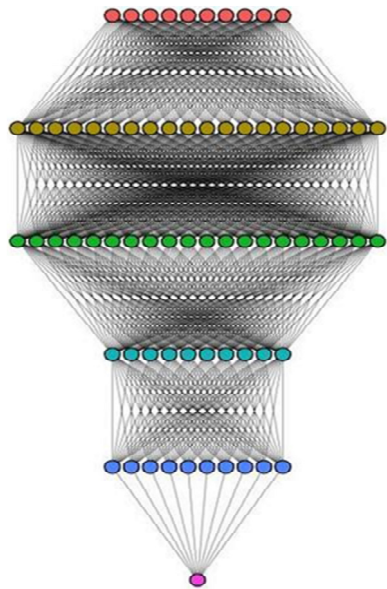


Perceptron



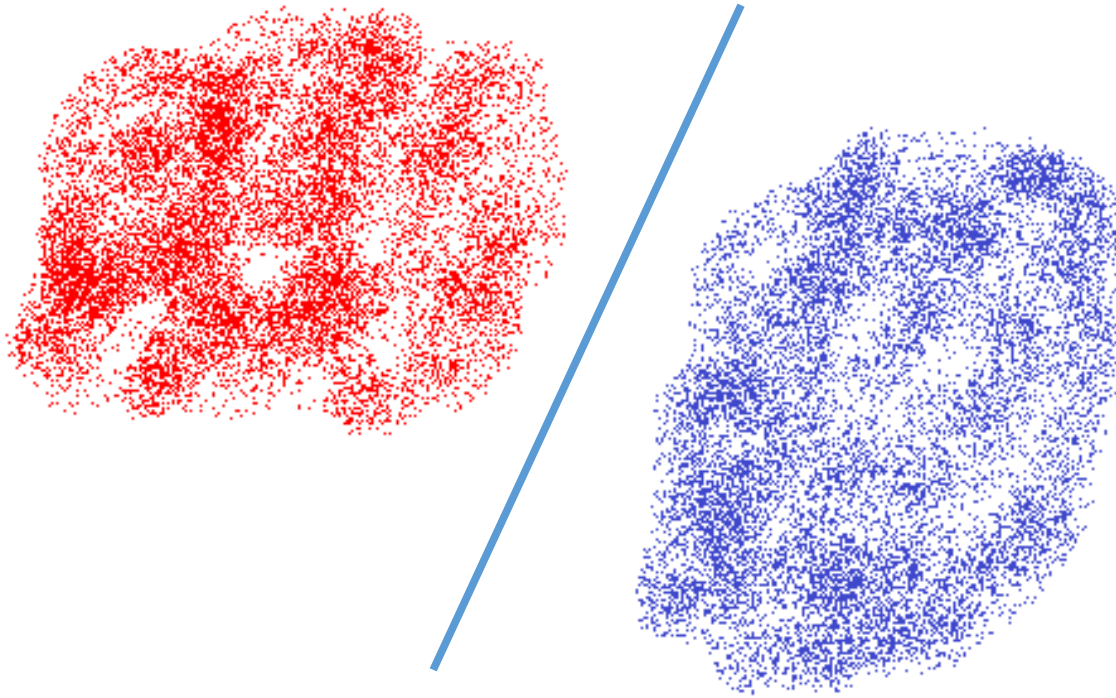
Prithwijit Guha
Dept. of EEE, IIT Guwahati

Supervised Learning



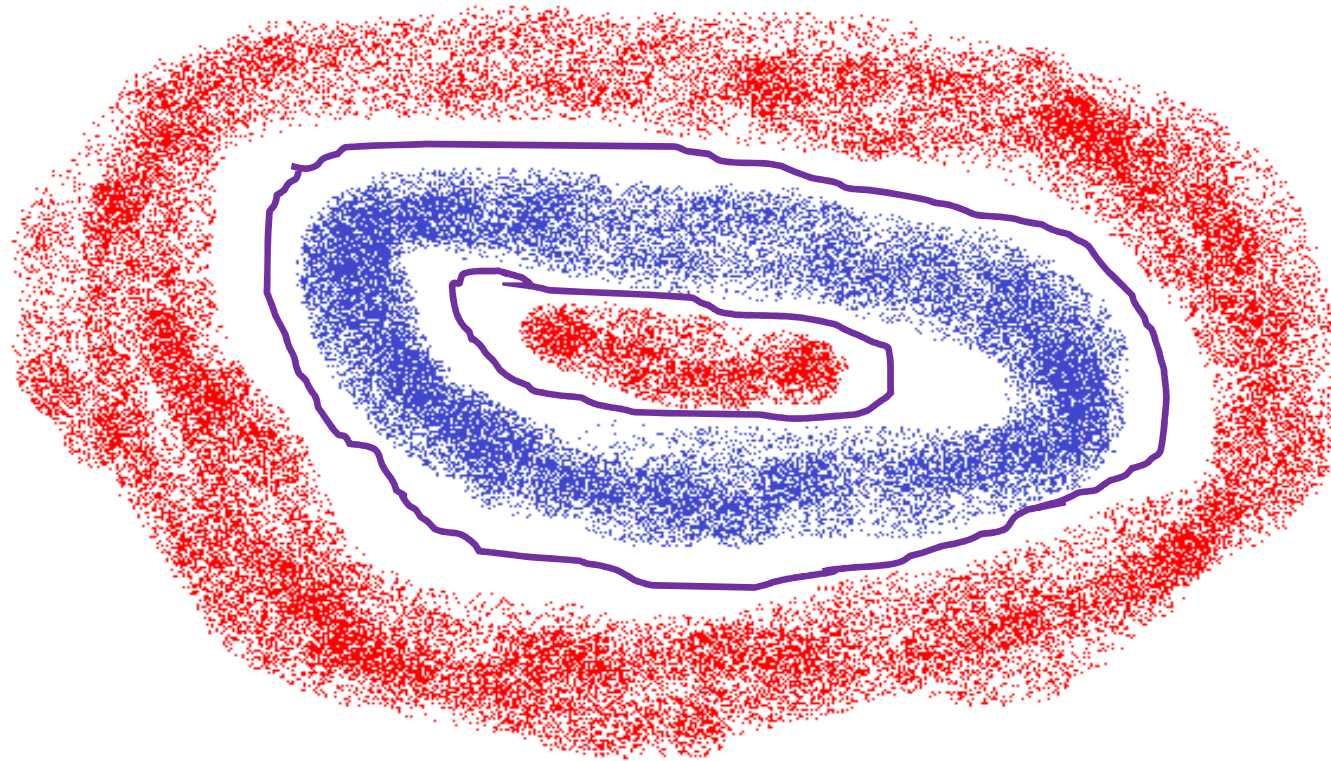
- Bayesian Classification, MAP, Chebyshev Inequality
- Performance Measures, Confusion Matrix, ROC Curves
- Logistic Regression
- Perceptron
- Multi-Layer Perceptron (MLP), ELM
- MLP Architectures, Learning, Interpretations
- Non-parametric Methods and K-NN
- Radial Basis Function Neural Networks
- Data Balancing; SMOTE & Weighted Loss Functions
- Classification & Regression Trees
- Support Vector Machines & Multiple Kernel Learning
- Ensemble Methods, Bagging and Boosting

