

Supervised Learning Methods: SVM, Neural Networks

ECE/CS 498 DS U/G

Lecture 20

Ravi K. Iyer

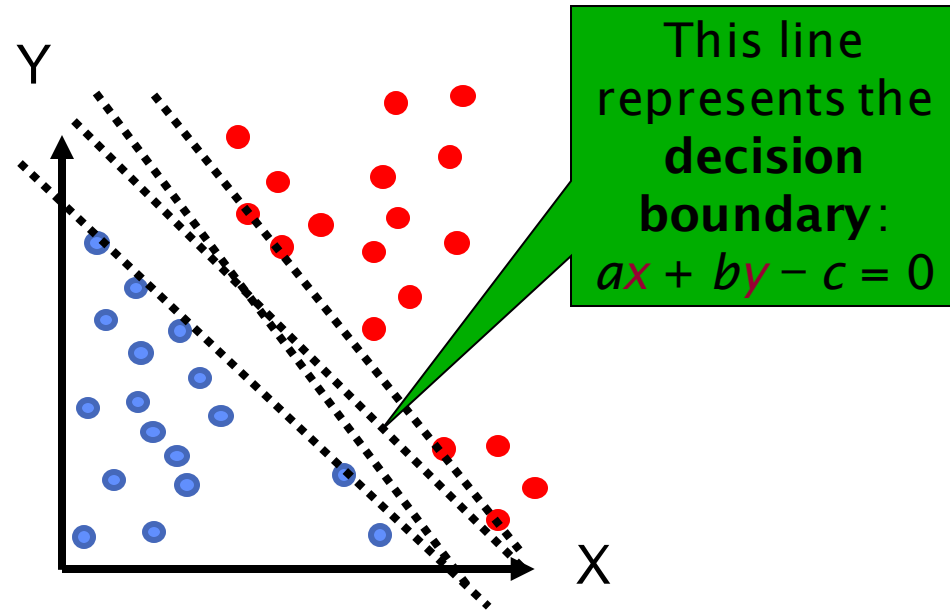
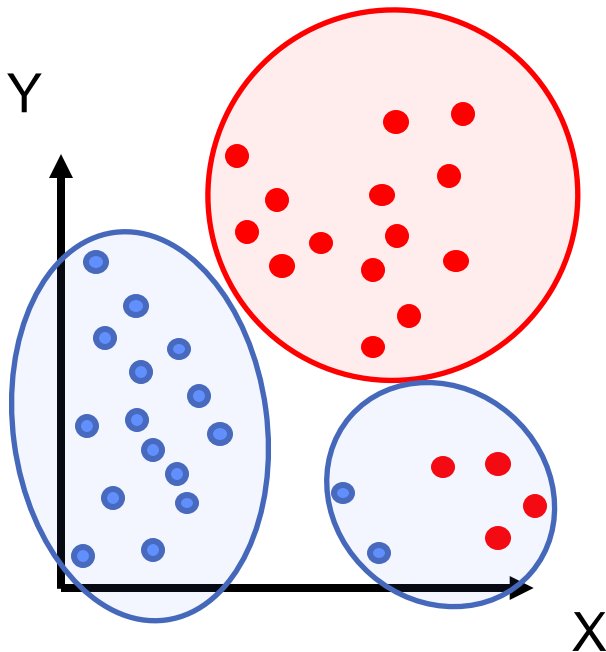
Dept. of Electrical and Computer Engineering
University of Illinois at Urbana Champaign

Announcements

- MP3 Checkpoint 2 is due on Wednesday, Apr 17
- HW 4 on Factor graphs and HMM will be released on Wednesday, Apr 17

