

# Data Analytics

## EEE 4774 & 6777

### Module 4 - Classification

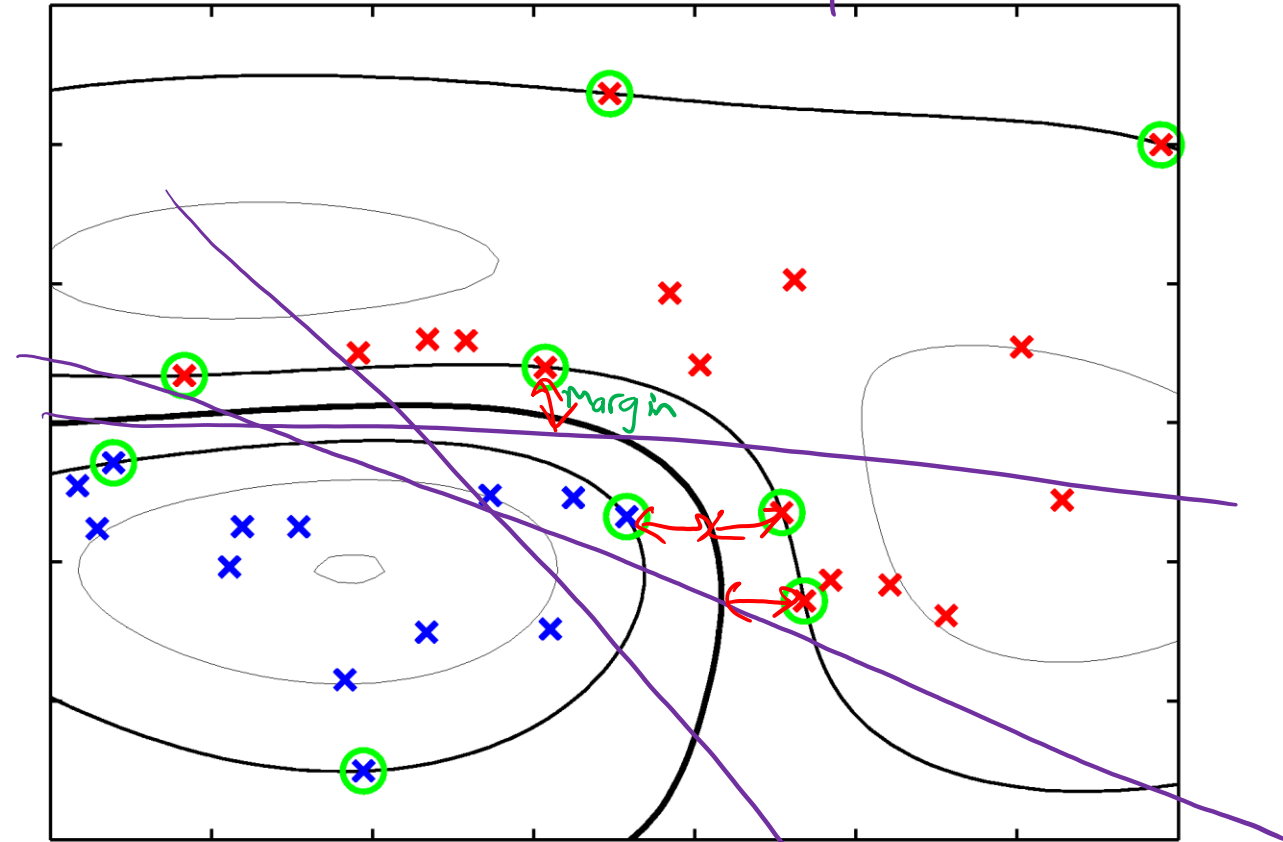
#### Support Vector Machine (SVM)

Spring 2022

# Support Vector Machine (SVM)

*Nonlinear Classification*

- Aims at maximizing the margin between decision boundary and data points
- Used for both classification and regression
- Determination of model parameters corresponds to a convex optimization problem, so any local solution is also a global optimum
- SVM makes extensive use of the Lagrange Multipliers concept from the Optimization Theory
- SVM is a decision machine, so does not provide posterior probabilities (unfortunately).
- Relevance Vector Machine (RVM) is based on Bayesian formulation, and provides posterior probabilities.