Linearly Separable



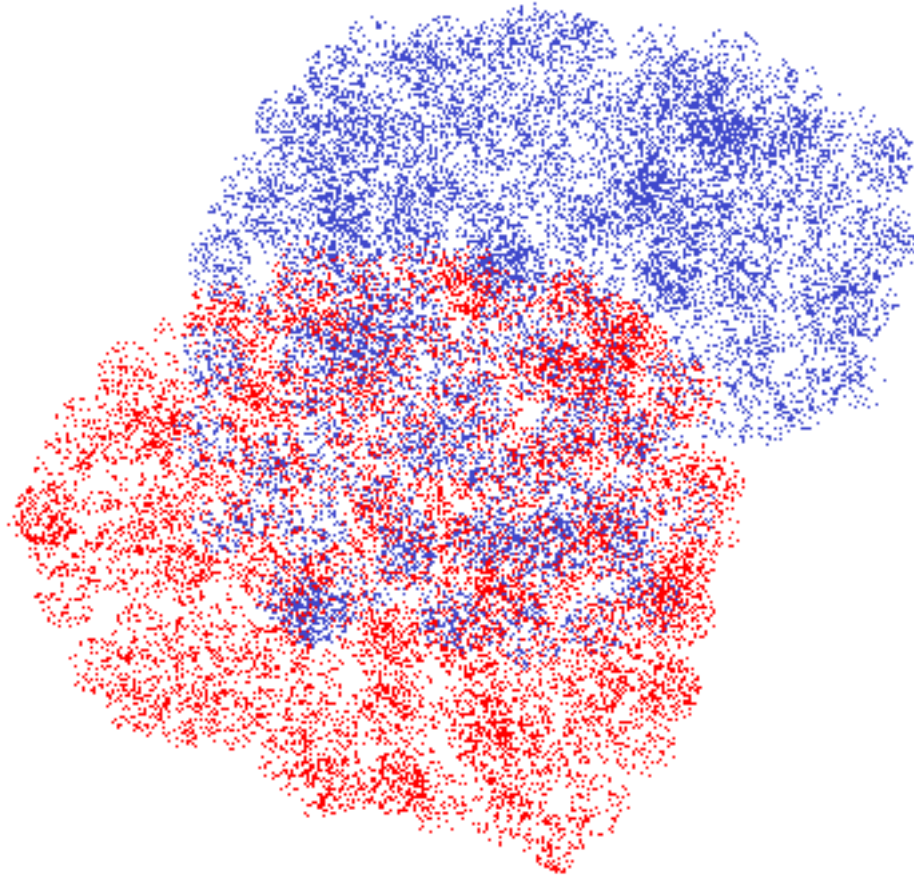
Hyperplane(s) can Separate Features Extracted from
Two (Or More) Classes

Not Linearly Separable



Nonlinear Hypersurfaces Separate Features Extracted from Two
(Or More) Classes

Not Separable



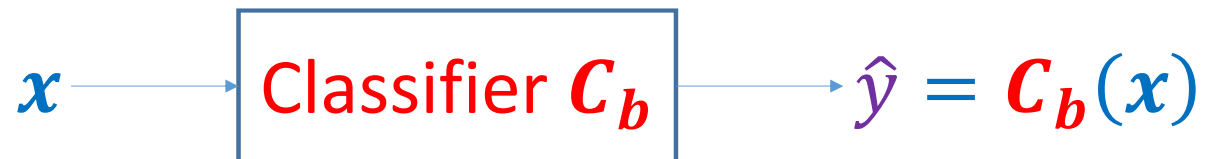
Require Further Transformations for Enhancing Separability

Binary Classification: Input-Output Specification

$$\mathcal{S}_b = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$$y \in \{0,1\}$$