Unsupervised to Supervised Learning

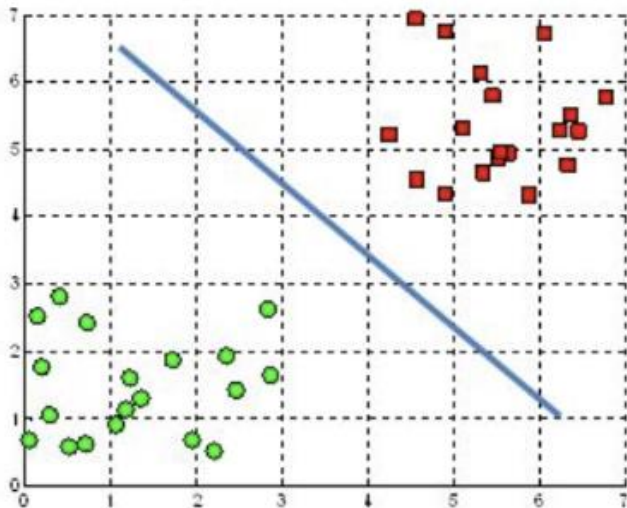
Can we do better? **Yes**



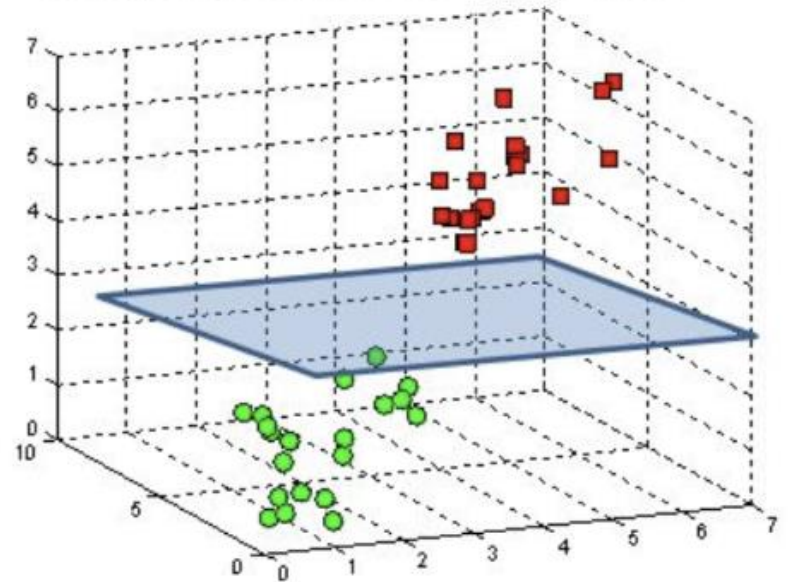
If in a classification task the data is linearly separable, a decision boundary with the same dimensionality as the data can be used to separate the data

Hyperplane as decision boundary

A hyperplane in \mathbb{R}^2 is a line



A hyperplane in \mathbb{R}^3 is a plane



Example of decision boundary in 2D space (line) and 3D space (plane).

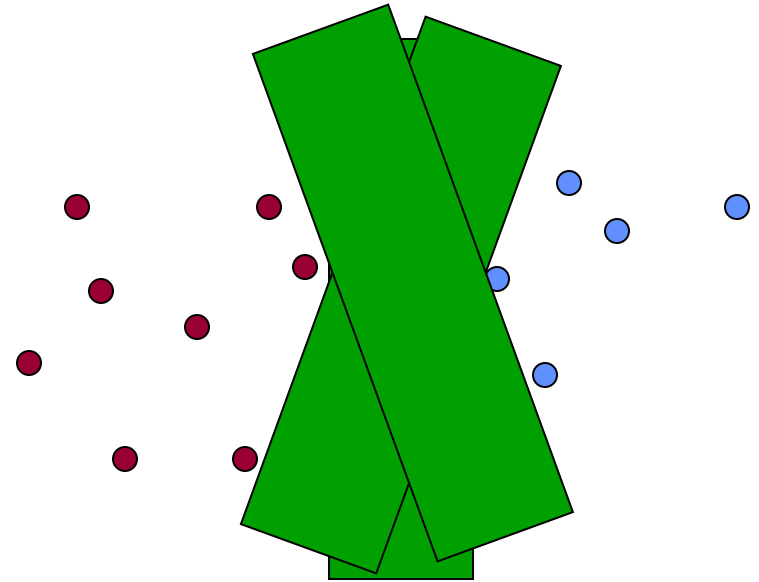
In d dimensional space, hyperplane is given by:

$$a_0 + a_1x_1 + a_2x_2 + \cdots + a_dx_d = 0$$

Image Source: <https://towardsdatascience.com/support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47>

Linear classifiers: Which Hyperplane?

- Lots of possible solutions for a , b , c .
- Some methods find a separating hyperplane, but not the optimal **one [according to some criterion of goodness]**
 - E.g., perceptron
- Support Vector Machine (SVM) finds an optimal* solution.
 - Maximizes the distance between the hyperplane and the “difficult points” close to decision boundary
 - SVMs maximize the *margin* around the separating hyperplane.
 - A.k.a. large margin classifiers
 - The decision function is fully specified by a subset of training samples, *the support vectors*.



