



BITS Pilani
Pilani Campus

Support Vector Machines

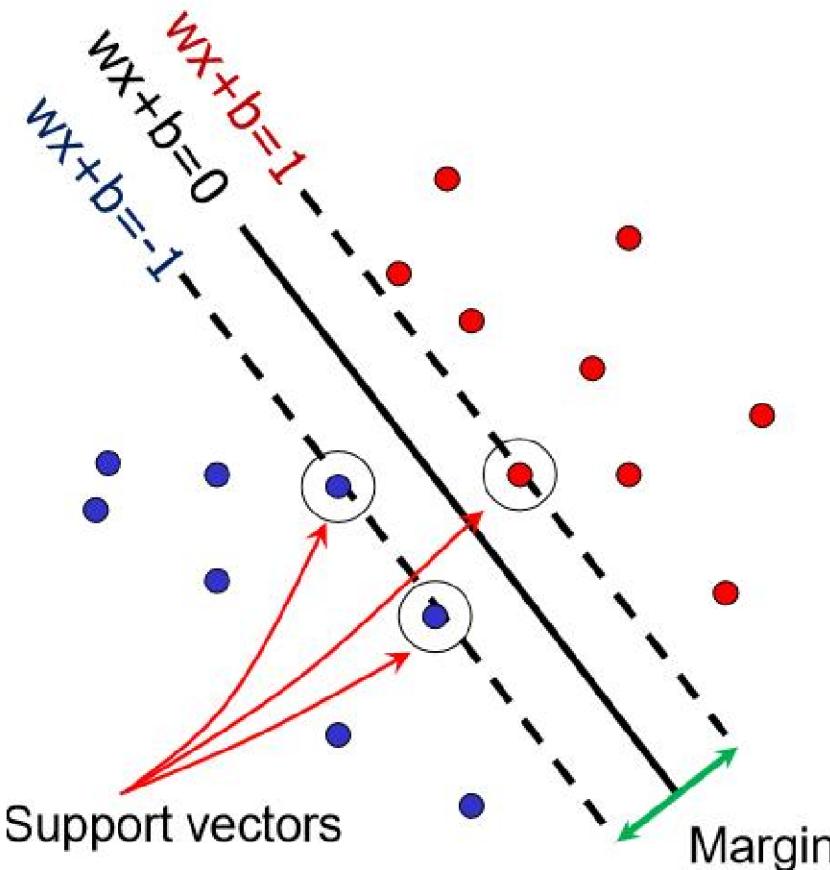
MFML Team

Topics to be covered

- Linear Classifiers
 - Maximum Margin Classification
 - Linear SVM
 - SVM optimization problem
 - Soft Margin SVM
-

Support Vector Machines

- Want line that maximizes the margin.



\mathbf{x}_i positive ($y_i = 1$): $\mathbf{x}_i \cdot \mathbf{w} + b \geq 1$

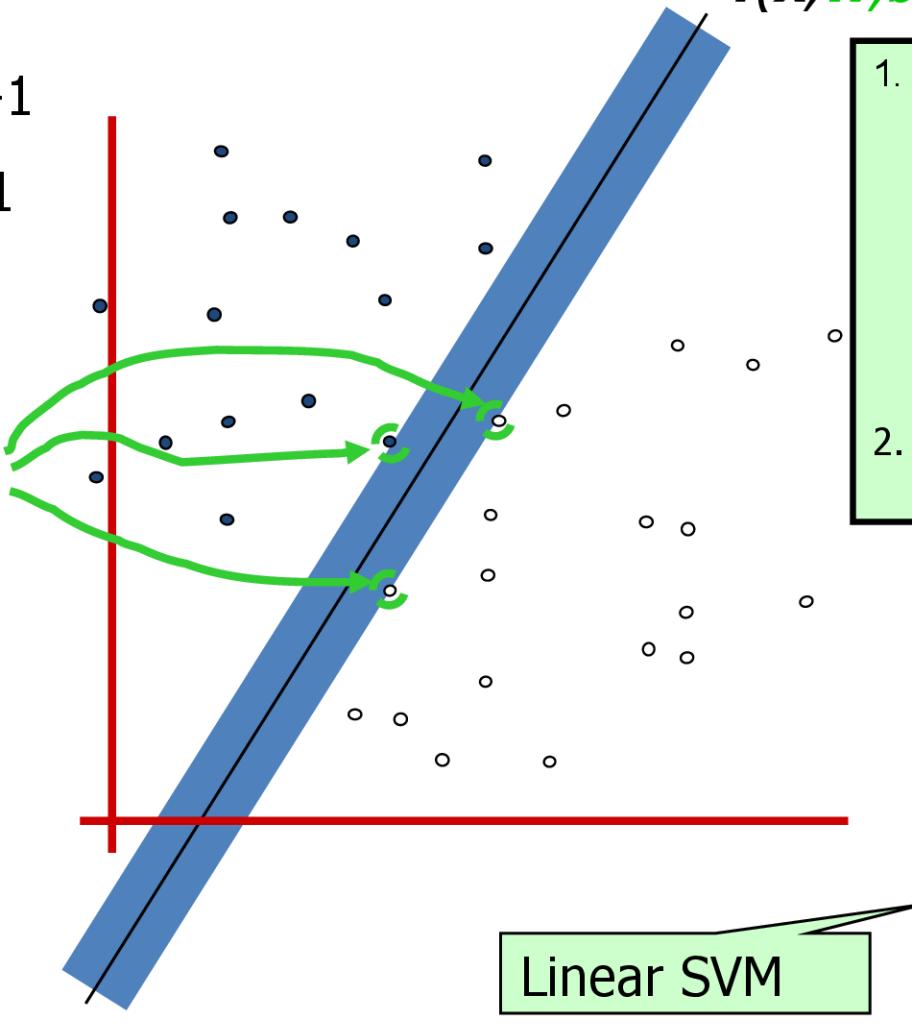
\mathbf{x}_i negative ($y_i = -1$): $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$

