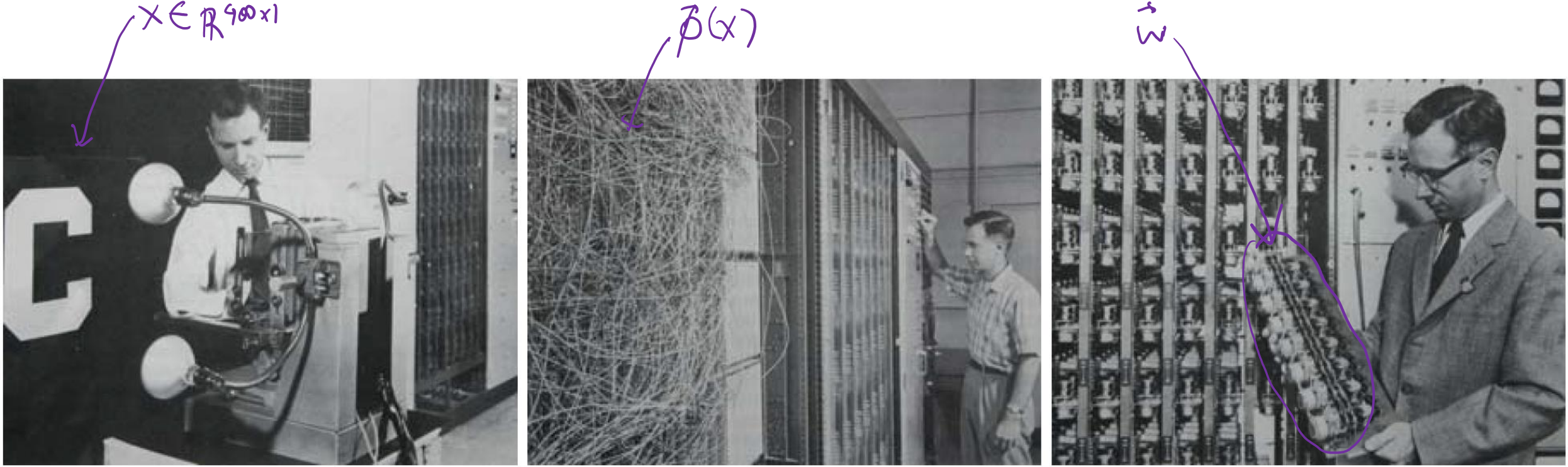- $x \in \mathbb{R}^{400 \times 1}$
- $\beta(x)$
- $\hat{w}$

**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 $20 \times 20$ 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)

*Biological Neural System*

- 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 $\phi(x)$ adaptive to the training data

$$\vec{w}_m^{(1)} = \begin{bmatrix} w_{m1}^{(1)} \\ w_{m2}^{(1)} \\ \vdots \\ w_{mD}^{(1)} \end{bmatrix} \in \mathbb{R}^D \qquad h\left(\vec{w}_m^{(1)T}\vec{x} + b\right) = h\left(w_{m1}^{(1)}x_1 + \cdots + w_{mD}^{(1)}x_D + b\right)$$

*output* (e.g. class label prediction)
$$y_k(x) = f\left(w_{2k}^T \phi(x) + b_{2k}\right) = \sigma\left(w_k^{(2)T}\phi(x) + b_{2k}\right), k = 1, \dots, K$$

*neuron index*  *layer index*  *offset parameter*

$$\phi_m(x) = h\left(w_{1m}^T x + b_{1m}\right) = \sigma\left(w_m^{(1)T}x + b_{1m}\right), \; m = 1, \dots, M$$

*nonlinear activation func.*  *parameter vector*  *input vector slice*

*e.g. image pixels*  *slice 1*  *slice k*

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

Perceptron : One info. processing unit (IPU)

Perceptron $\vec{w}$   $\phi(x) \to$ → Class prediction

inputs

hidden units

$x_D$   $z_M$   $\phi_M$   $w_{KM}^{(2)}$   $y_K$

$w_{MD}^{(1)}$   $w_{1M}^{(2)}$

$x_1$   outputs

$z_1$   $\phi_1$   $w_{10}^{(2)}$   $y_1 = f(w_i^{(2)} + b_{out})$

$x_0$   $z_0$

hidden layer

Neuron firing   Neuron output   stimuli thr.