If you have to place a fat separator between classes, you have less choices, and so the capacity of the model has been decreased

Linear Support Vector Machine (SVM)

Sec 15.1

Assumption: Data is linearly separable

Let (x_0, y_0) be a point on $w^T x + b = 1$

Then its distance to the separating plane $w^T x + b = 0$ is:

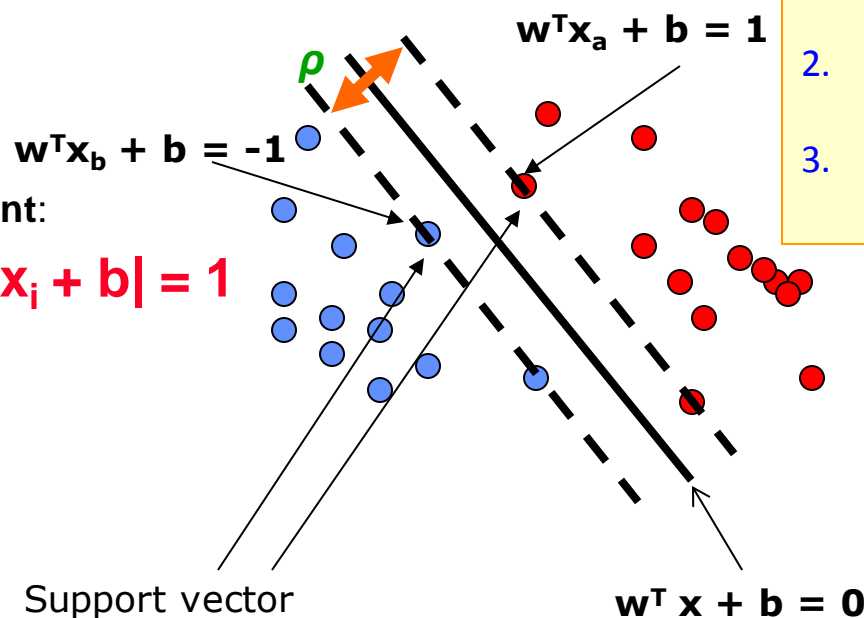
$$|w^T (x_0, y_0) + b| / ||w|| = 1 / ||w||$$

- **Hyperplane**

$$w^T x + b = 0$$

- **Extra scale constraint:**

$$\min_{i=1, \dots, n} |w^T x_i + b| = 1$$



Distance between

$w^T x + b = +1$ and -1 is $\rho = 2 / ||w||$

What we did:

1. Consider all possible w with different angles
2. Scale w such that the constraints are tight
3. Pick the one with largest margin/minimal size

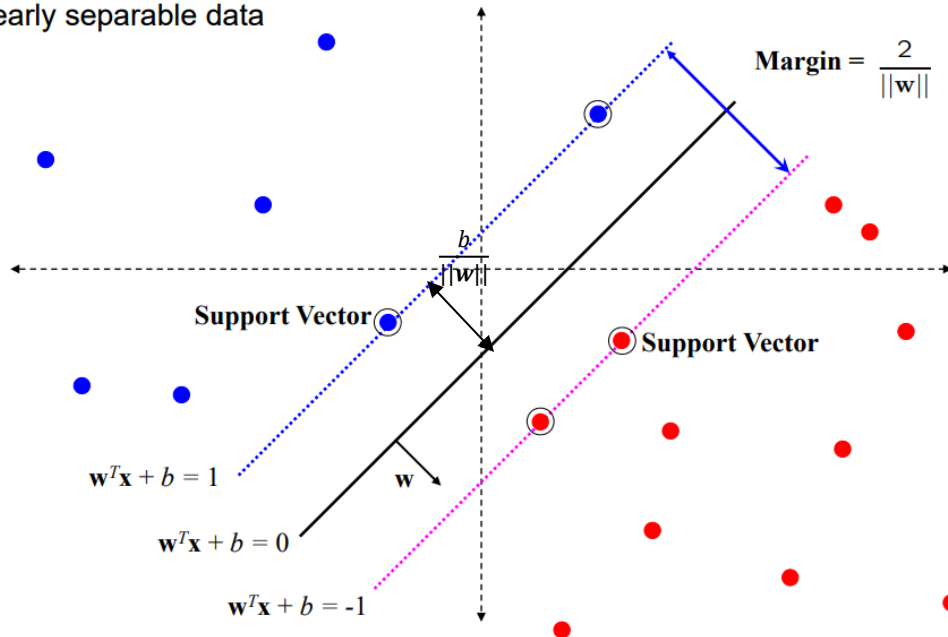
SVM Optimization:

Maximize: ρ or

Minimize: $\frac{1}{2} ||w||^2$

SVM : How to find such hyperplane?

linearly separable data



Each training point: $x_i \in R^d \forall i \in \{1, \dots, n\}$

Label: $y \in \{-1, 1\}$

$w \in R^d$

Maximum margins linear classifier:

$$\operatorname{argmin}_w \frac{1}{2} ||w||^2 \text{ such that } 1 \leq y_i(w^T x_i + b) \quad \forall i \in \{1, \dots, N\}$$