For support vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Maximum Margin

denotes +1
 denotes -1

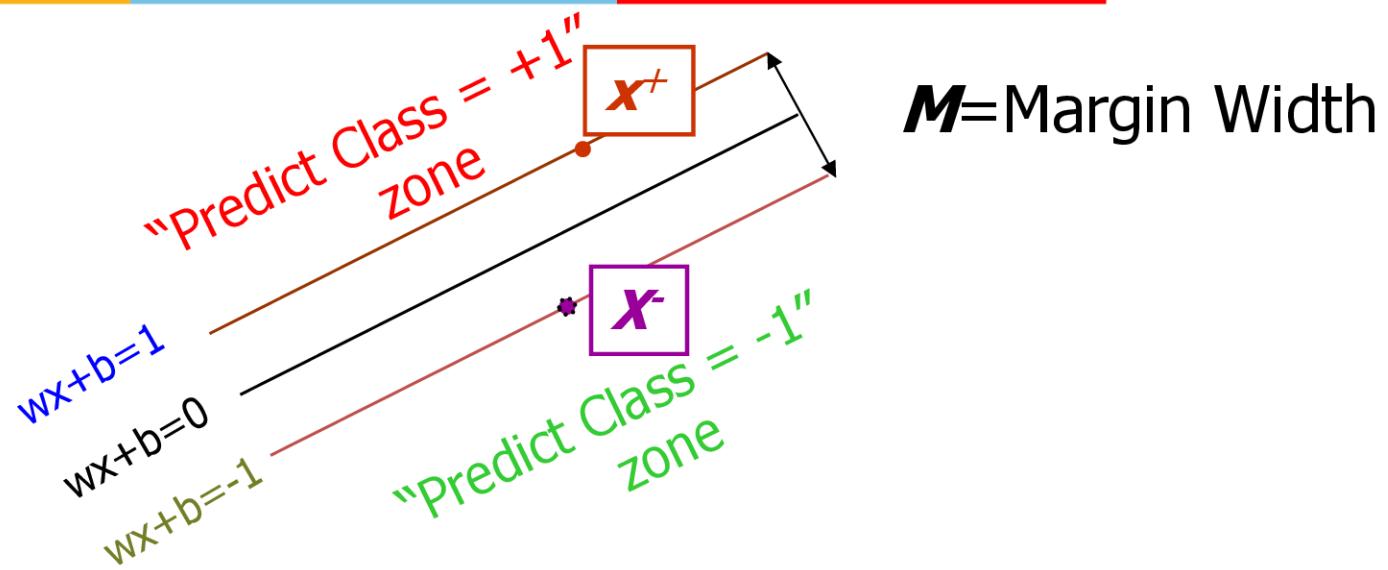
Support Vectors
 are those
 datapoints that
 the margin
 pushes up
 against



Support Vectors

- Geometric description of SVM is that the max-margin hyperplane is completely determined by those points that lie nearest to it.
- Points that lie on this margin are the support vectors.
- The points of our data set which if removed, would alter the position of the dividing hyperplane

Linear SVM Mathematically



Distance between lines given by solving linear equation:

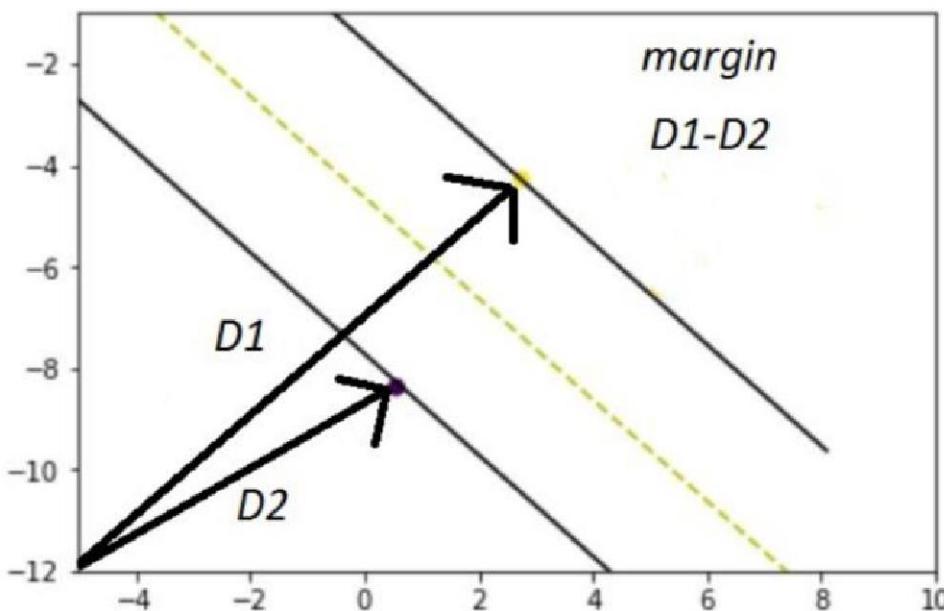
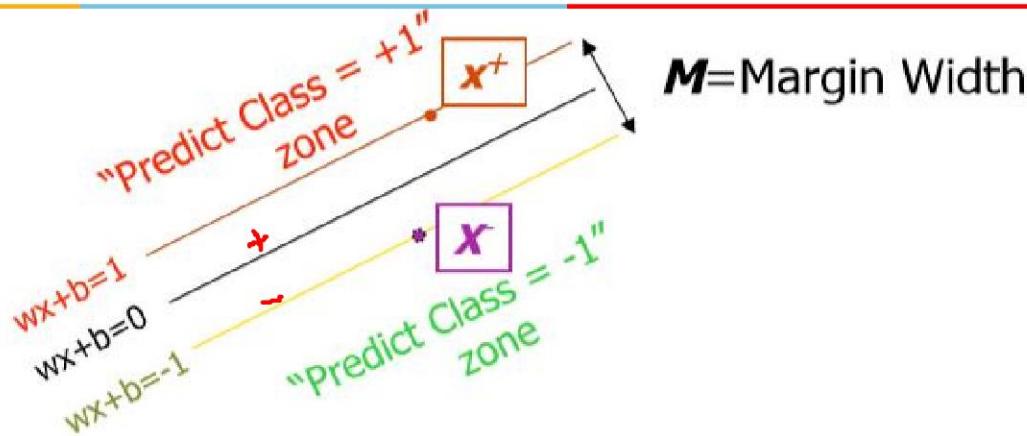
What we know:

- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$

$$\text{Maximize margin: } M = \frac{2}{\|\mathbf{w}\|}$$

$$\text{Equivalent to minimize: } \frac{1}{2} \|\mathbf{w}\|^2$$

Linear SVM Mathematically



$$D1 = w^T x + b = 1 \quad w^T x + b - 1 = 0$$

$$D2 = w^T x + b = -1 \quad w^T x + b + 1 = 0$$

$$w^T x + b - 1 - w^T x + b + 1$$



Solve algebraically

$$\frac{2}{|w|}$$

Optimization Problem

1. Maximize margin $2/\|\mathbf{w}\|$
2. Correctly classify all training data points:

$$\mathbf{x}_i \text{ positive } (y_i = 1) : \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1) : \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

Quadratic optimization problem:

Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \frac{1}{2}\|\mathbf{w}\|^2$ is minimized;

and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

$$+1(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

$$-1(\mathbf{w}^T \mathbf{x}_i + b) \leq 1$$

same as $(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Solving the Optimization Problem

Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$ is minimized; Type equation here.

and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

← Primal

- Need to optimize a *quadratic function subject to linear inequality* constraints.
- All constraints in SVM are linear
- Quadratic optimization problems are a well-known class of mathematical programming problems, and many (rather intricate) algorithms exist for solving them.
- The solution involves constructing a *unconstrained problem* where a *Lagrange multiplier α_i* is associated with every constraint in the primary problem:

Recall : Karush–Kuhn–Tucker (KKT) theorem

- $\max f(x)$ subject to $g_i(x) = 0$ and $h_j(x) \geq 0$ for all i, j
- Make the Lagrangian function

$$\mathcal{L} = f(x) - \sum_i \lambda_i g_i(x) - \sum_j \mu_j h_j(x)$$

- Necessary conditions to have a minimum are

$$\nabla_x \mathcal{L}(x^*, \lambda^*, \mu^*) = 0$$

$$g_i(x^*) = 0 \text{ for all } i$$

$$h_j(x^*) \geq 0 \text{ for all } j$$

$$\mu_j \geq 0 \text{ for all } j$$

$$\mu_j^* h_j(x^*) = 0 \text{ for all } j$$

Solving the Optimization Problem

- The solution involves constructing a *dual problem* where a *Lagrange multiplier* α_i is associated with every constraint in the primary problem:

$$L(w, b, \alpha_i) = \frac{1}{2}||w||^2 - \sum \alpha_i [y_i (w^T x_i + b) - 1]$$

- Taking partial derivative with respect to w , $\frac{\partial L}{\partial w} = 0$
 - $w - \sum \alpha_i y_i x_i = 0$
 - $w = \sum \alpha_i y_i x_i$
- Taking partial derivative with respect to b , $\frac{\partial L}{\partial b} = 0$
 - $-\sum \alpha_i y_i = 0$
 - $\sum \alpha_i y_i = 0$

