

Introduction to Data Science

Support Vector Machine ←



Arijit Mondal

Dept. of Computer Science & Engineering

Indian Institute of Technology Patna

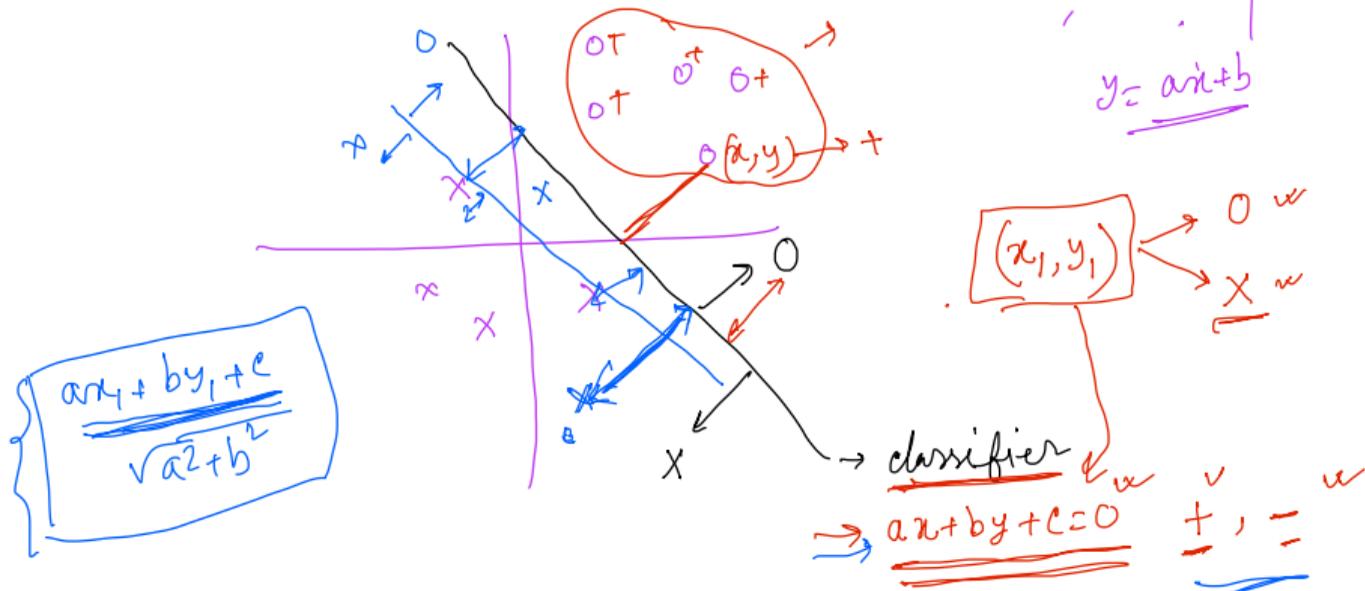
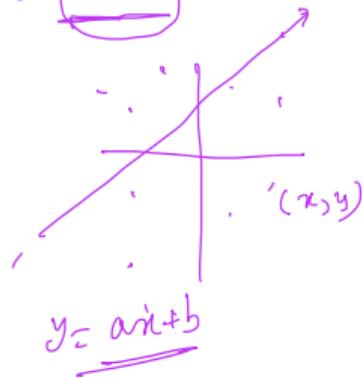
`arijit@iitp.ac.in`

Support Vector Machine ✓

- An approach for classification ✓
- Developed in 1990s ✓
- Generalization of maximum margin classifier | ←
 - Mostly limited to linear boundary
- Support vector classifier — broad range of classes
- SVM — Non-linear class boundary

Hyperplane

- In n dimensional space a hyperplane is a flat affine subspace of dimension $n - 1$
- Mathematically it is defined as
 - For 2 dimensions — $w_0 + w_1x_1 + w_2x_2 = 0$
 - For n dimensions — $w_0 + w_1x_1 + \dots + w_nx_n = 0$