(Softmax) Logistic Reg.   IPU for class k

Softmax Regression

Log. Reg. $\sigma(w^T x)$

# Neural Network for Classification



- Two hidden units with activation functions

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

- Single output with activation function

$$f\left(\boldsymbol{w}_{2}^{T}\boldsymbol{h}(\boldsymbol{x})\right) = \sigma\left(\boldsymbol{w}_{2}^{T}\boldsymbol{h}(\boldsymbol{x})\right)$$

- 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



$y = x^2$     (annotation: $y = x^2$)

$y = \sin x$

$y = |x|$

*step function*

- $N = 50$ data points

- 1 hidden layer   MLP
- 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

Adjust parameters $\{ w_1^{(1)}, w_2^{(1)}, w_3^{(1)}, w^{(2)}, b_1^{(1)}, b_2^{(1)}, b_3^{(1)}, b^{(2)} \}$ in training

# Artificial Neural Network

An ANN is typically defined by three types of parameters:

$\sigma, \tanh, ReLU, \ldots$

Fixed (Not updated in training)

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

Fully connected, Convolutional (MLP) (CNN)

3. $w_j, b_j$: The weights of the interconnections and offset, which are updated in the learning process.

Supervised Learning:

Training

$x_i \rightarrow O \; y_i = P(\text{class } 1 \mid x_i) = \pi$

$w^T x + b = (w_1 \cdots w_D)\begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} + b$

$= \sum_{i=1D}^{D} w_i x_i + b \cdot 1$

$= \sum_{i=0} \tilde{w}_i \tilde{x}_i$

$\tilde{w} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_D \end{bmatrix} \quad \tilde{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix}$

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

cross entropy loss      squared error loss $\rightarrow$ min MSE

# Learning ANN coefficients

- Mean squared error (MSE): 
$$E(\boldsymbol{w}) = \frac{1}{N} \sum_{i=1}^{N} [\, t_i - y_i(\boldsymbol{w}, \boldsymbol{x_i})\,]^2$$
- MAE, Huber loss, etc. *Cross-Entropy*

- **Gradient Descent** based on **Backpropagation**
  - Efficient technique for computing the gradient of error function $E(\boldsymbol{w})$.

  - Local message passing scheme in which information is sent alternately forward and backward through the network.

  - Given input x, the information flows forward to produce prediction $y$ and scalar cost $E(\boldsymbol{w})$.

  - 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.

*Testing*

$w_1^{(i)}, \dots, w_D^{(i)}, b^{(i)}$

*fixed after training*

*Use input and trained network params. to predict output*

$w_1, \dots, w_D, b$

# Backpropagation Algorithm

Error / Loss Func.

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

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

$$\frac{dE}{da_L} \text{ o } f_L' \, W_L \text{ o} f_{L-1}' \, W_{L-1} \dots f_1' \, W_1$$

$$\nabla_x E = W_1^{\mathrm{T}} f_1' \dots \text{o } W_{L-1}^T f_{L-1}' \text{ o } W_L^T f_L' \text{ o } \frac{dE}{da_L}$$

$\delta_1$

$\delta_{L-1}$: error at level L-1

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

$$\boxed{\delta_k = f_k' \text{ o } W_{k+1}^T \delta_{k+1}}$$

Computing $\delta_k$ in terms $\delta_{k+1}$ avoids duplicate operations.



Forword Propagation 1

Error Estimation 2

Output Layer

Hidden Layer

Input Layer

Backward Propagation 3