Solving the Optimization Problem

$$L(w, b, \alpha_i) = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i (w \cdot x_i + b) - 1]$$

- Expanding above equation:

$$L(w, b, \alpha_i) = \frac{1}{2} \|w\|^2 - \sum \alpha_i y_i w \cdot x_i - \sum \alpha_i y_i b + \sum \alpha_i$$

- Substituting $w = \sum \alpha_i y_i x_i$ and $\sum \alpha_i y_i = 0$ in above equation

$$L(w, b, \alpha_i) = \frac{1}{2} (\sum_i \alpha_i y_i x_i)(\sum_j \alpha_j y_j x_j) - (\sum_i \alpha_i y_i x_i)(\sum_j \alpha_j y_j x_j) + \sum \alpha_i$$

$$L(w, b, \alpha_i) = \sum \alpha_i - \frac{1}{2} (\sum_i \alpha_i y_i x_i)(\sum_j \alpha_j y_j x_j)$$

$$L(w, b, \alpha_i) = \sum \alpha_i - \frac{1}{2} (\sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i \cdot x_j)$$

Support Vectors

Using KKT conditions : $\alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1] = 0$

For this condition to be satisfied either $\alpha_i = 0$ and $y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 > 0$

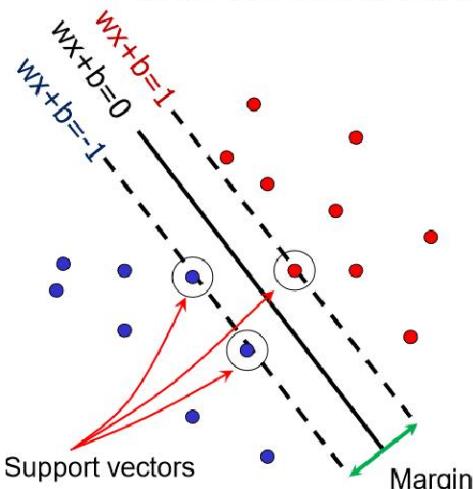
OR $y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 = 0$ and $\alpha_i > 0$

For support vectors: $y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 = 0$

$$L(\mathbf{w}, b, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1]$$

For all points other than support vectors: $\alpha_i = 0$

- Want line that maximizes the margin.



$$\mathbf{x}_i \text{ positive } (y_i = 1) : \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1) : \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

$$\text{For support vectors, } \mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$$

Solving the Optimization Problem

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$

Learned
weight

Support
vector

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

Learned
weight

Support
vector

Solving the Optimization Problem

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$ (for any support vector)

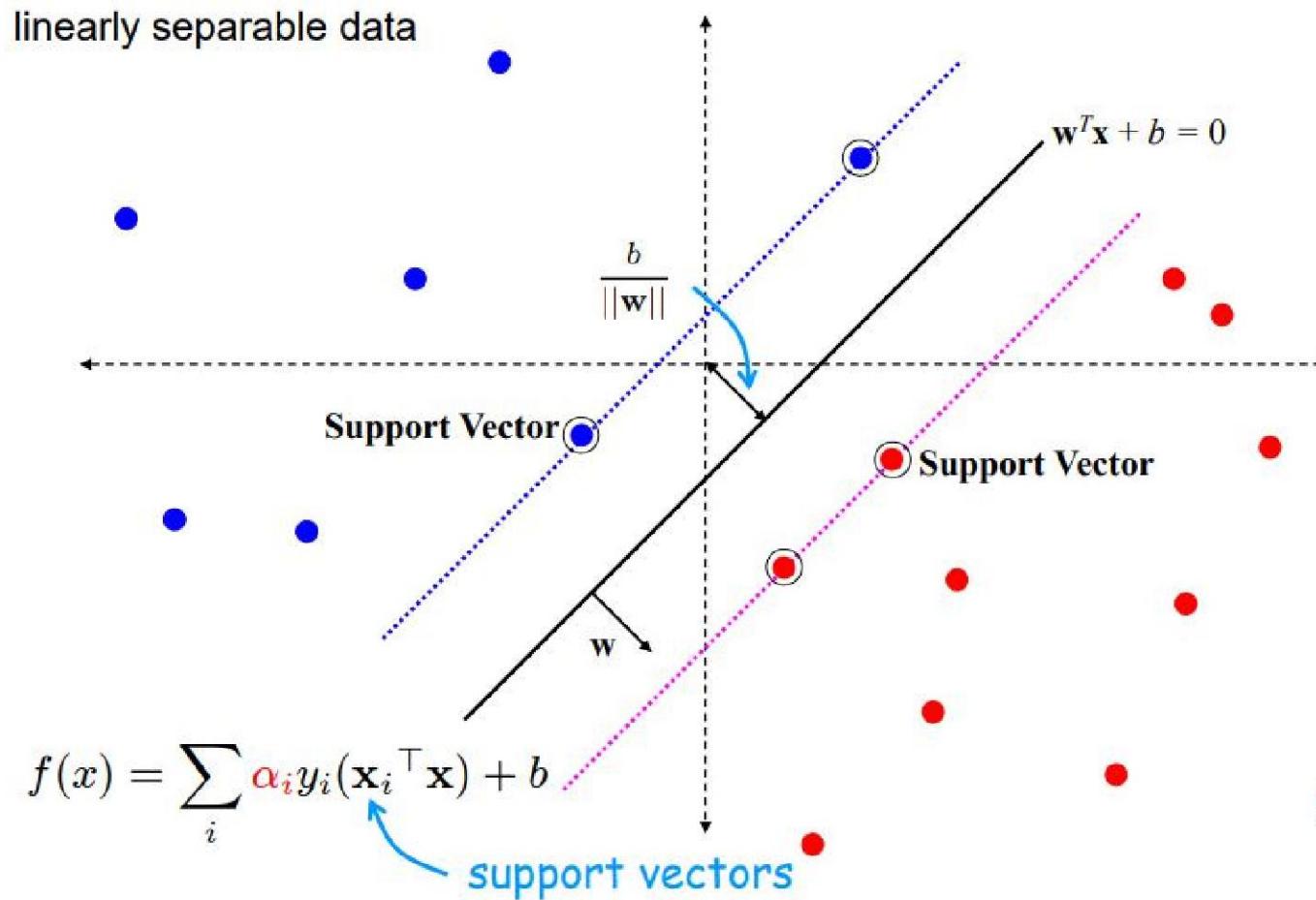
- Classification function:

$$\begin{aligned}f(x) &= \text{sign}(\mathbf{w} \cdot \mathbf{x} + b) \\&= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)\end{aligned}$$

If $f(x) < 0$, classify as negative, otherwise classify as positive.

- Notice that it relies on an *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i
- (Solving the optimization problem also involves computing the inner products $\mathbf{x}_i \cdot \mathbf{x}_j$ between all pairs of training points)

Substituting w in support vectors function



Linear SVM: Numerical Problem

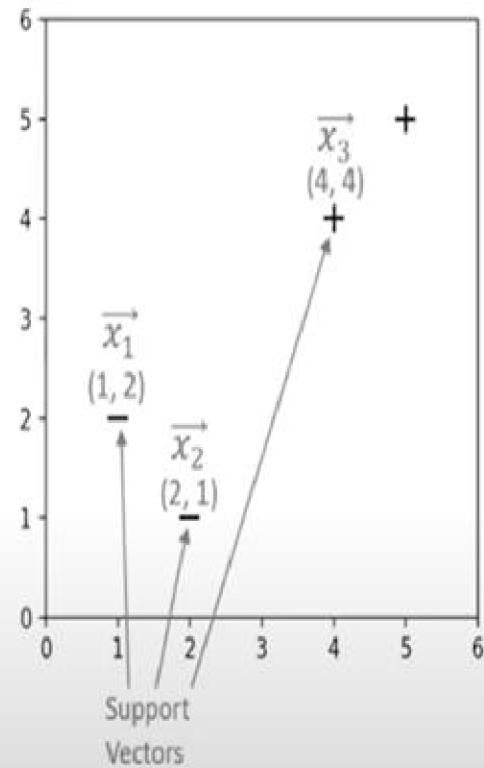
$$\sum_i \alpha_i y_i \vec{x}_i \cdot \vec{\chi} + b = 1 \quad \text{if } \vec{\chi} \text{ is +ve support vector}$$

$$\sum_i \alpha_i y_i \vec{x}_i \cdot \vec{\chi} + b = -1 \quad \text{if } \vec{\chi} \text{ is -ve support vector}$$

$$\sum_i \alpha_i y_i = 0$$

$$\text{Let } \beta_i = \alpha_i y_i$$

$$\begin{cases} \beta_1 \vec{x}_1 \cdot \vec{x}_1 + \beta_2 \vec{x}_2 \cdot \vec{x}_1 + \beta_3 \vec{x}_3 \cdot \vec{x}_1 + b = -1 & \vec{\chi} = \vec{v}_1 \\ \beta_1 \vec{x}_1 \cdot \vec{x}_2 + \beta_2 \vec{x}_2 \cdot \vec{x}_2 + \beta_3 \vec{x}_3 \cdot \vec{x}_2 + b = -1 & \vec{\chi} = \vec{v}_2 \\ \beta_1 \vec{x}_1 \cdot \vec{x}_3 + \beta_2 \vec{x}_2 \cdot \vec{x}_3 + \beta_3 \vec{x}_3 \cdot \vec{x}_3 + b = 1 & \vec{\chi} = \vec{v}_3 \\ \beta_1 + \beta_2 + \beta_3 = 0 \end{cases}$$