Constraint and Objective:

1. Define the hyperplane via the following rule:

$$y_i(w^T x_i + b) \geq 1$$

-- **Constraint**

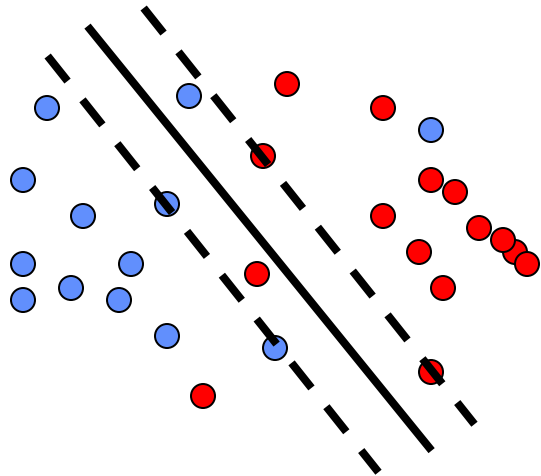
2. Maximize margin distance:

$$\text{Margin distance} = \frac{w^T x + b + 1 - (w^T x + b - 1)}{||w||} = \frac{2}{||w||}$$

$$\max \frac{2}{||w||} \Rightarrow \min ||w|| \Rightarrow \min \frac{||w||^2}{2}$$

-- **Objective**

Dataset with noise



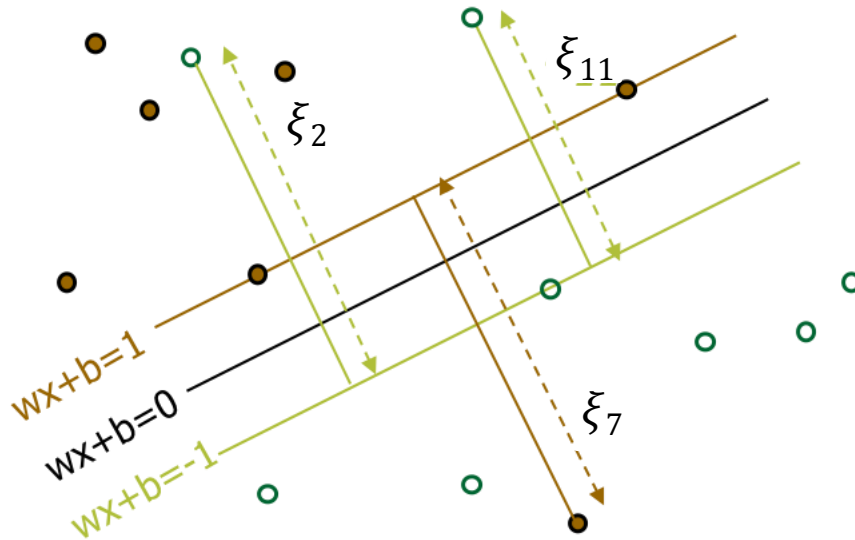
- **Hard Margin:** So far we require all data points be classified correctly
 - No training error
- **What if the training set is noisy?**
 - Solution: Use **Soft Margin**

The two red points and two blue points do not satisfy the criteria $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Soft Margin Classification

Slack variables $\xi_i \geq 0$ can be added to allow misclassification of difficult or noisy examples.

What should our optimization criterion be?



Minimize

$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{j=1}^k \xi_j$$

The slack variables will be zero for the samples which can be classified correctly and a positive value for the remaining points.

Hard Margin v.s. Soft Margin

- **The old formulation:**

Find \mathbf{w} and b such that

$\frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$ is minimized and for all $\{(\mathbf{x}_i, y_i)\}$
 $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