Multi-Category Classification: Input-Output Specification

$$\mathcal{S}_m = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$$


$$\mathbf{x} \in \mathbb{R}^d$$

$$y \in \{1, \dots, m\}$$

Classifier \mathcal{C}_m

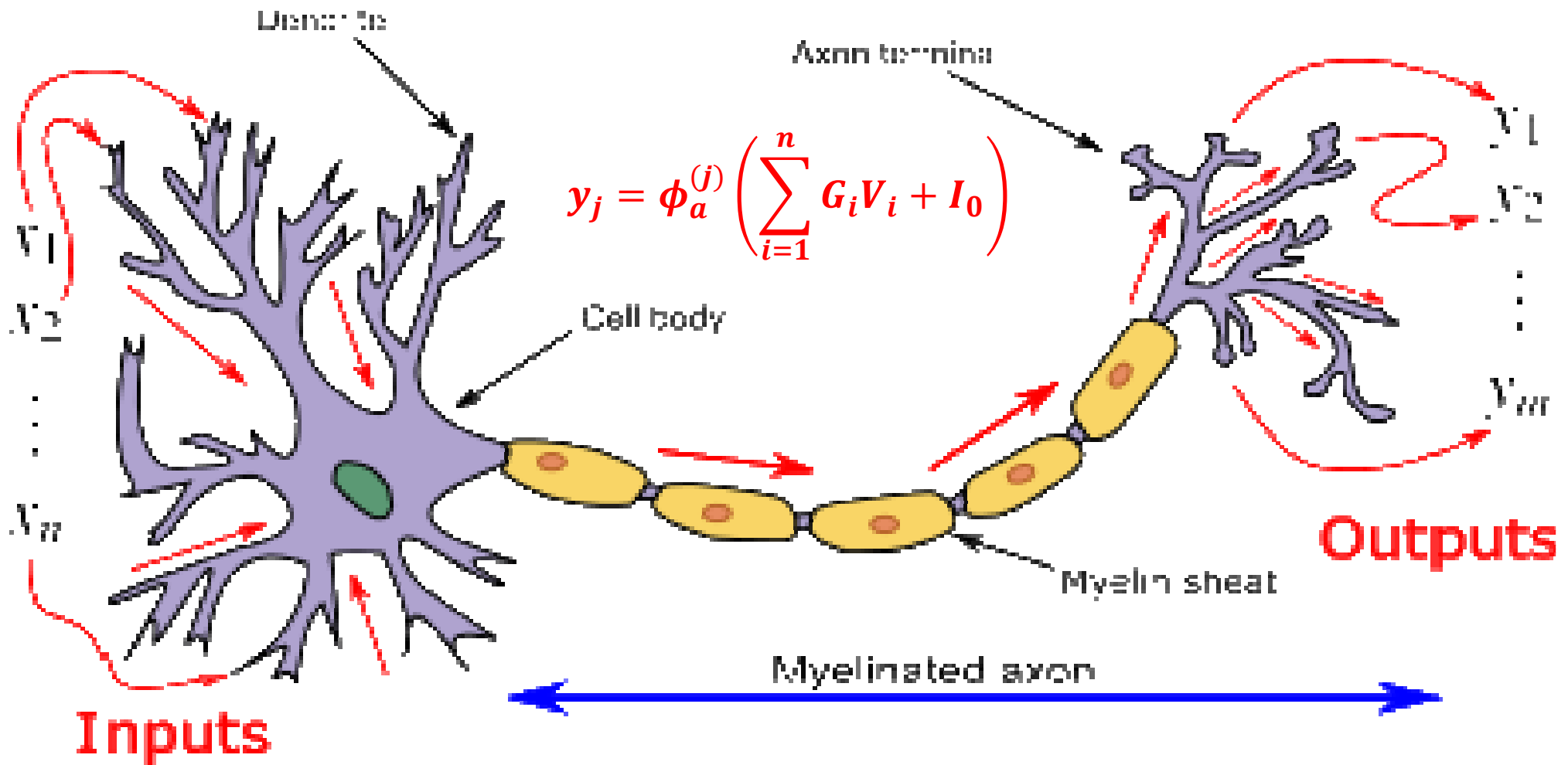


$$\hat{\mathbf{y}} = \mathcal{C}_m(\mathbf{x})$$

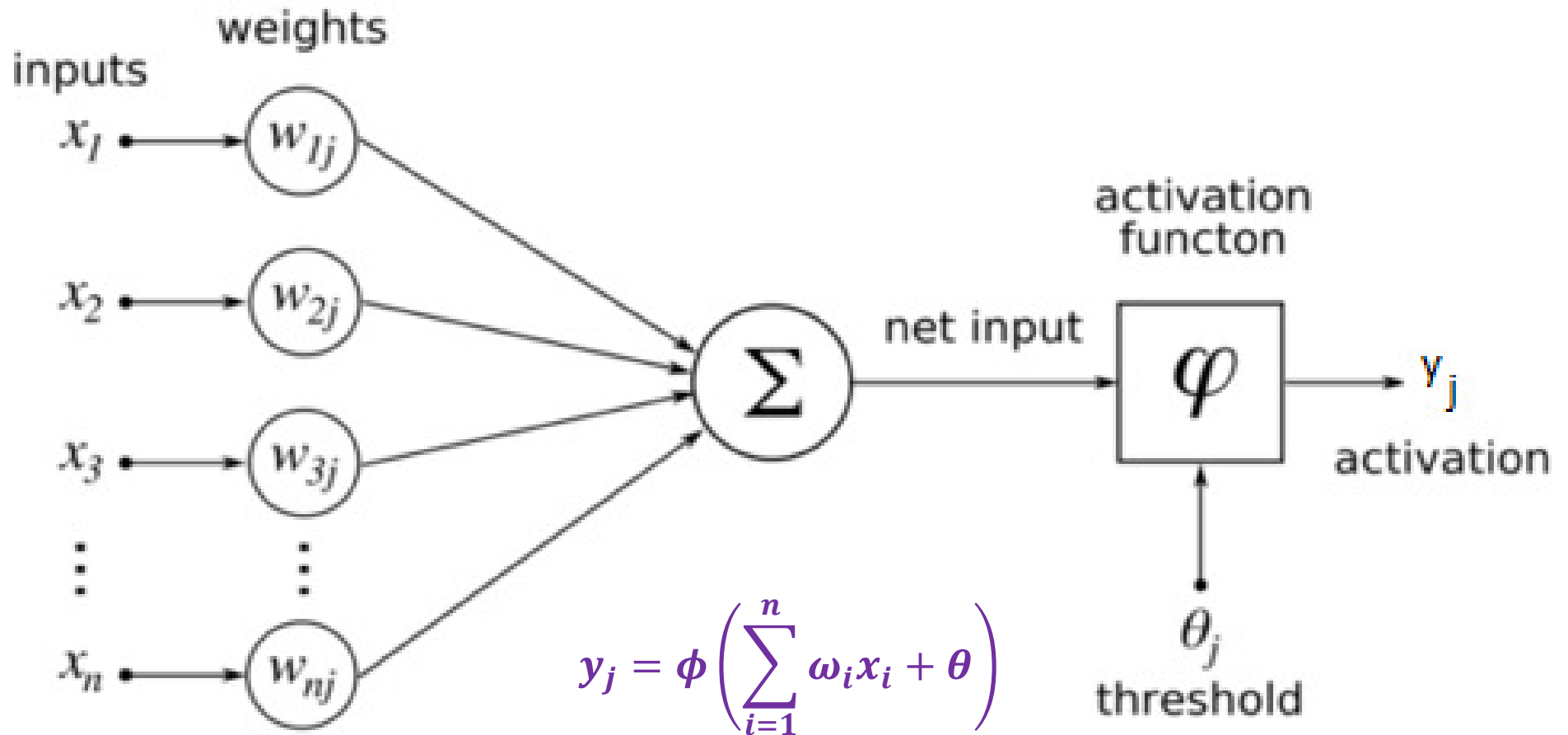
One Hot Output Encoding

$$y_i[j] = \begin{cases} 1, & y_i = j \\ 0, & y_i \neq j \end{cases}$$

Perceptron: Biological Motivation



Perceptron: Mathematical Formulation



Perceptron: Classifier & Feature Extractor

$$\mathcal{S}_b = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$$

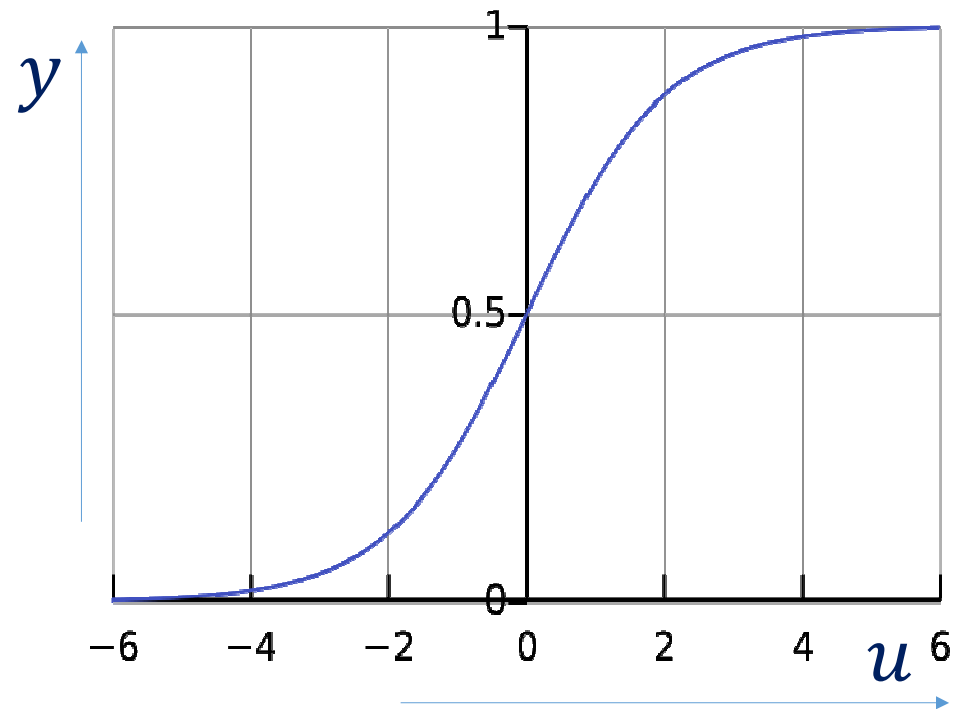
$$u_i = \boldsymbol{\omega}^T \mathbf{x}_i + b$$