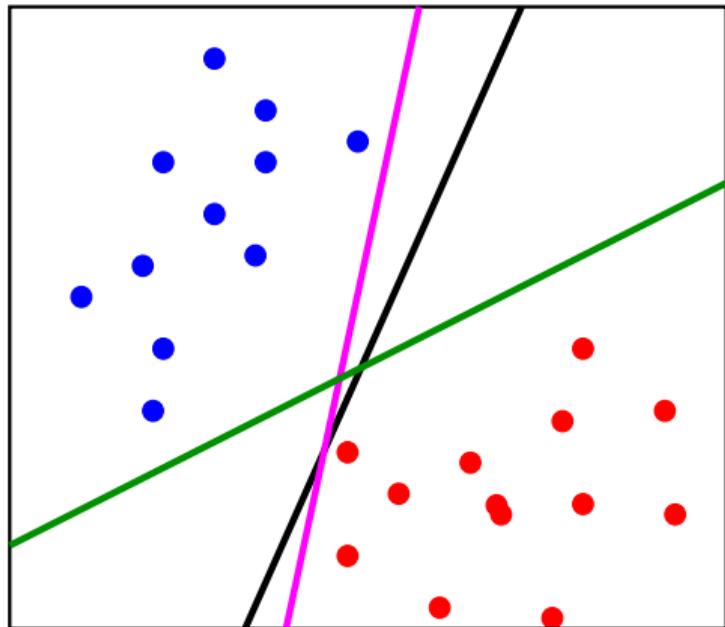
$$\frac{ax_1 + by_1 + c}{\sqrt{a^2 + b^2}}$$

classifier

\Rightarrow $ax + by + c = 0$ $\begin{matrix} + \\ - \end{matrix}$, $\begin{matrix} + \\ - \end{matrix}$

Classification using Hyperplane

- Assume, m training observation in n dimensional space
- Separating hyperplane has the property
 - $w_0 + w_1x_1 + \dots + w_nx_n > 0$ if $y_i = 1$
 - $w_0 + w_1x_1 + \dots + w_nx_n < 0$ if $y_i = -1$



Classification using Hyperplane

- Assume, m training observation in n dimensional space

- Separating hyperplane has the property

- $w_0 + w_1x_1 + \dots + w_nx_n > 0$ if $y_i = 1$
- $w_0 + w_1x_1 + \dots + w_nx_n < 0$ if $y_i = -1$

- Hence, $y_i(w_0 + w_1x_1 + \dots + w_nx_n) > 0$

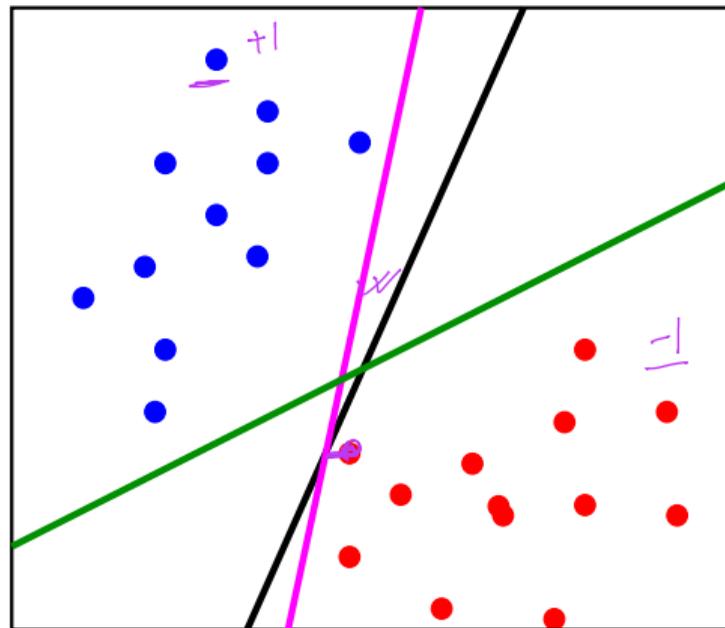
- Classification of test observation x^* is done based on the sign of