Numerical: Linear SVM

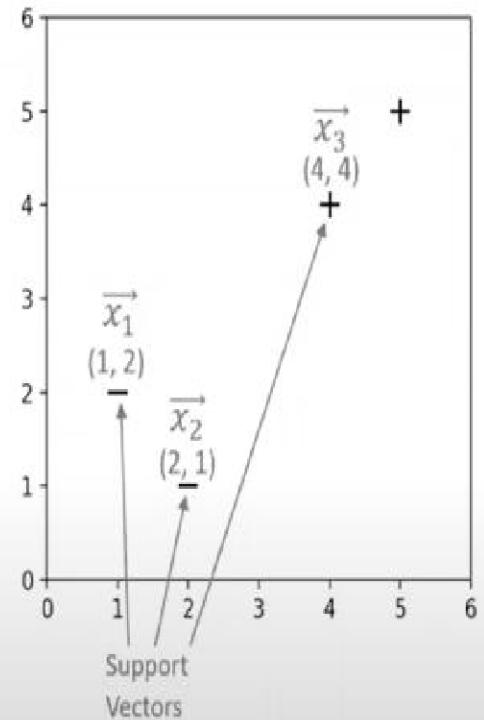
$$\begin{cases} \beta_1 \vec{x}_1 \cdot \vec{x}_1 + \beta_2 \vec{x}_2 \cdot \vec{x}_1 + \beta_3 \vec{x}_3 \cdot \vec{x}_1 + b = -1 \\ \beta_1 \vec{x}_1 \cdot \vec{x}_2 + \beta_2 \vec{x}_2 \cdot \vec{x}_2 + \beta_3 \vec{x}_3 \cdot \vec{x}_2 + b = -1 \\ \beta_1 \vec{x}_1 \cdot \vec{x}_3 + \beta_2 \vec{x}_2 \cdot \vec{x}_3 + \beta_3 \vec{x}_3 \cdot \vec{x}_3 + b = 1 \\ \beta_1 + \beta_2 + \beta_3 = 0 \end{cases}$$

Plug in the values of all the vectors

$$\begin{cases} 5\beta_1 + 4\beta_2 + 12\beta_3 + b = -1 \\ 4\beta_1 + 5\beta_2 + 12\beta_3 + b = -1 \\ 12\beta_1 + 12\beta_2 + 32\beta_3 + b = 1 \\ \beta_1 + \beta_2 + \beta_3 = 0 \end{cases}$$



$$\begin{cases} \beta_1 = -0.08 \\ \beta_2 = -0.08 \\ \beta_3 = 0.16 \\ b = -2.2 \end{cases}$$



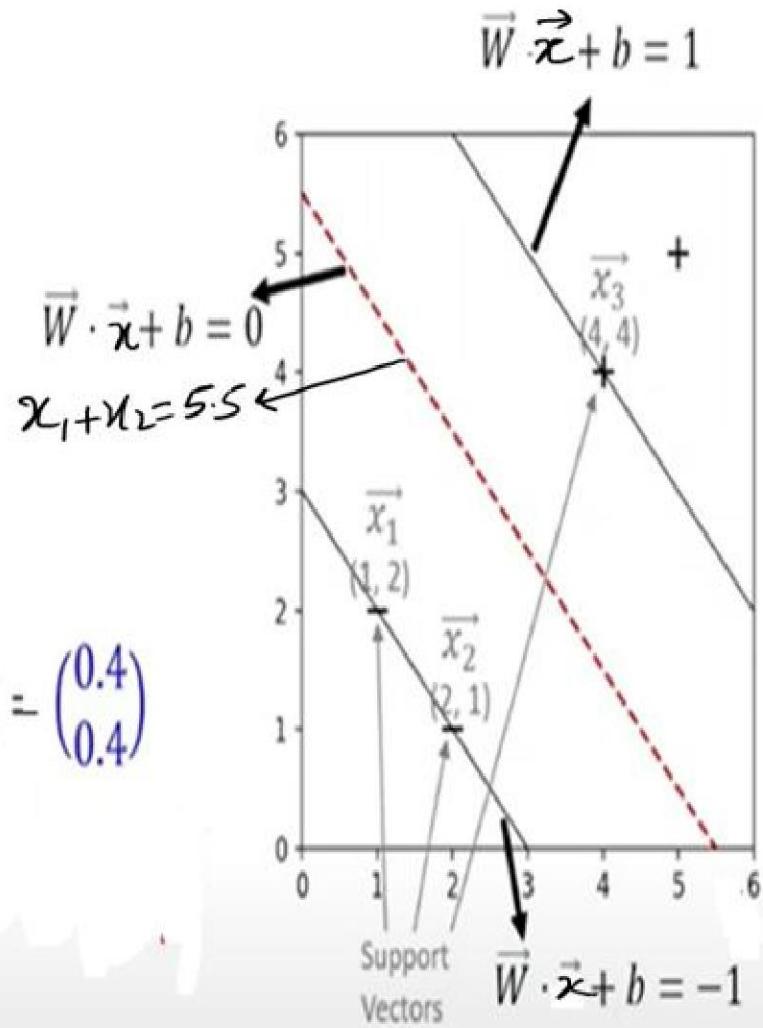
$$\begin{cases} \beta_1 = -0.08 \\ \beta_2 = -0.08 \\ \beta_3 = 0.16 \end{cases}$$