$$u_i = \sum_{j=1}^d \boldsymbol{\omega}_j \mathbf{x}_{ij} + b$$

$$\hat{y}_i = \frac{1}{1 + e^{-u_i}}$$

$$\hat{y} = P(\mathbf{x}; \boldsymbol{\omega}, b)$$

$$P : \mathbb{R}^d \rightarrow (0,1)$$



Perceptron: Linear Algebra Interpretation

$$u(x_1, x_2) < 0$$



$$P([x_1, x_2]; [\omega_1, \omega_2], b) < 0.5$$

$$u(x_1, x_2) = \omega_1 x_1 + \omega_2 x_2 + b = 0$$

$$u(x_1, x_2) > 0$$

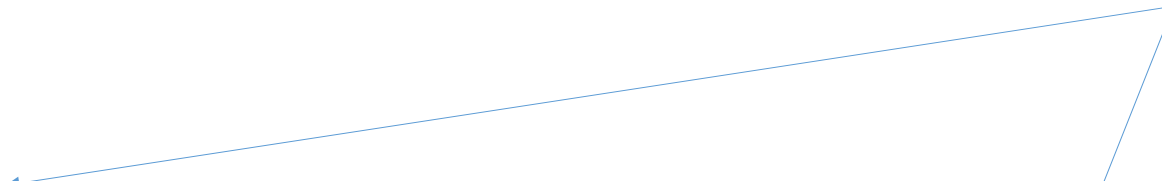



$$P([x_1, x_2]; [\omega_1, \omega_2], b) > 0.5$$

Perceptron: Learning by Gradient Descent

$$\mathcal{S}_b = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$$

$$e_i = y_i - \hat{y}_i = y_i - P(\mathbf{x}_i; \omega, b) \quad \longrightarrow \quad E = \frac{1}{n} \sum_{i=1}^n e_i^2$$


$$\omega_j^{(k+1)} = \omega_j^{(k)} - \eta_{jk} \frac{\partial E}{\partial \omega_j}$$


$$b^{(k+1)} = b^{(k)} - \eta_{jk} \frac{\partial E}{\partial b}$$

Perceptron: Learning by Gradient Descent

$$E = \frac{1}{n} \sum_{i=1}^n e_i^2 \quad \longrightarrow \quad \frac{\partial E}{\partial \omega_j} = \frac{1}{n} \sum_{i=1}^n 2e_i \frac{\partial e_i}{\partial \omega_j}$$

$$e_i = y_i - \hat{y}_i \quad \longrightarrow \quad \frac{\partial e_i}{\partial \omega_j} = -\frac{\partial \hat{y}_i}{\partial \omega_j}$$

$$\frac{\partial \hat{y}_i}{\partial \omega_j} = \frac{\partial \hat{y}_i}{\partial u_i} \times \frac{\partial u_i}{\partial \omega_j}$$

Perceptron: Learning by Gradient Descent

$$\hat{y}_i = \frac{1}{1 + e^{-u_i}} \quad \longrightarrow \quad \frac{\partial \hat{y}_i}{\partial u_i} = \hat{y}_i(1 - \hat{y}_i)$$

$$u_i = \sum_{r=1}^d \omega_r x_{ir} + b \quad \longrightarrow \quad \frac{\partial u_i}{\partial \omega_j} = x_{ij}$$

$$\frac{\partial E}{\partial \omega_j} = \frac{1}{n} \sum_{i=1}^n 2e_i \frac{\partial e_i}{\partial \hat{y}_i} \times \frac{\partial \hat{y}_i}{\partial u_i} \times \frac{\partial u_i}{\partial \omega_j}$$

Perceptron: Learning by Gradient Descent

$$\frac{\partial E}{\partial \omega_j} = -\frac{2}{n} \sum_{i=1}^n e_i \hat{y}_i (1 - \hat{y}_i) \mathbf{x}_{ij}$$

$$\frac{\partial E}{\partial b} = -\frac{2}{n} \sum_{i=1}^n e_i \hat{y}_i (1 - \hat{y}_i)$$

Perceptron: Learning by Pseudo-Inverse

$$\mathcal{S}_b = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$$

$$\begin{array}{ll} y_i = 0 & \longrightarrow \tilde{y}_i^{(0)} = 0 + \epsilon \\ y_i = 1 & \longrightarrow \tilde{y}_i^{(1)} = 1 - \epsilon \end{array} \quad \begin{array}{l} \text{Small} \\ \text{Number} \end{array}$$

$$\begin{array}{l} \text{Example: } \epsilon = 0.001 \\ \begin{array}{l} \nearrow \tilde{y}_i^{(0)} = 0.001 \\ \searrow \tilde{y}_i^{(1)} = 0.999 \end{array} \end{array}$$

Perceptron: Learning by Pseudo-Inverse

$$\tilde{y}_i = \frac{1}{1 + e^{-u_i}} \Rightarrow e^{-u_i} = \frac{1}{\tilde{y}_i} - 1 \Rightarrow e^{u_i} = \frac{\tilde{y}_i}{1 - \tilde{y}_i}$$

$$u_i = \log_e \left(\frac{\tilde{y}_i}{1 - \tilde{y}_i} \right)$$
$$u_i = \boldsymbol{\omega}^T \mathbf{x}_i + b = \log_e \left(\frac{\tilde{y}_i}{1 - \tilde{y}_i} \right)$$

