

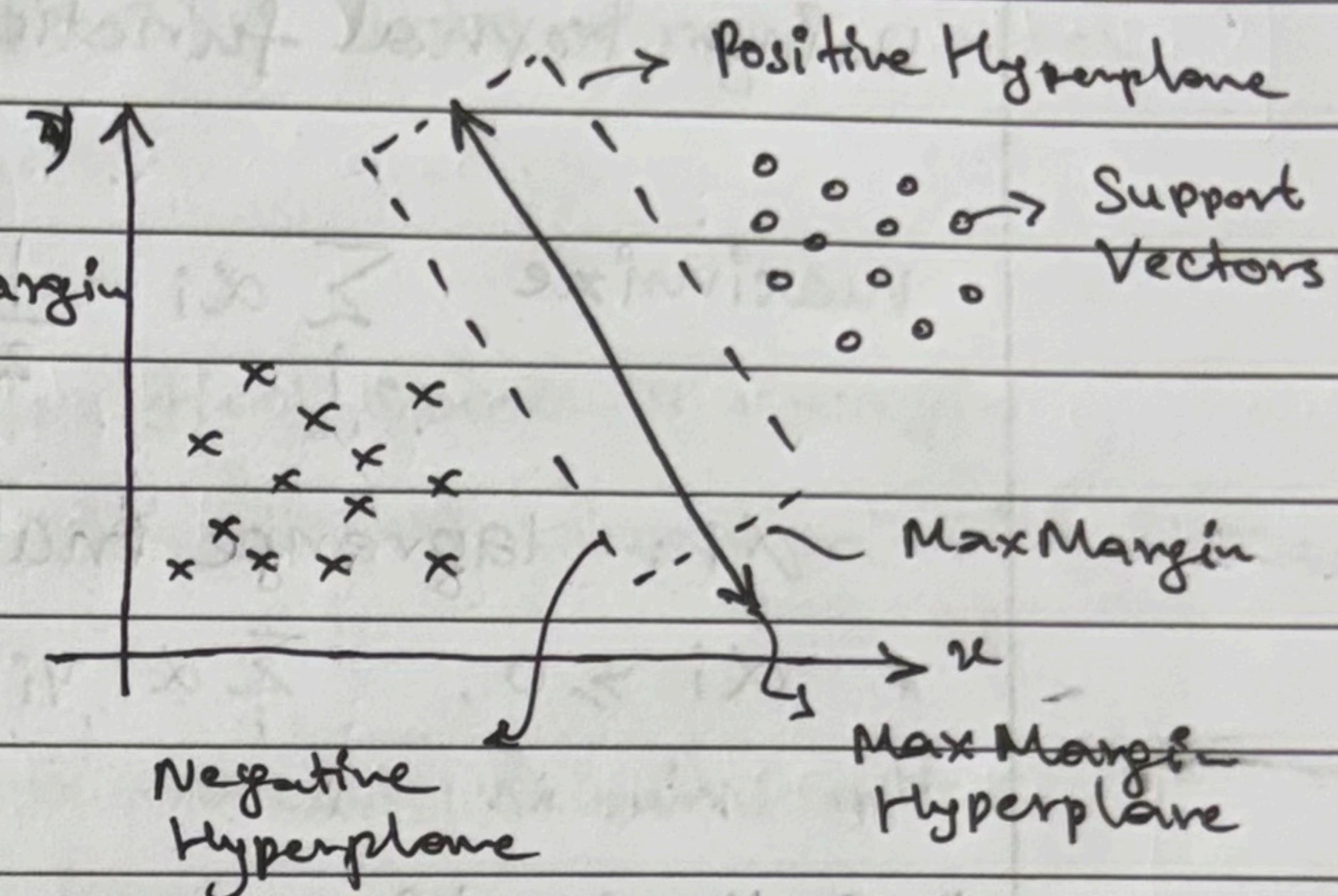
MAXIMISE MARGIN

⑪ Support Vector Machines (SVM)

Goal: Is to maximise the margin b/w the two classes.

Larger the margin, better the model performs on unseen data.

1. The best hyperplane - Hard Margin is the one that maximises the distance b/w hyperplane and nearest data points



Find w, b to: max margin = $2 / \|w\|$ Euclidean Norm

Equivalent Optimization: $\min_{w, b} \frac{1}{2} \|w\|^2$

subject to $y_i (w^T x_i + b) \geq 1, y_i \in \{-1, 1\}$

- * No errors allowed; Ensures each data point is correctly classified.
- * Sensitive to outliers; classified and lies outside margin.