# Support Vector Machine (SVM)

$\phi$ : Nonlinear func. of  $x$

- $y(x) = \mathbf{w}^T \phi(x) + b$  : score for  $x$

- Nonlinear fixed feature space mapping  $\phi(x)$

- 2-class model,

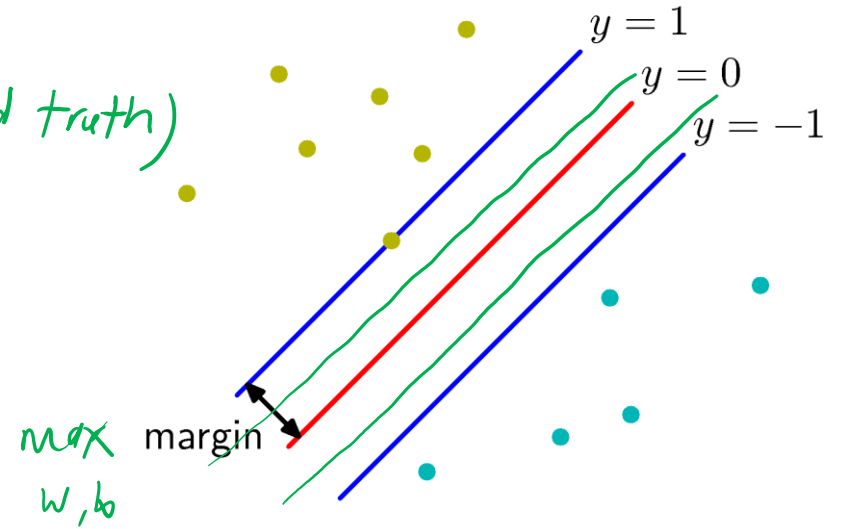
$t_n \in \{-1, +1\}$  actual class label (Ground truth)

class label estimate (SVM output)

- $\hat{t}_n = \begin{cases} +1, & y(x) \geq 0 \\ -1, & y(x) < 0 \end{cases}$  Class 1  
Class 2

- If linearly separable, many solutions exist

- SVM: Maximum margin classifier



Correct Class:  $t_n = \hat{t}_n = +1$  }  $t_n \hat{t}_n = 1$   
 $t_n = \hat{t}_n = -1$

# Support Vector Machine (SVM)

Objective:

Find  $w, b$  which give the maximum margin

- Correctly classified points:  $t_n y(x_n) > 0$

$$\max \frac{t_n y(x_n)}{\|w\|} = \max \frac{t_n (w^T \phi(x_n) + b)}{\|w\|}$$

$w^T w = \|w\|^2$

$$\arg \max_{w, b} \left\{ \frac{1}{\|w\|} \min_n [t_n (w^T \phi(x_n) + b)] \right\}$$