So, hyperplane is essentially those
that satisfy the equality.

$$\vec{W} = \sum_i \alpha_i y_i \vec{x}_i \quad \sum_i \beta_i \vec{x}_i = -0.08 \begin{pmatrix} 1 \\ 2 \end{pmatrix} - 0.08 \begin{pmatrix} 2 \\ 1 \end{pmatrix} + 0.16 \begin{pmatrix} 4 \\ 4 \end{pmatrix} = \begin{pmatrix} 0.4 \\ 0.4 \end{pmatrix}$$

$$: \vec{W} \cdot \vec{x} + b = 0$$

$$\text{So, } \begin{pmatrix} 0.4 \\ 0.4 \end{pmatrix} \cdot \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \end{pmatrix} - 2.2 = 0 \rightarrow \kappa_1 + \kappa_2 = 5.5$$



Hinge Loss

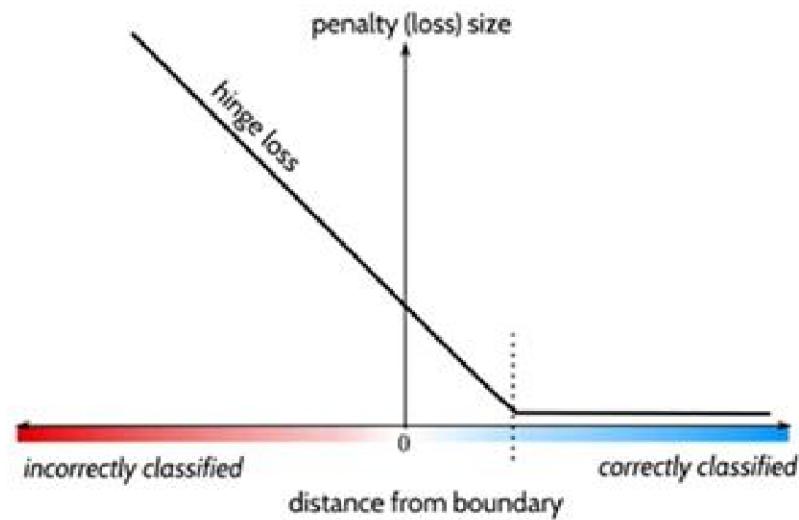
Hinge Loss is one of the types of Loss Function, mainly used for maximum margin classification models.

Hinge Loss incorporates a margin or distance from the classification boundary into the loss calculation. Even if new observations are classified correctly, they can incur a penalty if the margin from the decision boundary is not large enough.

$$L = \max (0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))$$

0 - for correct classification

1 - for wrong classification



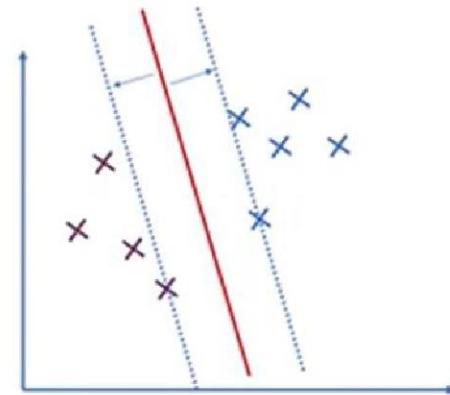
Hinge loss

HINGE LOSS – Numerical Example 1

Correctly classification

- Actual output $y = +1$. Predicted output $y' = 0.5.$ ✓
- Hinge Loss = $\max(0, 1 - y * y')$
 $= \max(0, 1 - 1 * 0.5)$
 $= 0.5.$

$$y' = w^T x + b$$



Hinge loss is **close to zero** for correctly classified sample.

Misclassification

- Actual output $y = -1.$ Predicted output $y' = 0.5.$
- Hinge Loss = $\max(0, 1 - y * y')$
 $= \max(0, 1 - (-1) * 0.5)$
 $= 1.5.$
- Hinge loss is **higher** for misclassified sample.



SVM : Estimating w and b

n: no. of data points

- Want to minimize $J(w, b)$:

$$J(w, b) = \frac{1}{2} \sum_{j=1}^d (w^{(j)})^2 + C \sum_{i=1}^n \max \left\{ 0, 1 - y_i \left(\sum_{j=1}^d w^{(j)} x_i^{(j)} + b \right) \right\}$$

Empirical loss $L(x_i, y_i)$

- Compute the gradient $\nabla J(w, b)$ w.r.t. $w^{(j)}$

$$\frac{\partial J(w, b)}{\partial w^{(j)}} = w^{(j)} + C \sum_{i=1}^n \frac{\partial L(x_i, y_i)}{\partial w^{(j)}}$$

$$\frac{\partial J}{\partial b} = C \sum_{i=1}^n \frac{\partial L(x_i, y_i)}{\partial b}$$

