



CS-F415: Data Mining

08

Support Vector Machine (SVM)



Dr. Kamlesh Tiwari
Associate Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA
Feb 19, 2024 **MW/F 4:00pm** 6101 @ BITS-Pilani [Jan-May 2024]

<http://ktiwari.in/dm>

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 L_2 norm¹

Classification

- Hyperplane have two sides (say +ve and -ve)
- which side a point $X = (x_1, x_2, \dots, x_n)$ lies?

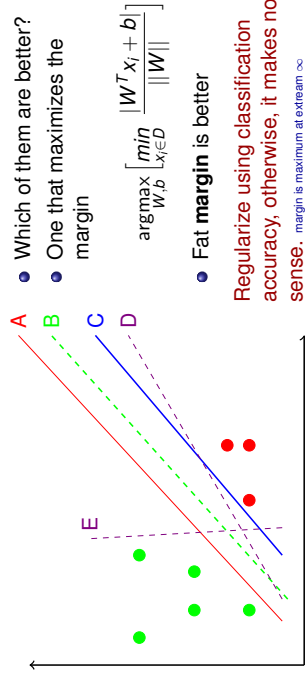
Substitute coordinates in the equation $W^T X + b$ and check the sign

Note: W defines same hyperplane as $2W, 3W, \dots, kW$

¹ $\|W = (w_1, w_2, \dots, w_n)$, and the L_2 norm is $\|W\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$

Which decision boundary is better

Many decision boundaries with perfect classification are possible



- Which of them are better?
- One that maximizes the margin

$$\operatorname{argmax}_{W,b} \left[\min_{X_i \in D} \frac{|W^T X_i + b|}{\|W\|} \right]$$

- Fat margin is better

Regularize using classification accuracy, otherwise, it makes no sense. margin is maximum at extreme ∞

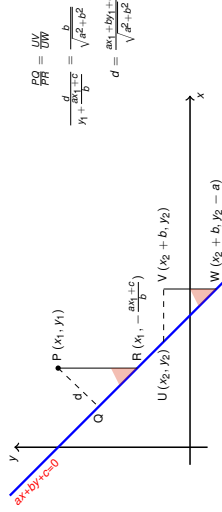
Note: The separating hyperplane would be in middle

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)



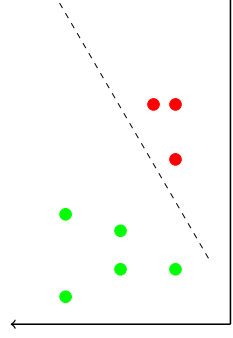
$$\frac{PO}{\sin \theta} = \frac{PW}{1}$$

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

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

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

$$\min_{X_i \in D} \frac{|W^T X_i + b|}{\|W\|}$$

- 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

Let's Look Closer

$$\operatorname{argmax}_{W,b} \left[\min_{X_i \in D} \frac{|W^T X_i + b|}{\|W\|} \right] = \operatorname{argmax}_{W,b} \frac{1}{\|W\|} \left[\min_{X_i \in D} |W^T X_i + b| \right] = \operatorname{argmax}_{W,b} \frac{1}{\|W\|}$$

- Hypothesis is a hyperplane represented by W ; scaled parameter $k \cdot W$ also represents the same hyperplane
- For the points on hyperplane $W^T X_i + b = 0$
- For points NOT on hyperplane $W^T X_i + b$ changes if $k \cdot W$ is used instead of W . Leading to different $\min_{X_i \in D} |W^T X_i + b|$
- One 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 $\min_{X_i \in D} |W^T X_i + b| = 1$

Support Vector Machine (SVM)

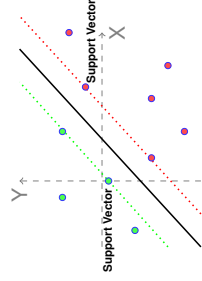
SVM is a linear decision machine; uses $\text{sign}(w^T x^{(i)} + b)$ for decision

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

- We need $w^T x + b \geq \gamma \|w\|$
let $\gamma \|w\| = 1$

- $w^T x + b \geq 1$ if x is +1
- $w^T x + b \leq -1$ if x is -1

- It leads to $y^{(i)}(w^T x^{(i)} + b) \geq 1$
- Points with $y^{(i)}(w^T x^{(i)} + b) = 1$ are called **support vector**



Support Vector Machine (SVM)

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

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

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

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

Lagrangian with optimized values

$$L(w, b) = \frac{1}{2} w \cdot w - \sum \alpha_i (y^{(i)}(w x^{(i)} + b) - 1)$$

$$L(w, b) = \frac{1}{2} w \cdot w - \sum \alpha_i y^{(i)} w x^{(i)} - \sum \alpha_i y^{(i)} b + \sum \alpha_i$$

$$L(w, b) = \frac{1}{2} w \cdot w - \sum \alpha_i y^{(i)} w x^{(i)} + \sum \alpha_i$$

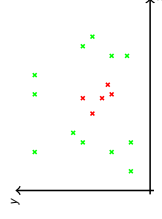
$$L(w, b) = \frac{1}{2} \sum \alpha_i y^{(i)} y^{(j)} x^{(i)} x^{(j)} - \sum \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)} x^{(j)} + \sum \alpha_i$$

$$L(w, b) = \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)} x^{(j)}$$

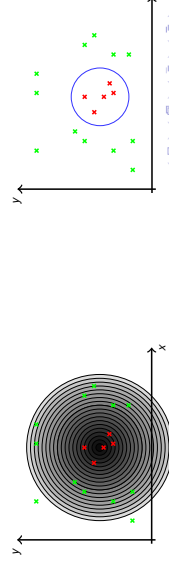
- One have to maximize $L(w, b)$ subject to the constraints $\alpha_i \geq 0$ and $\sum_{i=1}^m \alpha_i y^{(i)} = 0$
- A training sample is support vector, if corresponding α_i is large. Otherwise, when α_i is near to 0 it is not a support vector
- Optimized α_i provides the value of $w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$ to be used in linear decision boundary
- This w can further be used in equation $y^{(i)}(w^T x^{(i)} + b) = 1$ to get the value of b using a data point

Thank You!

Transformation: An Example



- Transform all 2D points $(x^{(i)}, y^{(i)})$ in 3D as $(x^{(i)}, y^{(i)}, z)$ where $z = (x^{(i)} - x_c)^2 + (y^{(i)} - y_c)^2$



Thank you very much for your attention!

Queries ?