

# BITS F464: Machine Learning

# 25

## Support Vector Machine (SVM)



Dr. Kamlesh Tiwari  
Assistant Professor, Department of CSE,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA



March 22, 2021    ONLINE    (Campus @ BITS-Pilani Jan-May 2021)

<http://ktiware.in/ml>

### Geometry

- Essentially; distance of a point  $X = (x_1, x_2, \dots, x_n)$  from a hyperplane represented by  $(b, w_1, w_2, \dots, w_n)$  is given by

$$\frac{W^T X + b}{\|W\|}$$

where  $\|W\|$  is norm <sup>1</sup>

#### Classification

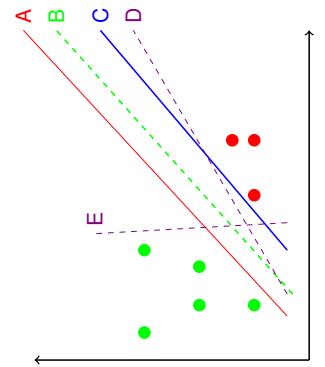
- Hyperplane have two sides (say +ve and -ve)
  - which side a point  $X = (x_1, x_2, \dots, x_n)$  lies?
- Substitute coordinates in the equation  $W^T X + b$  and check the sign

<sup>1</sup> $W = (w_1, w_2, \dots, w_n)$ , and  $\|W\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$

Machine Learning (BITS F464)    MWF (10-11 AM) online@BITS-Pilani    Lecture-25(March 22, 2021)    3/12

### Which decision boundary is better

Many decision boundaries with perfect classification are possible



Note: The separating hyperplane would be in middle

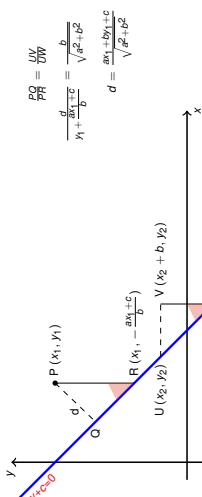
Machine Learning (BITS F464)    MWF (10-11 AM) online@BITS-Pilani    Lecture-25(March 22, 2021)    5/12

### Geometry

- Determine the length of the perpendicular drawn on a line  
 $ax + by + c = 0$  from a point  $(x_1, y_1)$

$$d = \frac{|ax_1 + by_1 + c|}{\sqrt{a^2 + b^2}}$$

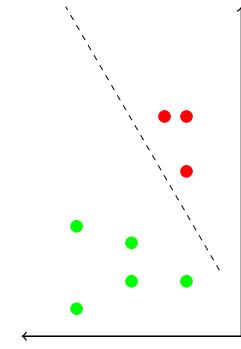
- Proof: (by geometry)



Machine Learning (BITS F464)    MWF (10-11 AM) online@BITS-Pilani    Lecture-25(March 22, 2021)    2/12

### Special Case of Classification

Consider a linearly separable dataset



- Hyperplane is defined by  $W$
- $Margin$  is the distance of nearest data point from the separating hyperplane
- Most of the real-world problems are NOT linearly separable
- Sometime data could be transformed to a high-dimensional space where classes may be linearly separable
- Caution:** this could lead to over-fitting

Machine Learning (BITS F464)    MWF (10-11 AM) online@BITS-Pilani    Lecture-25(March 22, 2021)    4/12

### Look Closer

- $\arg\max_{W, b} \left[ \min_{x_i \in D} \frac{|W^T x_i + b|}{\|W\|} \right] = \arg\max_{W, b} \frac{1}{\|W\|} \left[ \min_{x_i \in D} |W^T x_i + b| \right] = \arg\max_{W, b} \frac{1}{\|W\|} \left[ \min_{x_i \in D} |W^T x_i + b| \right]$
- Hypothesis is a hyperplane represented by  $W$ ; scaled parameter  $k.W$  also represents the same hyperplane
  - For the points on hyperplane  $W^T x_i + b = 0$
  - For points NOT on hyperplane  $W^T x_i + b = 0$  changes if  $k.W$  used instead of  $W$ . Ultimately different  $\min_{x_i \in D} |W^T x_i + b|$
  - You can get any value for  $\min_{x_i \in D} |W^T x_i + b|$  by changing the value of  $k$  without changing the hyperplane

So without loss of generality, let us fix  $\min_{x_i \in D} |W^T x_i + b| = 1$

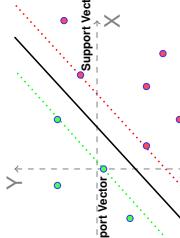
Machine Learning (BITS F464)    MWF (10-11 AM) online@BITS-Pilani    Lecture-25(March 22, 2021)    6/12