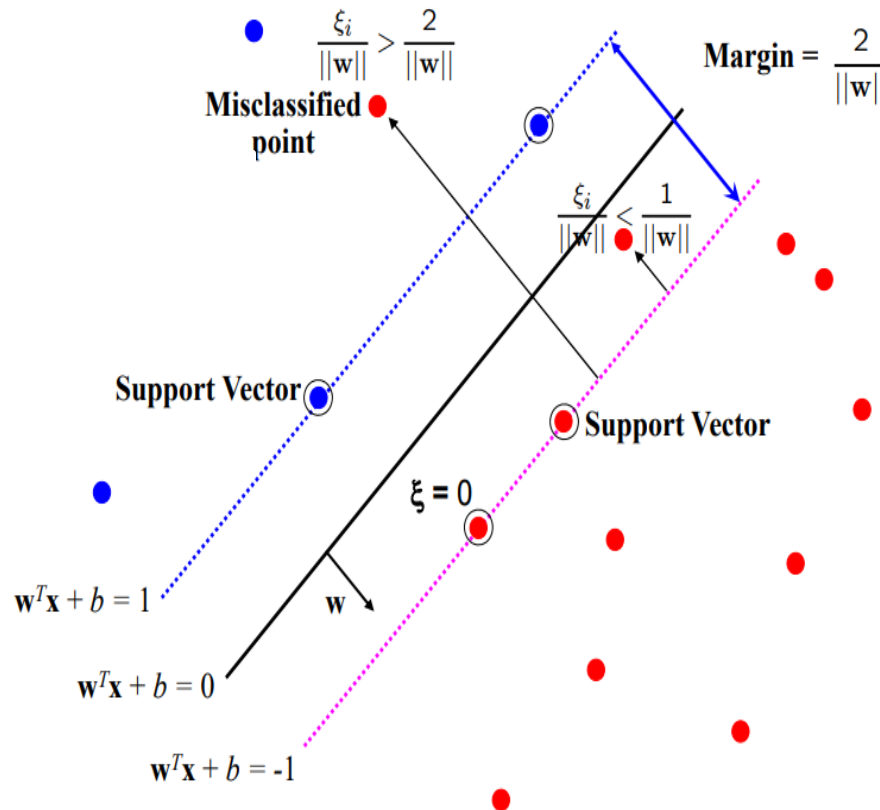
- **The new formulation incorporating slack variables:**

Find \mathbf{w} and b such that

$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i$ is minimized and for all $\{(\mathbf{x}_i, y_i)\}$
 $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$ for all i

- **Parameter C can be viewed as a way to control overfitting.**

SVM: Non-separable case



- $\xi_i \geq 0$
- $0 < \xi_i \leq 1$: point is between margin and correct side of hyper plane
—margin violation
- $\xi_i > 1$: point is on the wrong side of margin
—misclassification

Maximum margins linear classifier:

$$\operatorname{argmin}_{w, \xi_i \geq 0} \frac{1}{2} \|w\|^2 + C \sum_i^N \xi_i \text{ such that } 1 - \xi_i \leq y_i(w^T x_i + b) \quad \forall i \in \{1, \dots, N\}$$

Transformation to data

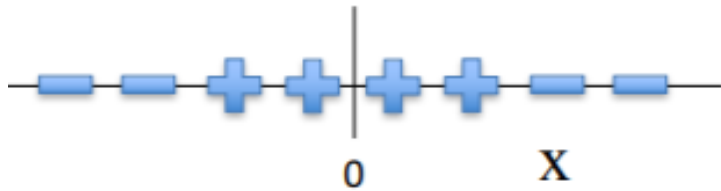


Fig 1

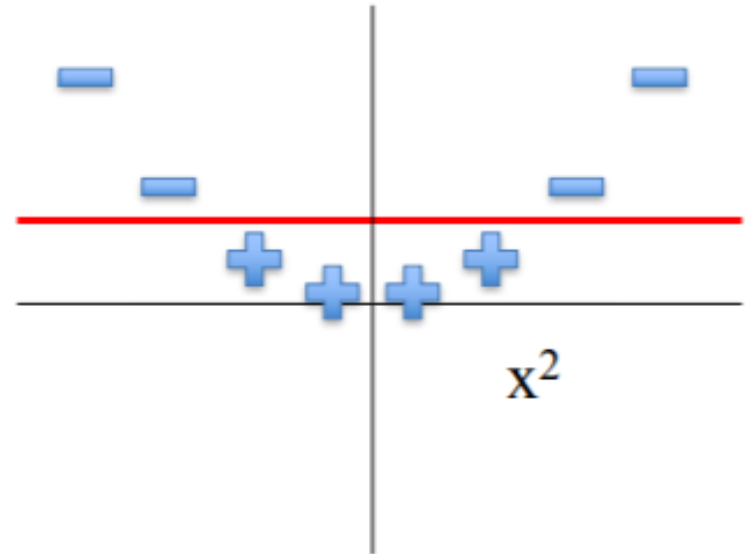
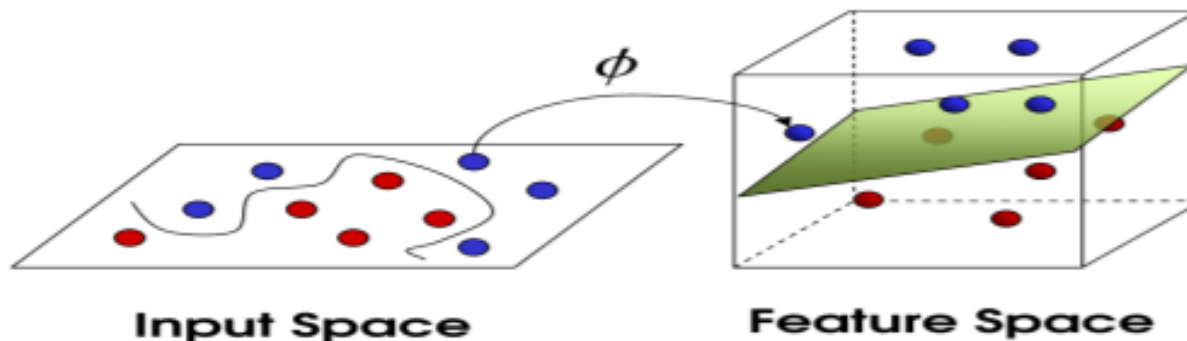