Perceptron: Learning by Pseudo-Inverse

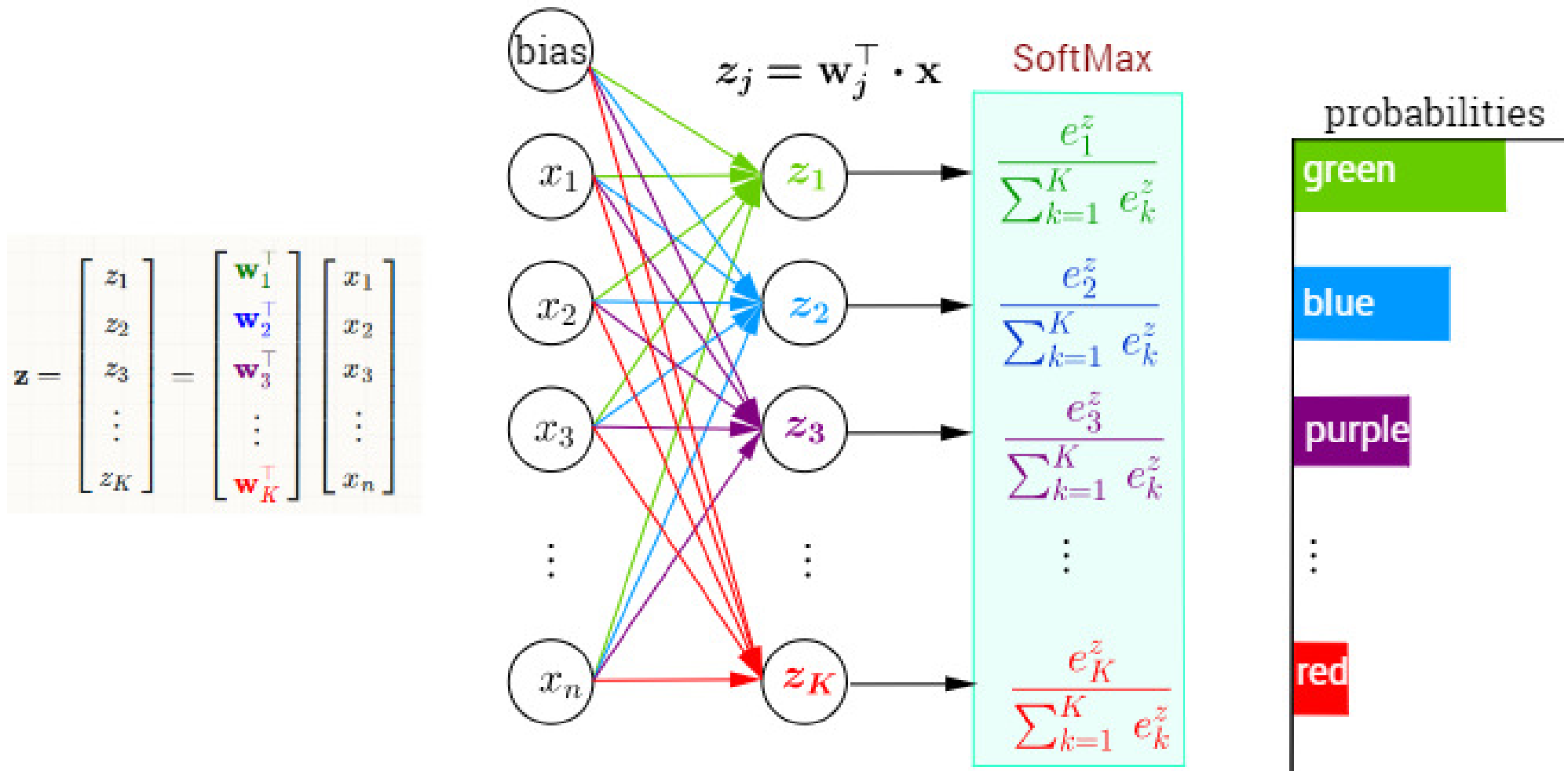
$$A = \begin{bmatrix} \mathbf{x}_1^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_i^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_n^T & 1 \end{bmatrix}_{n \times (d+1)}$$

$$P = \begin{bmatrix} \boldsymbol{\omega} \\ b \end{bmatrix}_{(d+1) \times 1}$$

$$B = \begin{bmatrix} \log_e \left(\frac{\tilde{y}_1}{1 - \tilde{y}_1} \right) \\ \vdots \\ \log_e \left(\frac{\tilde{y}_i}{1 - \tilde{y}_i} \right) \\ \vdots \\ \log_e \left(\frac{\tilde{y}_n}{1 - \tilde{y}_n} \right) \end{bmatrix}_{n \times 1}$$

$$P = (A^T A)^{-1} A^T B = A^\# B$$

Perceptron: Multi-Category Classification



Perceptron: Multi-Category Classification

$$\mathcal{S}_m = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$$

$$\mathbf{x} \in \mathbb{R}^d$$

$$y \in \{1, \dots, m\}$$

$$u_{ij} = \boldsymbol{\omega}_j^T \mathbf{x}_i + b_j$$



$$\mathbf{u}_i = W \mathbf{x}_i + \mathbf{b}$$

$$\hat{y}_{ij} = \frac{e^{u_{ij}}}{\sum_{r=1}^m e^{u_{ir}}}$$

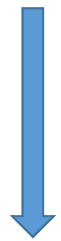
$$W_{m \times d} = \begin{bmatrix} \boldsymbol{\omega}_1^T \\ \vdots \\ \boldsymbol{\omega}_m^T \end{bmatrix}$$

$$\mathbf{b}_{m \times 1} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

Soft-Max Activation Function

Soft-Max Activation Function

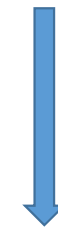
$$\hat{y}_{ij} = \frac{e^{u_{ij}}}{\sum_{r=1}^m e^{u_{ir}}}$$



Soft-Max

$$\hat{y}_{ij} = \frac{1}{1 + \sum_{\substack{r=1 \\ r \neq j}}^m e^{(u_{ir} - u_{ij})}}$$

$$y = \frac{e^u}{1 + e^u}$$



Sigmoid

$$y = \frac{1}{1 + e^{-u}}$$

Multi-Class Perceptron: Learning by Gradient Descent

$$\mathcal{S}_m = \{(\mathbf{x}_i, \mathbf{y}_i); i = 1, \dots, n\}$$

$$\mathbf{x} \in \mathbb{R}^d$$

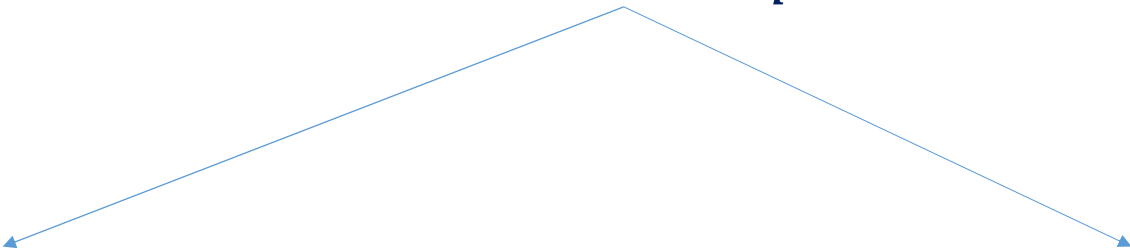
$$\mathbf{y} \in \{0,1\}^m$$

$$u_{iq} = \sum_{r=1}^d \omega_{qr} x_{ir} + b_q \quad \leftarrow q = 1, \dots, m \quad \rightarrow \quad \hat{\mathbf{y}}_{iq} = \phi_q(u_{iq})$$

$$\mathbf{e}_i = \mathbf{y}_i - \hat{\mathbf{y}}_i \quad \longrightarrow \quad \mathbf{e}_{iq} = \mathbf{y}_{iq} - \hat{\mathbf{y}}_{iq}$$