$$f(x^*) = w_0 + w_1x_1^* + \dots + w_nx_n^*$$

- Magnitude of $f(x^*)$

- Far from 0 — Confident about prediction

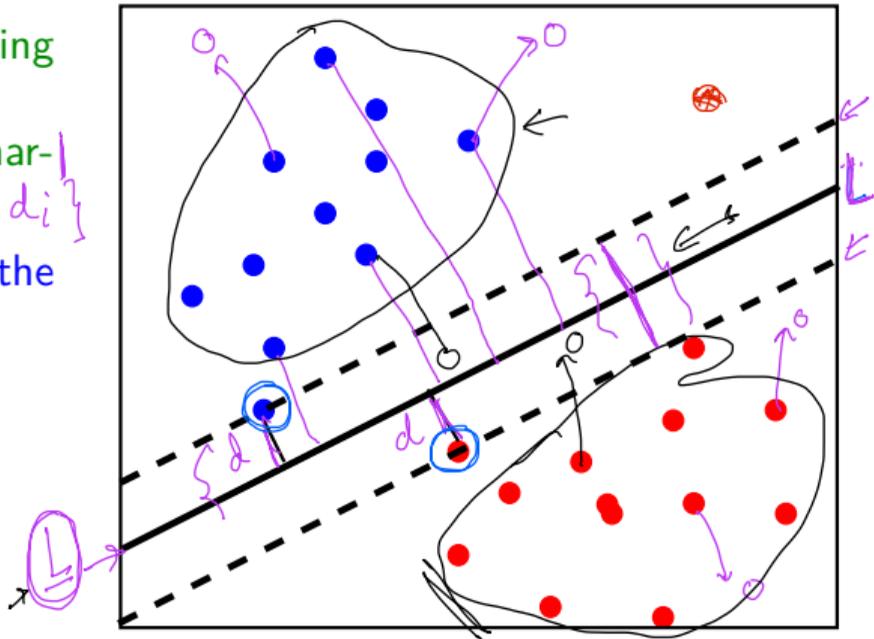
- Close to 0 — Less certain



Maximal margin classifier

- Also known as optimal separating hyperplane
- Separating hyperplane farthest from training observation
 - Compute perpendicular distance from training point to the hyperplane
 - Smallest of these distances represents the margin
- Target is to find the hyperplane for which the margin is the largest

$$\max_L \left\{ \min \{ d_i \} \right\}$$
$$M = \min \{ d_i \}$$



Construction of maximal margin classifier

- Input — m points in n dimension space ie. x_1, x_2, \dots, x_m
- Input — labels y_1, y_2, \dots, y_m for each point x_i where $y_i \in \{-1, 1\}$
- Need to solve the following optimization problem

$$\max_{w_0, w_1, \dots, w_n, M}$$

M

subject to

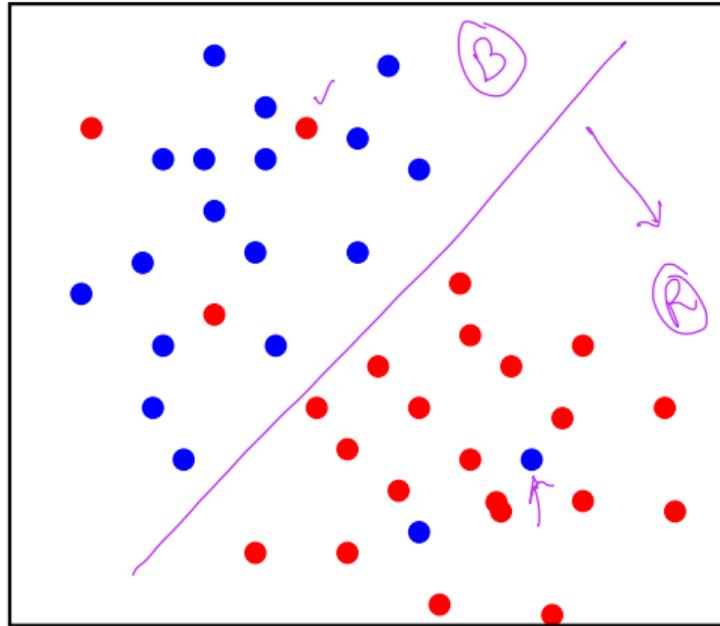
$$y_i(w_0 + w_1x_{i1} + w_2x_{i2} + \dots + w_nx_{in}) \geq M \quad \forall i = 1, \dots, m$$