Fig 2

- Data points in Fig 1 cannot be separated using SVM
- Applying **transformation** $\Phi(x) = x^2$ gives data points in Fig 2
 - Can be separated by a line
- Apply SVM on transformed data

Mapping data to new feature space



$$\Phi : \mathcal{X} \mapsto \hat{\mathcal{X}} = \Phi(\mathbf{x})$$

- For example, if $x_i = [x_{i1}, x_{i2}]$, i.e., $\mathbf{x} \in \mathbb{R}^2$,
 $\Phi(\mathbf{x}) = [1, x_{i1}, x_{i2}, x_{i1}x_{i2}, x_{i1}^2, x_{i2}^2]$
- Run SVM on $\Phi(\mathbf{x})$ instead of \mathbf{x}

Kernels

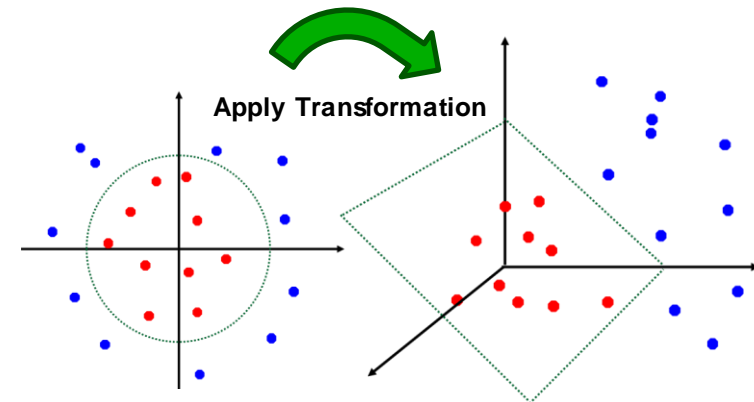
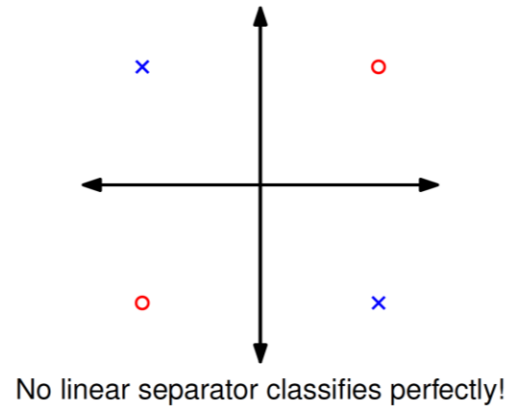
- Define Kernel $K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$
- Use $K(x_i, x_j)$ in SVM
- Because of the way the solution of SVM optimization problem is computed, it is easy to use $K(x_i, x_j)$

Example of commonly used Kernels ($x_i, x_j \in \mathbb{R}$)

- Polynomial kernel of order 2: $K(x_i, x_j) = 1 + x_i + x_j + x_i^2 + x_j^2 + 2x_i x_j$
- Radial Basis Function: $K(x_i, x_j) = \exp\left(-\frac{(x_i - x_j)^2}{2\gamma^2}\right)$