Multi-Class Perceptron: Learning by Gradient Descent

$$e_{iq} = \mathbf{y}_{iq} - \hat{\mathbf{y}}_{iq}$$

$$E = \frac{1}{m \times n} \sum_{i=1}^n \sum_{q=1}^m e_{iq}^2$$


$$\omega_{jk}^{(t+1)} = \omega_{jk}^{(t)} - \eta_{jk}^{(t)} \frac{\partial E}{\partial \omega_{jk}}$$

$$b_j^{(t+1)} = b_j^{(t)} - \eta_j^{(t)} \frac{\partial E}{\partial b_j}$$

Multi-Class Perceptron: Learning by Gradient Descent

$$E = \frac{1}{mn} \sum_{i=1}^n \sum_{q=1}^m e_{iq}^2 \quad \longrightarrow \quad \frac{\partial E}{\partial \omega_{jk}} = \frac{1}{mn} \sum_{i=1}^n 2e_{ij} \frac{\partial e_{ij}}{\partial \omega_{jk}}$$

$$e_{ij} = y_{ij} - \hat{y}_{ij} \quad \longrightarrow \quad \frac{\partial e_{ij}}{\partial \omega_{jk}} = -\frac{\partial \hat{y}_{ij}}{\partial \omega_{jk}}$$

$$\frac{\partial \hat{y}_{ij}}{\partial \omega_{jk}} = \frac{\partial \hat{y}_{ij}}{\partial u_{ij}} \times \frac{\partial u_{ij}}{\partial \omega_{jk}}$$

Multi-Class Perceptron: Learning by Gradient Descent

$$\hat{\mathbf{y}}_{ij} = \phi_j(u_{ij}) \quad \longrightarrow \quad \frac{\partial \hat{\mathbf{y}}_{ij}}{\partial u_{ij}} = \phi_j^{(1)}(u_{ij})$$

$$u_{ij} = \sum_{r=1}^d \boldsymbol{\omega}_{jr} \mathbf{x}_{ir} + b \quad \longrightarrow \quad \frac{\partial u_{ij}}{\partial \boldsymbol{\omega}_{jk}} = \mathbf{x}_{ik}$$

$$\frac{\partial E}{\partial \boldsymbol{\omega}_{jk}} = \frac{1}{mn} \sum_{i=1}^n 2\mathbf{e}_{ij} \frac{\partial \mathbf{e}_{ij}}{\partial \hat{\mathbf{y}}_{ij}} \times \frac{\partial \hat{\mathbf{y}}_{ij}}{\partial u_{ij}} \times \frac{\partial u_{ij}}{\partial \boldsymbol{\omega}_{jk}}$$

Multi-Class Perceptron: Learning by Gradient Descent

$$\frac{\partial E}{\partial \omega_{jk}} = -\frac{2}{mn} \sum_{i=1}^n \mathbf{e}_{ij} \phi_j^{(1)}(u_{ij}) \mathbf{x}_{ik}$$

$$\frac{\partial E}{\partial b_j} = -\frac{2}{mn} \sum_{i=1}^n \mathbf{e}_{ij} \phi_j^{(1)}(u_{ij})$$

Multi-Class Perceptron: Learning by Pseudo-Inverse

$$\mathcal{S}_m = \{(\mathbf{x}_i, y_i); i = 1, \dots, n\} \quad \mathbf{x} \in \mathbb{R}^d \quad y \in \{1, \dots, m\}$$

$$\tilde{\mathbf{y}}_i[j] = \begin{cases} 1 - \epsilon, & y_i = j \\ \frac{\epsilon}{m - 1}, & y_i \neq j \end{cases} \quad \epsilon : \text{Small Number}$$

Example: $m = 3; y_i = 2; \epsilon = 0.001 \longrightarrow \tilde{\mathbf{y}}_i = \begin{bmatrix} 0.0005 \\ 0.999 \\ 0.0005 \end{bmatrix}$

Multi-Class Perceptron: Learning by Pseudo-Inverse

$$\tilde{y}_{ij} = \frac{1}{1 + e^{-u_{ij}}} \Rightarrow e^{-u_{ij}} = \frac{1}{\tilde{y}_{ij}} - 1 \Rightarrow e^{u_{ij}} = \frac{\tilde{y}_{ij}}{1 - \tilde{y}_{ij}}$$

$$u_{ij} = \log_e \left(\frac{\tilde{y}_{ij}}{1 - \tilde{y}_{ij}} \right) \Rightarrow W \mathbf{x}_i + \mathbf{b} = \mathbf{u}_i$$

$$u_{ij} = \boldsymbol{\omega}_j^T \mathbf{x}_i + b_j$$

Multi-Class Perceptron: Learning by Pseudo-Inverse

$$P = [W \quad \mathbf{b}]_{m \times (d+1)}$$

$$B = [\mathbf{u}_1 \quad \cdots \quad \mathbf{u}_i \quad \cdots \quad \mathbf{u}_n]_{m \times n}$$

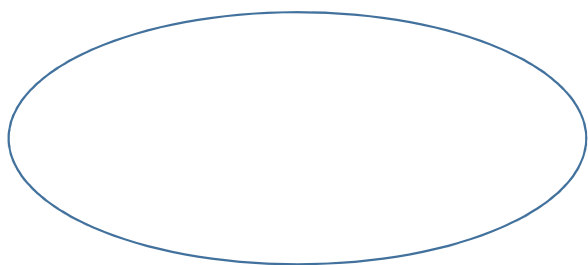
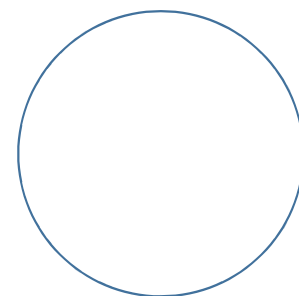
$$A = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_i & \cdots & \mathbf{x}_n \\ 1 & \cdots & 1 & \cdots & 1 \end{bmatrix}_{(d+1) \times n}$$

$$P = BA^T(AA^T)^{-1}$$

Moore-Penrose Generalized Pseudo-Inverse

Perceptron: Input Normalization

$$(x - c_x)^2 + (y - c_y)^2 = r^2$$



$$\left(\frac{x - c_x}{s_x}\right)^2 + \left(\frac{y - c_y}{s_y}\right)^2 = 1$$

Perceptron: Input Normalization

$$\tilde{x}_i = Q^T (x_i - \mu)$$

Matrix of Eigen Vectors of
Covariance Matrix

A diagram illustrating the components of the input normalization equation. The equation $\tilde{x}_i = Q^T (x_i - \mu)$ is centered at the top. Below it, two arrows point from descriptive text to specific parts of the equation. The first arrow, originating from the text 'Matrix of Eigen Vectors of Covariance Matrix', points to the Q^T term. The second arrow, originating from the text 'Mean Vector', points to the μ term.