$\frac{\partial L(y_i, y_i)}{\partial b} = 0 \quad \text{if } y_i(w \cdot x_i + b) \geq 1$
 $= -y_i \quad \text{otherwise}$

$$\frac{\partial L(x_i, y_i)}{\partial w^{(j)}} = 0 \quad \text{if } y_i(w \cdot x_i + b) \geq 1$$

$$= -y_i x_i^{(j)} \quad \text{else}$$

Gradient Descent

Iterate until Convergence

$$\text{Compute : } \frac{\partial J}{\partial b} = C \sum_1^n \frac{\partial L(x_i, y_i)}{\partial b}$$

$$\text{Update : } b_{new} \leftarrow b - \eta \frac{\partial J}{\partial b},$$

For j=1.....d

$$\text{Compute : } \frac{\partial J}{\partial w^j} = w^j + C \sum_1^n \frac{\partial L(x_i, y_i)}{\partial w^j}$$

$$\text{Update : } w_{new}^j \leftarrow w^j - \eta \frac{\partial J}{\partial w^j},$$

$$b \leftarrow b_{new}, w \leftarrow w_{new}$$

⋮

C parameter in cost function

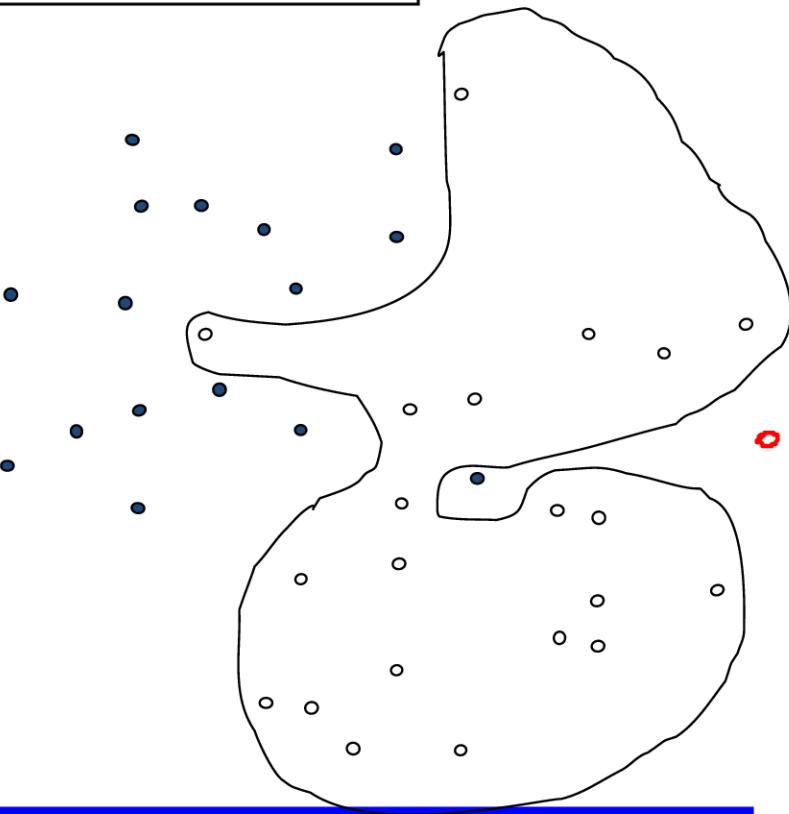
C parameter tells the SVM optimization how much you want to avoid misclassifying each training example.

For large values of C, the optimization will choose a smaller-margin hyperplane if that hyperplane does a better job of getting all the training points classified correctly.

Conversely, a very small value of C will cause the optimizer to look for a larger-margin separating hyperplane, even if that hyperplane misclassifies more points

Dataset with noise

- denotes +1
- denotes -1



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

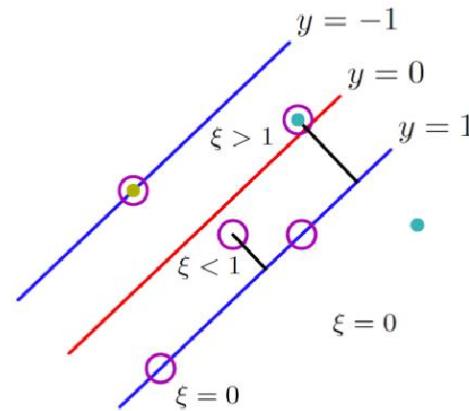
Soft Margin Classification

Slack variables ξ_i can be added to allow misclassification of difficult or noisy examples.

What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \xi_k$$



- **Slack variable** as giving the classifier some leniency when it comes to moving around points near the **margin**.
- When C is large, larger slacks penalize the objective function of SVM's more than when C is small.

Soft Margin

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

The \mathbf{w} that minimizes...

Misclassification cost $\sum_{i=1}^N \xi_i$

data samples N

Slack variable ξ_i

Maximize margin $\frac{1}{2} \|\mathbf{w}\|^2$

Minimize misclassification $C \sum_{i=1}^N \xi_i$

subject to $y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i,$
 $\xi_i \geq 0, \quad \forall i = 1, \dots, N$

Hard Margin versus Soft Margin

- Hard Margin:

Find \mathbf{w} and b such that

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

- Soft Margin incorporating slack variables:

Find \mathbf{w} and b such that

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

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