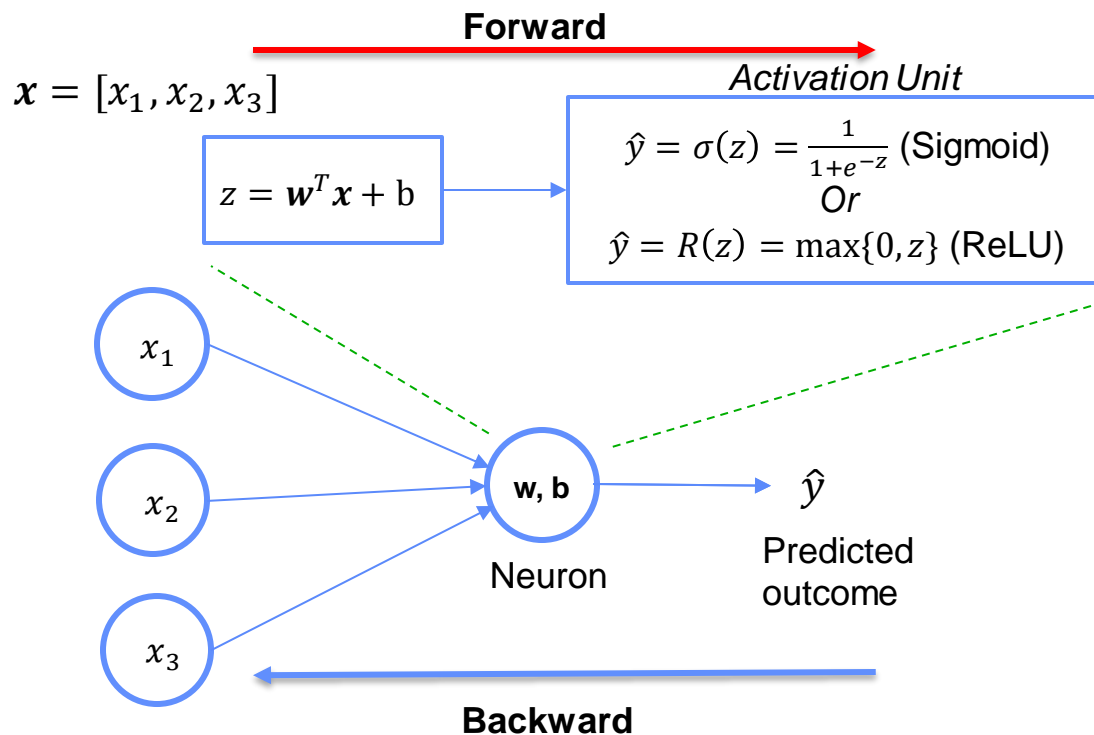
Neural Networks

- Limitations in SVM Models
 - Kernel-SVM can separate data using transformation $\phi(x)$, i.e., x is not linearly separable but $\phi(x)$ is linearly separable; think of $\phi(x)$ as the new feature
 - Identifying the appropriate kernel is difficult
 - Computation time for large problems
- Neural Network
 - Derive features of the data automatically as a constrained optimization problem

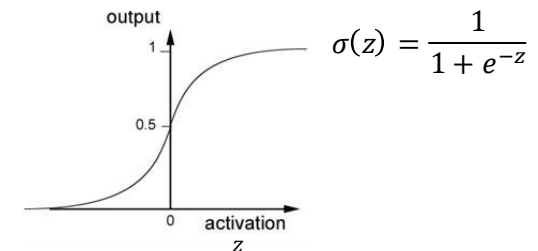


Perceptron Model

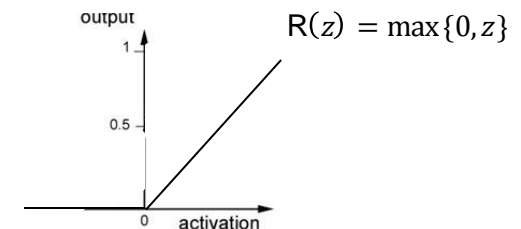
- The core of the neural network is perceptron model



Sigmoid Function



ReLU Function



Update Rule (Backward):

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla J(\mathbf{w}_t)$$

η : Learning rate

Loss

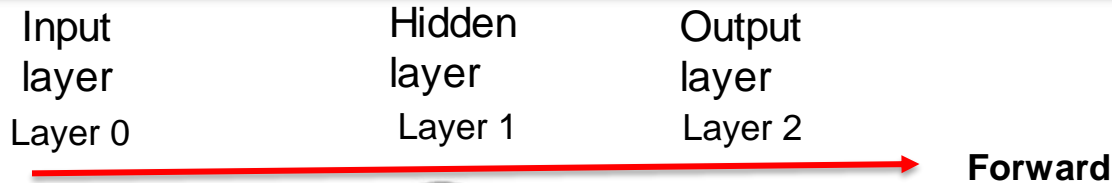
$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{w} \cdot \mathbf{x}^{(i)}, y^{(i)})$$

N : number of samples $\mathbf{x}^{(i)}$: feature of i^{th} sample

Computing Gradient

$$\nabla J(\mathbf{w}_0) = \left(\frac{\partial J(\mathbf{w})}{\partial w_0}, \frac{\partial J(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_n} \right)_{\mathbf{w}_0}$$