$$\sum_{i=1}^n w_i^2 = 1$$

$$\frac{ax_1 + bx_1 + c}{\sqrt{a^2 + b^2}} = 1$$

Issues

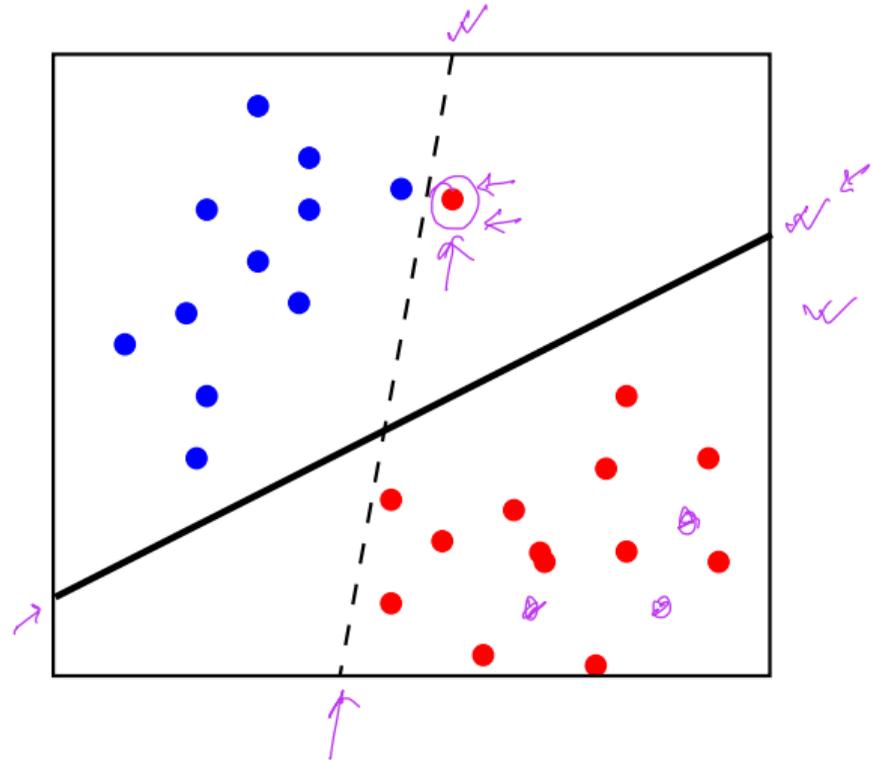
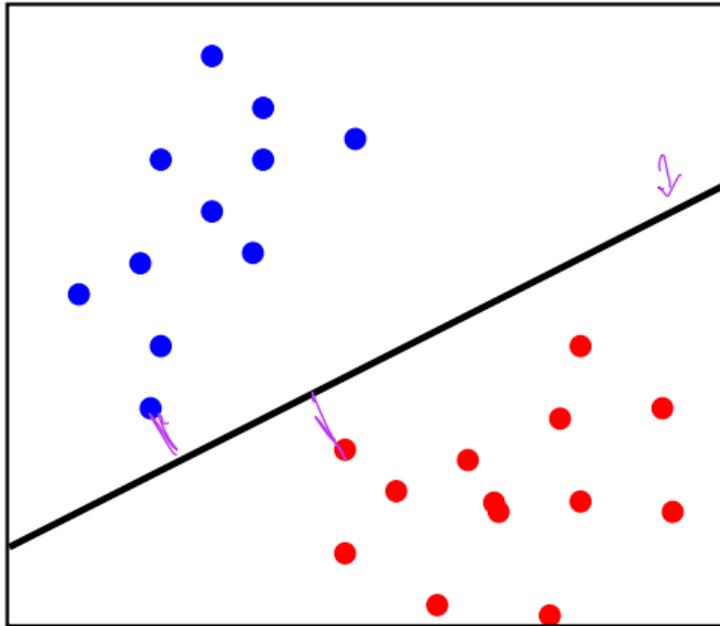
- Maximal margin classifier fails to provide classification in case of overlap