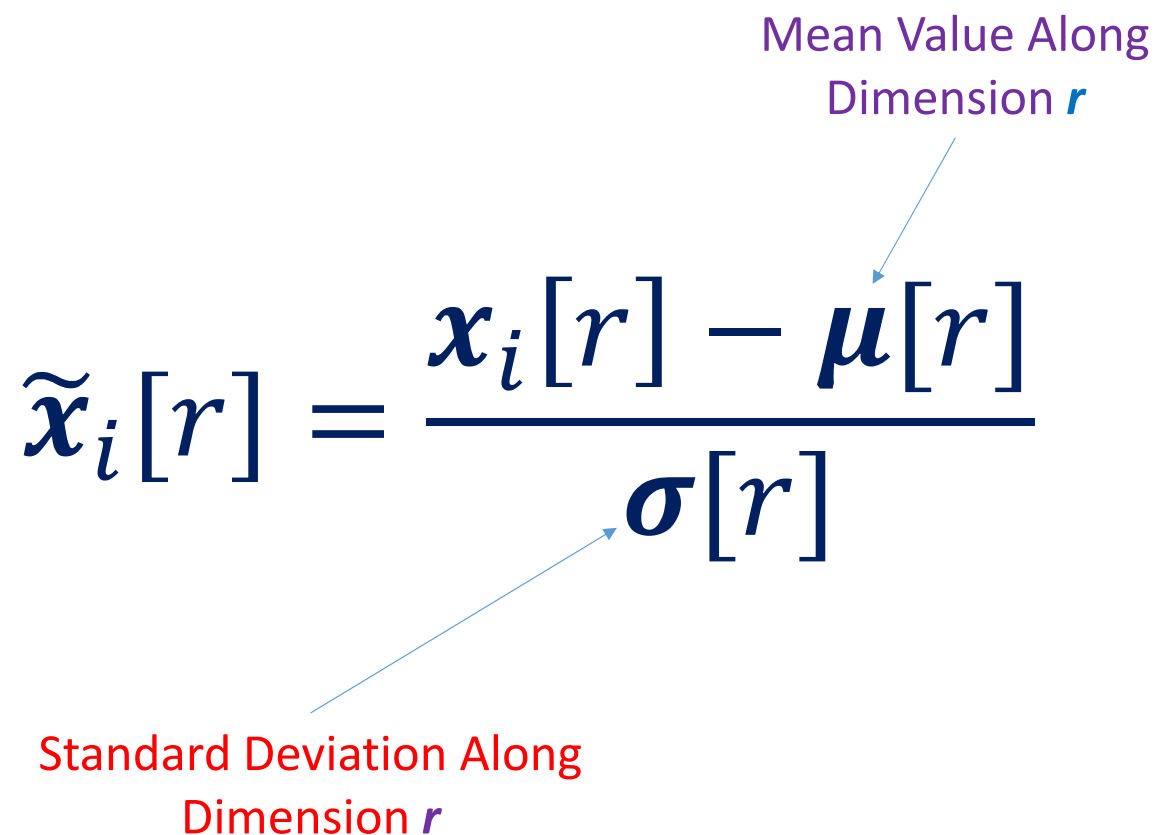
Mean Vector

Perceptron: Input Normalization

$$\tilde{x}_i[r] = \frac{x_i[r] - \mu[r]}{\sigma[r]}$$

Mean Value Along Dimension r

Standard Deviation Along Dimension r

The diagram illustrates the formula for input normalization in a perceptron. The formula is $\tilde{x}_i[r] = \frac{x_i[r] - \mu[r]}{\sigma[r]}$. A purple label 'Mean Value Along Dimension r ' with a blue arrow points to the $\mu[r]$ term in the numerator. A red label 'Standard Deviation Along Dimension r ' with a blue arrow points to the $\sigma[r]$ term in the denominator.

Perceptron: Input Normalization

$$\tilde{x}_i[r] = \frac{x_i[r] - x_{min}[r]}{x_{max}[r] - x_{min}[r]}$$

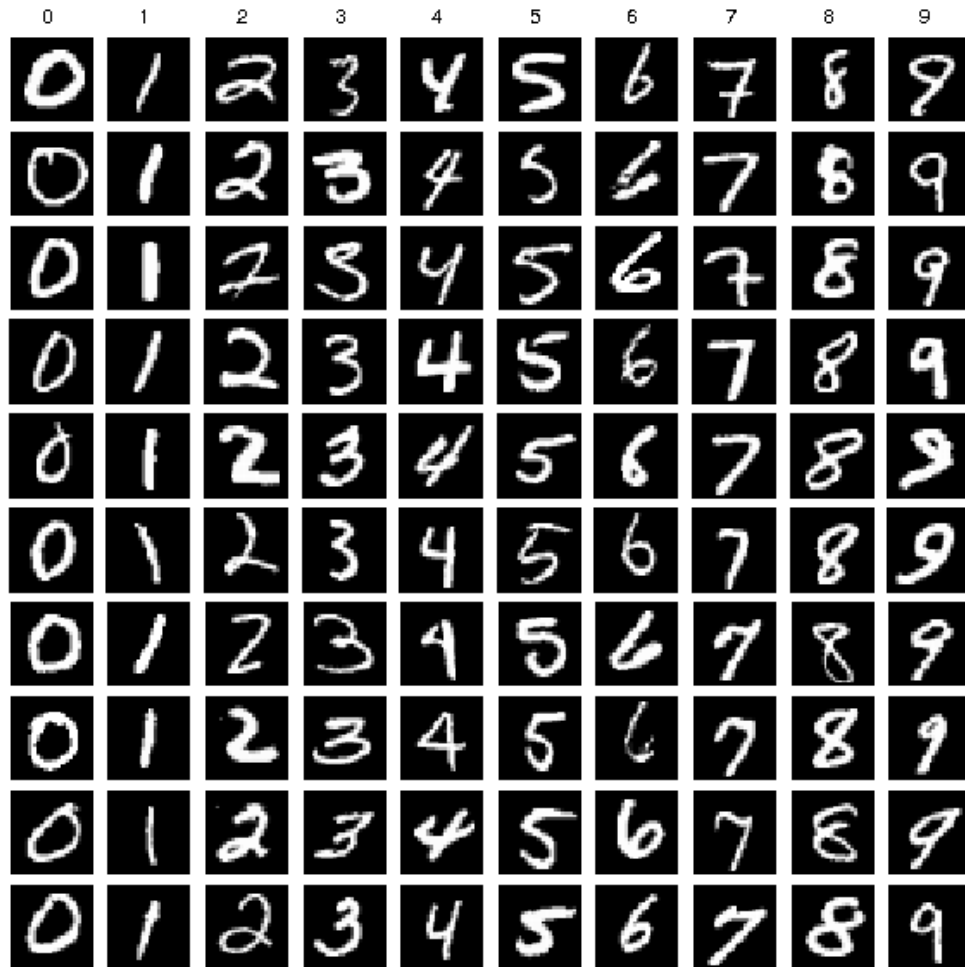
Diagram illustrating the input normalization formula for a Perceptron. The formula is:

$$\tilde{x}_i[r] = \frac{x_i[r] - x_{min}[r]}{x_{max}[r] - x_{min}[r]}$$

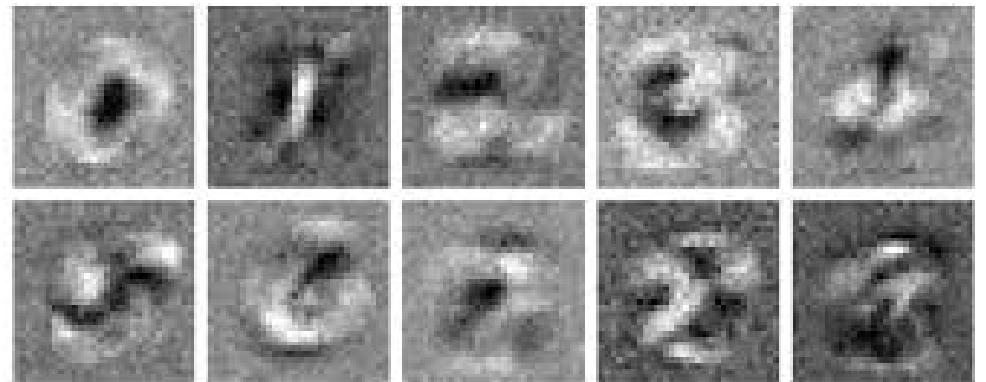
The diagram includes two annotations with arrows pointing to the formula:

- An arrow points from the text "Maximum Value Along Dimension r " to the term $x_{max}[r]$ in the denominator.
- An arrow points from the text "Minimum Value Along Dimension r " to the term $x_{min}[r]$ in the denominator.

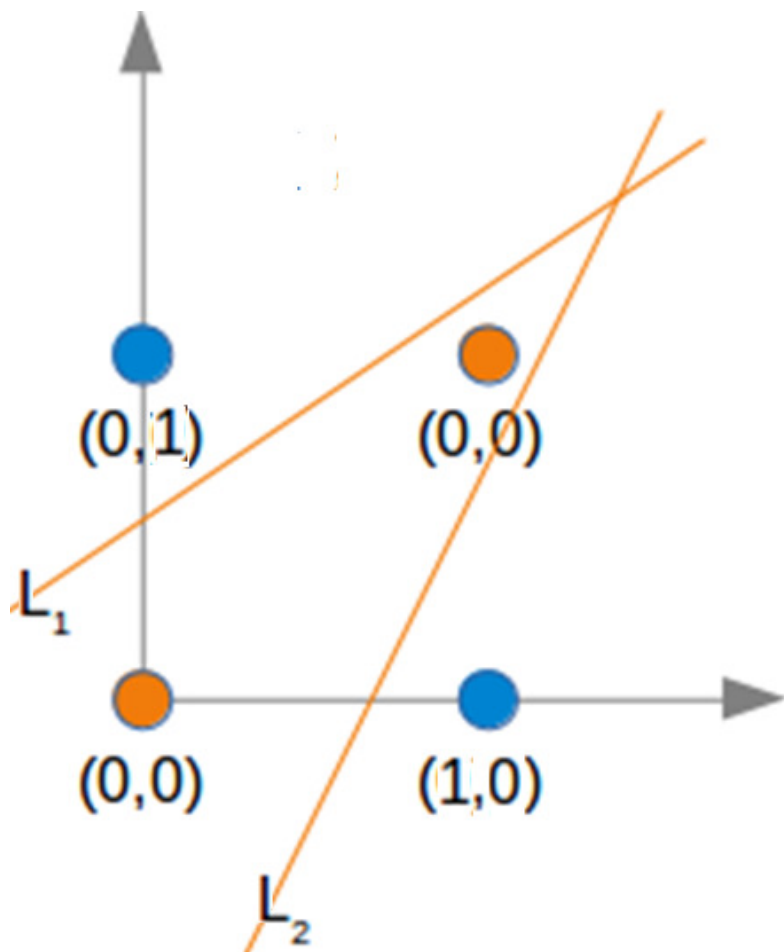
Application to Handwritten Digit Classification



The MNIST Dataset



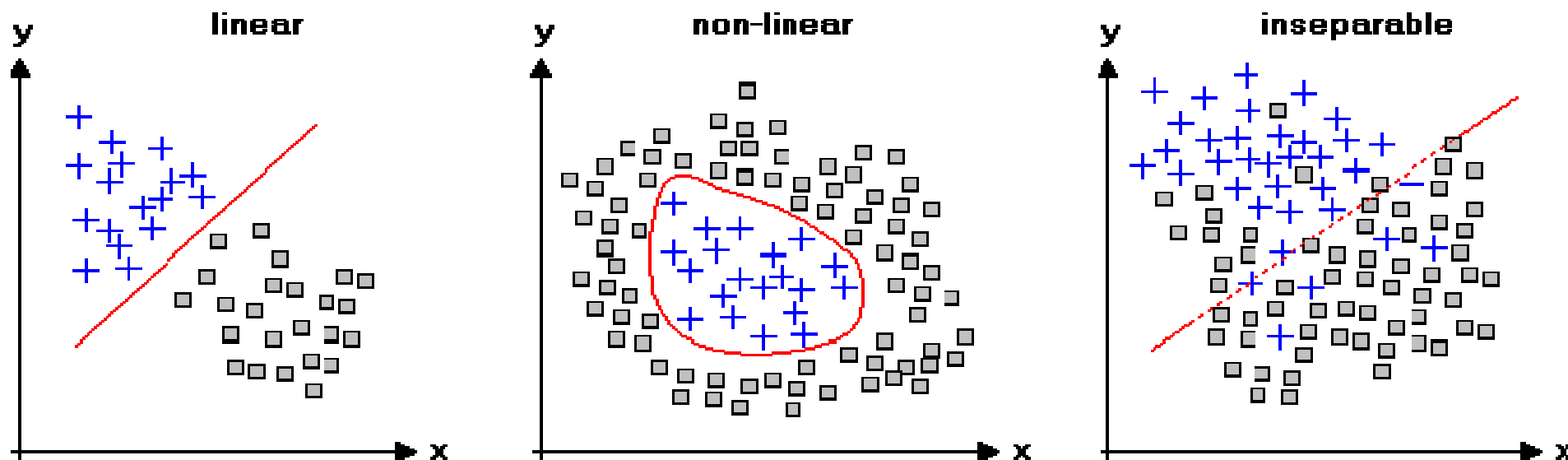
Perceptron: Failure Cases



The XOR Problem

$$y = x_1(1 - x_2) + x_2(1 - x_1)$$

Perceptron: Failure Cases



Covers Theorem of Separability

Nonlinear Transformation to a Higher Dimensional Space Increases the Probability of Linear Separability

Summary

- Separability of Data in Feature Space
- Input-Output Description for Classification Problems
- Perceptron: Motivation & Formulation
- Perceptron: Interpretations
- Perceptron: Learning from Data
- Perceptron: Multi-category Classification
- Perceptron: Applications
- Failure Cases: Towards Multi-Layered Perceptron



Thank You