Neural Network (Forward)



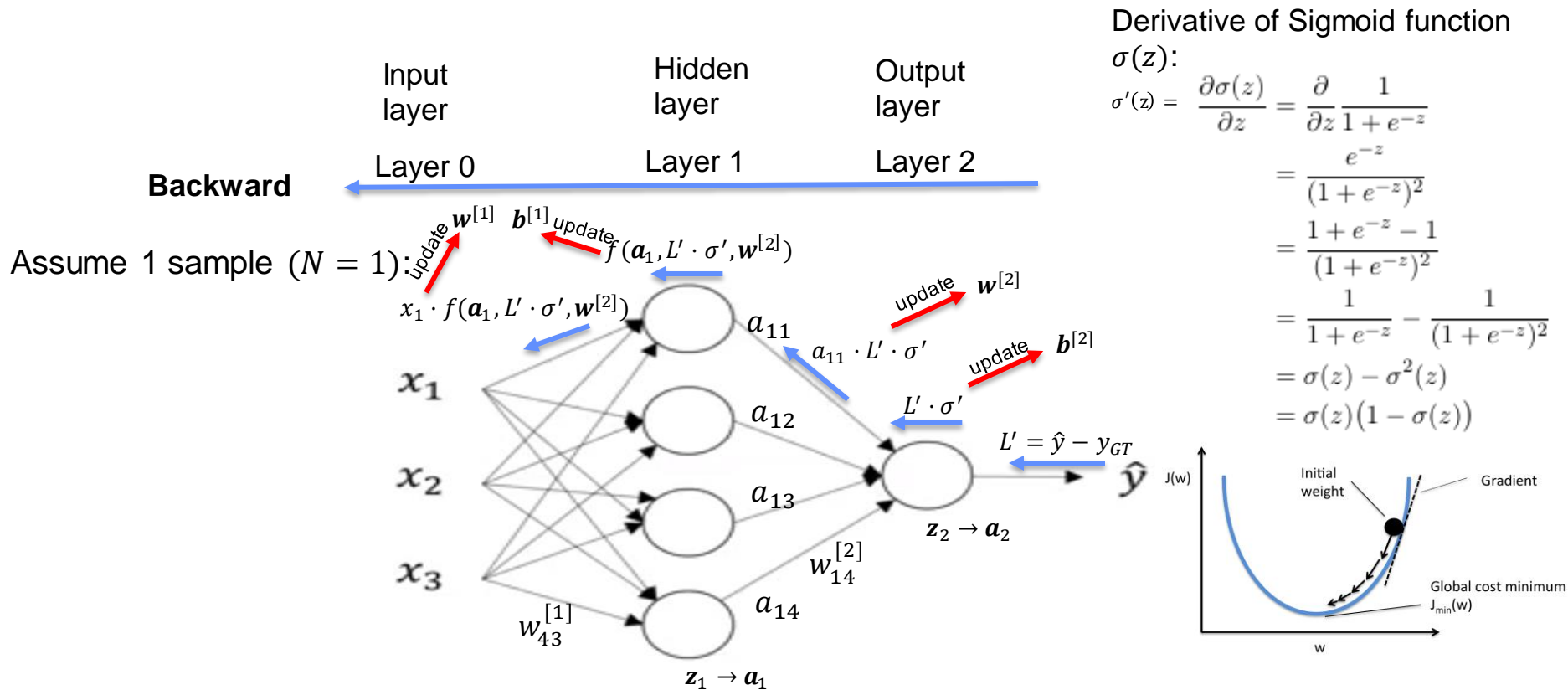
MSE Loss Function:

$$L = (\hat{y} - y)^2$$

$$[x_1 \quad x_2 \quad x_3] \cdot \begin{bmatrix} w_{11}^{[1]} & w_{21}^{[1]} & w_{31}^{[1]} & w_{41}^{[1]} \\ w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} & w_{42}^{[1]} \\ w_{13}^{[1]} & w_{23}^{[1]} & w_{33}^{[1]} & w_{43}^{[1]} \end{bmatrix} + [b_1^{[1]} \quad b_2^{[1]} \quad b_3^{[1]} \quad b_4^{[1]}] = \mathbf{z}_1 \rightarrow \mathbf{a}_1$$

Where, $\mathbf{z}_1 = [z_{11} \quad z_{12} \quad z_{13} \quad z_{14}]$
 $\mathbf{a}_1 = [a_{11} \quad a_{12} \quad a_{13} \quad a_{14}]$ $a_{ij} = \sigma(z_{ij})$

Backpropagation: Gradient Descent



Purpose of backpropagation:

Apply *chain rule of derivative* to update parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$