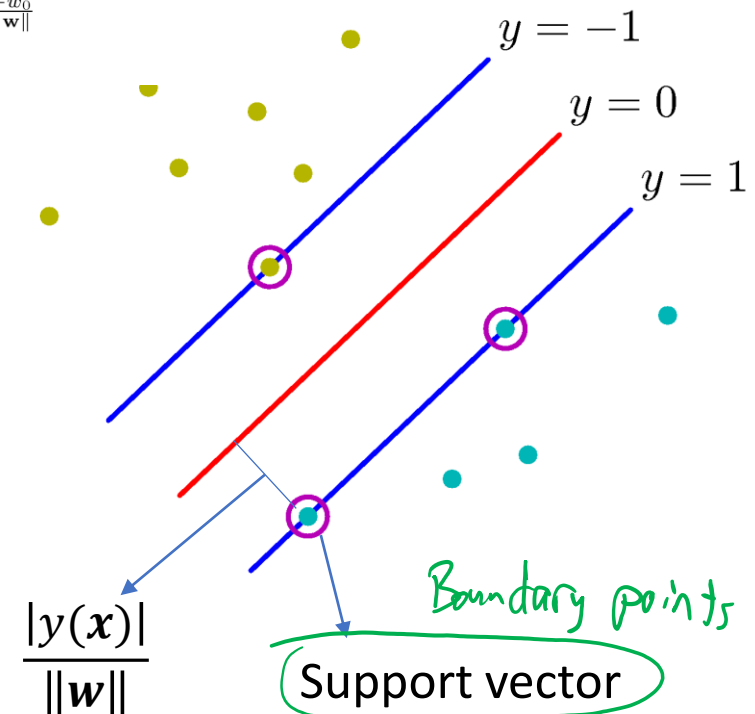
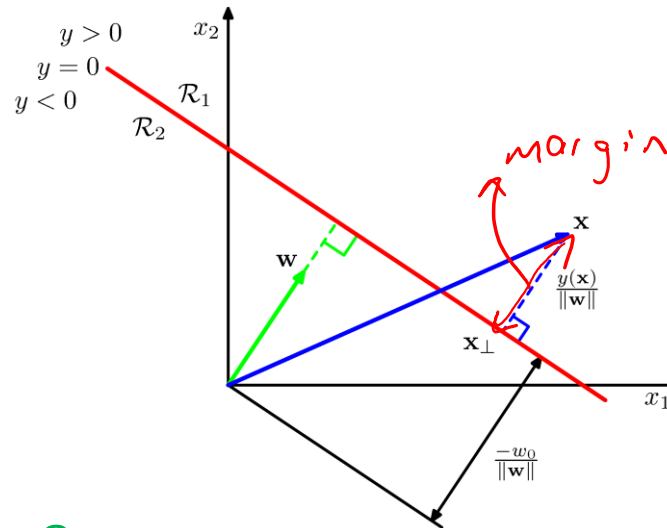
margin  $= \frac{t_n y(x_n)}{\|w\|}$  does not change when  $w \rightarrow \kappa w$  and  $b \rightarrow \kappa b$

$$\frac{t_n (w^T \phi(x_n) + b)}{\|w\|} = \frac{t_n y(x_n)}{\|w\|}$$

- Choose  $\kappa$  such that  $t_n (w^T \phi(x_n) + b) = 1$

$$t_n y(x_n) \geq 1 \quad \forall n$$

- $\arg \min_{w, b} \|w\|^2$  s.t.  $t_n (w^T \phi(x_n) + b) \geq 1, n = 1, \dots, N$



# Support Vector Machine (SVM)

kernel func. : similarity score  
 $k(x, x_n)$  betw.  $x$  and  $x_n$

Lagrange multiplier

- Using Lagrange multipliers  $a_n \geq 0$  we obtain  $w = \sum_{n=1}^N a_n t_n \phi(x_n)$

$$\phi(x_n)^T \phi(x)$$

- Hence,  $y(x) = \sum_{n=1}^N a_n t_n k(x, x_n) + b$   
 $\hookrightarrow$  label for  $n^{th}$  training instance

$$y(x_n) = w^T \phi(x_n) + b$$

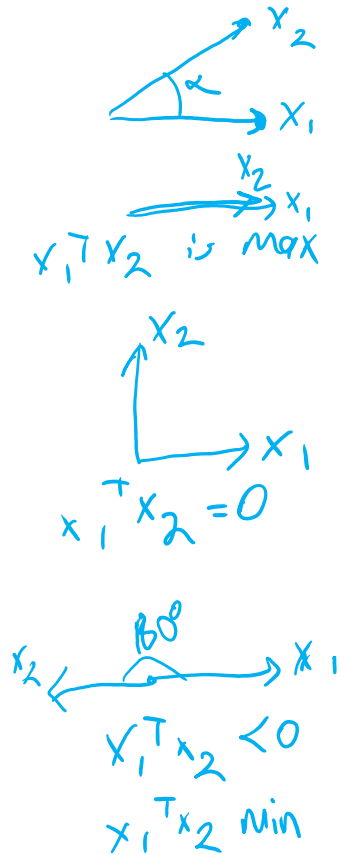
$$y(x_n) = \sum_{n=1}^N a_n t_n \underbrace{\phi(x_n)^T \phi(x_n)} + b$$