2. A soft margin allows for some misclassifications or violations of margin to improve generalization

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

subject to $y_i (w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$

ξ_i : slack variable - how much the point violates the margin

C : regularization parameter. (trade off b/w margin width and misclassification)

large C : behaves like hard margin

Small C : allows more misclassifications - better generalization

→ Dual Problem in SVM

Dual Problem involves maximizing the Lagrange multipliers associated with the support vectors.

This transformation allows solving the SVM optimization using kernel functions for non-linear classifications.

$$\text{Maximize } \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

α_i : Lagrange multipliers

Kernal function that

$\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$ computes similarity b/w datapoints

* Optimization depends only on dot products $x_i^T x_j$

↳ Allows kernal trick later.

* α_i tells how important each training point is in defining the decision boundary.

If $\alpha_i > 0$: that point lies inside/on the margin - support vector

If $\alpha_i = 0$: point is far away.

* Weight vector w : $\sum_i \alpha_i y_i x_i$

only support vectors ($\alpha_i > 0$) contribute.

Once we solve Dual problem ↗ * Decision Function (Prediction Rule)

↳ After training to classify a new point x ,

$$f(x) = \sum_i \alpha_i y_i (x_i^T x) + b \quad \text{and then}$$

$$\text{Predicted class} = \text{sign}(f(x))$$

The magnitude of $f(x)$ is proportional to the confidence of classification.

Large $|f(x)| \rightarrow$ far from boundary - confident pred.

Small $|f(x)| \rightarrow$ near boundary - uncertain pred

- Primal to Dual -

- 1. Primal Objective (minimize $\frac{1}{2} \|\mathbf{w}\|^2$)
- 2. Dual Form \rightarrow Reframe using Lagrange multipliers
- 3. Solution \sim find ~~key~~ support vectors
- 4. Decision Function, $f(\mathbf{x}) = \sum \alpha_i y_i (\mathbf{x}_i^T \mathbf{x}) + b$: Prediction
- 5. Classification $\sim \text{sign}(f(\mathbf{x}))$ Determines the label

\rightarrow Hinge loss: Penalizes misclassified points or margin violations and is combined with regularization with SVM.

for margin maximization and penalty minimization

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

If $y_i f(x_i) \geq 1$: correctly classified

If $y_i f(x_i) < 1$: inside margin / misclassified

\rightarrow Soft Margin with Hinge loss

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))$$

Margin Maximization \rightarrow Hinge loss

* Linear SVM is thus a linear model trained with hinge loss + L2 Regularization.

(12)

SVM Kernels

In dual form, classifier depends upon only dot products

$$f(\mathbf{x}) = \sum_i \alpha_i y_i (\mathbf{x}_i^T \mathbf{x}) + b$$

We can replace this dot product with a kernel function:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

where $\phi(\cdot)$ maps data to a larger dimension

A kernel is a function that maps data points into a higher dimensional space without explicitly computing the coordinates in that space.

1. Linear Kernel, when data is linearly separable.

$$K(x, y) = x \cdot y$$

2. Polynomial Kernel, allows SVM to model more complex relationships by introducing polynomial terms. $K(x, y) = (x \cdot y + c)^d$

3. RBF (Radial Basis Function) Kernel,

It maps data into an infinite-dimensional space making it highly effective for complex classification problems. $K(x, y) = e^{-\gamma \|x - y\|^2}$

γ : parameter that controls influence of each sample

4. Sigmoid Kernel ~ behave similarly to the activation function of a neuron.

$$K(x, y) = \tanh(\gamma \cdot x^T y + r)$$

→ Basis Expansion (Feature Transformation)

Sometimes we explicitly add non-linear features to make data linearly separable.

- suppose we've a 1D data that is not linearly separable.

We can't separate it using a line in 1D.

- Expand the features: $\phi(x) = [1, x, x^2]$

In 2D/3D it becomes linearly separable.

1. $\phi(x) = [1, x, x^2] \rightarrow$ Non-monotonic function \rightarrow

can make non-separable data, separable

2. $\phi(x) = [1, x^3] \rightarrow$ Monotonic \rightarrow Doesn't change separability.