## Support Vector Machine (SVM)

SVM is a liner decision machine; uses  $\text{sign}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$  for decision

- We want  $\mathbf{w}^T \mathbf{x}^{(i)} + b \geq \gamma$  for **+ve** (and  $< -\gamma$  for **-ve**)
- Distance of a point  $(x, y)$  from hyper plane  $\mathbf{w}^T \mathbf{x} + b = 0$  is  $\frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|}$
- Distance can be maximized, by either **maximizing  $b$**  or by **minimizing  $\|\mathbf{w}\|$**

- We needs  $\mathbf{w}^T \mathbf{x} + b \geq \gamma \|\mathbf{w}\|$   
let  $\gamma \|\mathbf{w}\| = 1$ 
  - $\mathbf{w}^T \mathbf{x} + b \geq 1$  if  $\mathbf{x}$  is **+1**
  - $\mathbf{w}^T \mathbf{x} + b \leq -1$  if  $\mathbf{x}$  is **-1**
- It leads to  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1$
- Points with  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1$  are called **support vector**



## Lagrangian with optimized values

$$\begin{aligned} L(\mathbf{w}, b) &= \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_i \alpha_i (y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1) \\ L(\mathbf{w}, b) &= \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_i \alpha_i y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} - \sum_i \alpha_i y^{(i)} b + \sum_i \alpha_i \\ L(\mathbf{w}, b) &= \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_i \alpha_i y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} + \sum_i \alpha_i \\ L(\mathbf{w}, b) &= \frac{1}{2} \sum_i \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} - \sum_i \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} + \sum_i \alpha_i \\ L(\mathbf{w}, b) &= \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} \\ \text{One have to maximize } L(\mathbf{w}, b) \text{ subject to } \alpha_i \geq 0 \text{ and} \\ \sum_{i=1}^m \alpha_i y^{(i)} &= 0 \end{aligned}$$

- If  $\alpha_i$  is large the corresponding training sample is support vector.
- Otherwise, when  $\alpha_i = 0$  it is not a support vector
- Optimized  $\alpha_i$  provides the value of  $\mathbf{w} = \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)}$  to be used in linear decision boundary

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 9/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 10/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 11/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 12/12

## Support Vector Machine (SVM)

### Advantages

- It works well with clear margin of separation
- Do not stuck to local minima
- Effective in high dimensional spaces
- It is effective in cases where number of dimensions is greater than the number of samples
- Memory efficient as it uses a subset of training points in decision (called support vectors)

### Disadvantages

- It doesn't perform well, when data set is large as the training time is higher
- Doesn't perform well when target classes are overlapping
- Probability estimates are not directly provide
- Higher classification confidence for points lying far from the decision boundary

## Support Vector Machine (SVM)

- Minimization of  $\mathbf{w}$  is same as minimization of  $\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w}$
- Other constraints are  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1$
- However, for support vectors  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1$
- Define a Lagrangian Multiplier to optimize
- $L(\mathbf{w}, b) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_i \alpha_i (y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1)$
- Derivative  $\frac{\partial L}{\partial b} = -\sum_i \alpha_i y^{(i)}$  that should be equated to zero

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

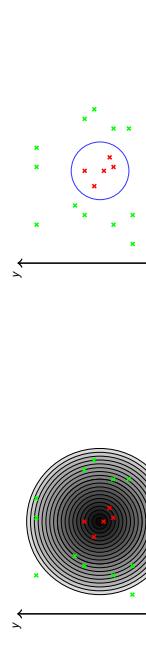
- Derivative  $\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y^{(i)} \mathbf{x}^{(i)}$  equated to zero gives

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

## Transformation: An Example



- Transform all 2D points  $(\mathbf{x}^{(i)}, y^{(i)})$  in 3D as  $(\mathbf{x}^{(i)}, y^{(i)}, z)$  where  $z = (\mathbf{x}^{(i)} - \mathbf{x}_0)^2 + (y^{(i)} - y_0)^2$



## Thank You!

(Reference<sup>2</sup>)

- Thank you very much for your attention!

## Queries ?

<sup>2</sup> [1] Text Book: Machine Learning by Tom M Mitchell [2] Mod-01 Lec-29 Support Vector Machine  
<https://www.youtube.com/watch?v=vSwvHtSQTE>

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 12/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 13/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 14/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 15/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 16/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 17/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 18/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 19/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 20/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 21/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 22/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 23/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 24/12

Machine Learning (BITS F464) MWF (10-11 AM) online@BITS-Pilani Lecture-25(March 22, 2021) 25/12