- Kernel function:  $k(x_1, x_2) = \phi(x_1)^T \phi(x_2)$

- Stationary kernels:  $k(x_1, x_2) = k(x_1 - x_2)$

- Homogeneous kernels (Radial basis functions):  $k(x_1, x_2) = k(\|x_1 - x_2\|)$

Rbf



# Support Vector Machine (SVM)

**Kernel Trick:** Compute the **similarity score**  $k(\mathbf{x}, \mathbf{x}_n)$  directly without defining  $\phi(\mathbf{x})$

- Linear Kernel:

$$k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2 = x_{11}x_{21} + x_{12}x_{22}$$

- Polynomial kernel

$$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{x}_2 + r)^d$$

e.g., (r=0, d=2) *2<sup>nd</sup> degree polynomial*

$$k(\mathbf{x}_1, \mathbf{x}_2) = x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11}x_{12}x_{21}x_{22}$$

(Hidden)  $\phi(\mathbf{x}_1) = [x_{11}^2, x_{12}^2, \sqrt{2}x_{11}x_{12}]^T$

$$k(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2)$$

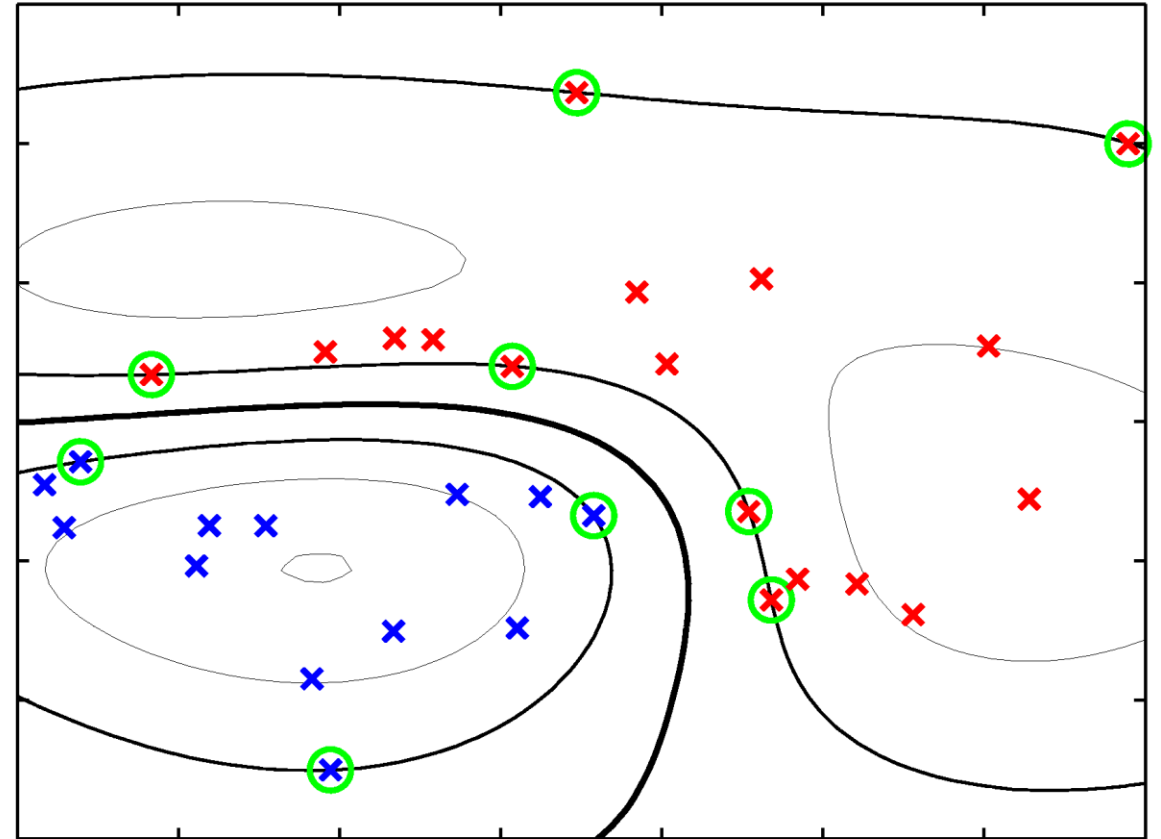
- Rbf (*Gaussian kernel*)

$$k(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2} = e^{-\gamma(\mathbf{x}_1 - \mathbf{x}_2)^T(\mathbf{x}_1 - \mathbf{x}_2)}$$

- Sigmoid

$$k(\mathbf{x}_1, \mathbf{x}_2) = \tanh(\mathbf{x}_1^T \mathbf{x}_2 + r)$$

*hyperbolic tangent*



**Custom Kernels:** You can define your own kernel function. Must be a valid kernel function!

# Support Vector Machine (SVM)

- When written in terms of minimization of a regularized error function, SVM has similarities with Logistic Regression and Perceptron:

in the figure,

**blue** for SVM (also for Perceptron by a shift of 1)

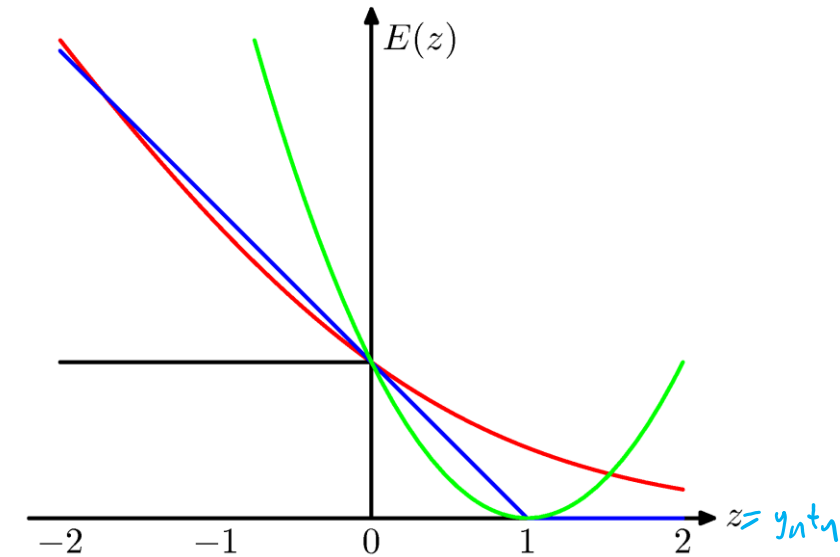
**red** for Logistic Regression

**black** for Misclassification error

**green** for Quadratic error

Hinge Loss:  $E(y_n t_n) = \max\{1 - y_n t_n, 0\}$

*< 0 if  $y_n t_n > 1 \Rightarrow$  correctly classified training pts.  $\Rightarrow E(y_n t_n) = 0$*   
*> 0 if  $y_n t_n < 1 \Rightarrow E(y_n t_n) > 0$*



- Weighted voting (compare to kNN) with weights coming from the similarity metric  $k(\mathbf{x}, \mathbf{x}_n)$   
*kernel func.*

- Extends to multiclass problems

- Used also for **anomaly detection (One-Class SVM)** and **regression (SVR)**

*Support Vector Regression*