✓

Issues

- Single observation point can change the hyperplane drastically

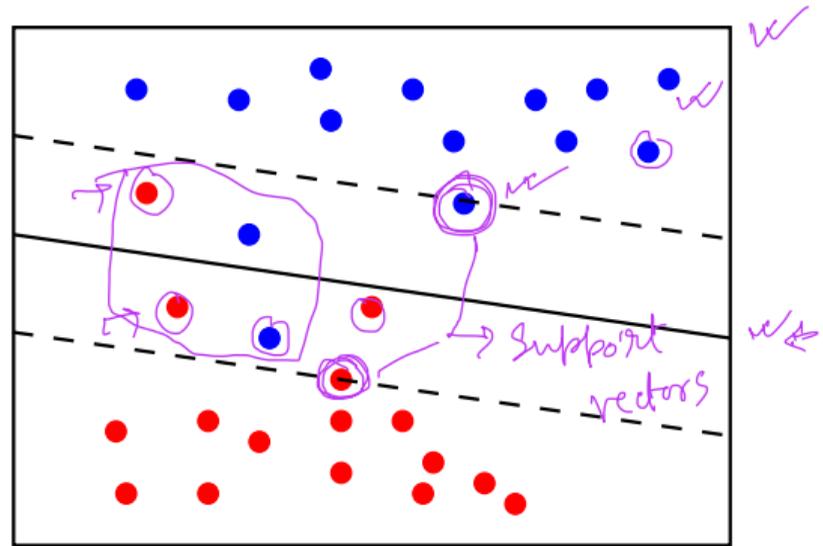
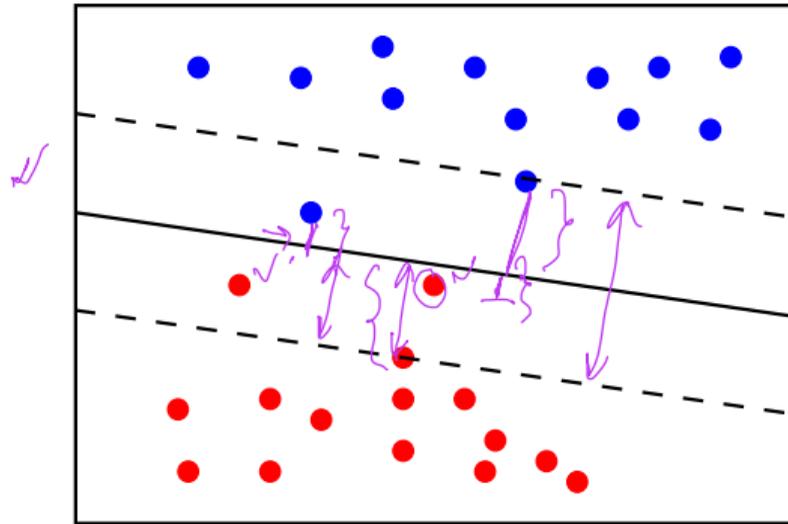


Support Vector Classifier

- Provides greater robustness to individual observations
- Better classification of most of the training observations
- Worthwhile to misclassify a few training observations |
- Also known as soft margin classifier

Support Vector Classifier

- Points can lie within the margin or wrong side of hyperplane

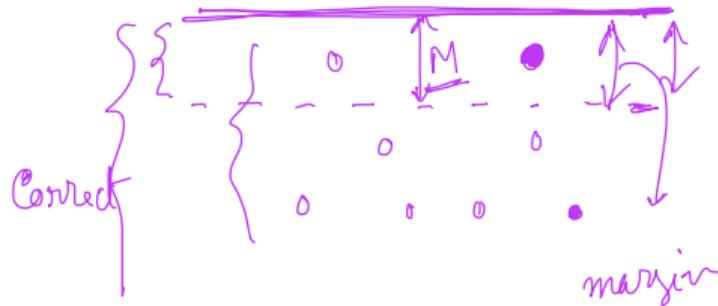
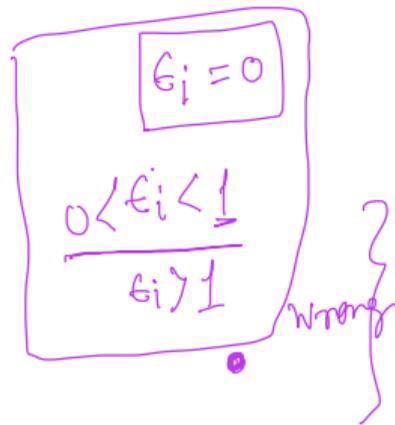


Optimization with misclassification

- Input — x_1, x_2, \dots, x_m and y_1, y_2, \dots, y_m
- Need to solve the following optimization problem

$$\begin{aligned} & \max_{w_0, w_1, \dots, w_n, M} M \\ & \text{subject to} \\ & y_i(w_0 + w_1 x_{i1} + \dots + w_n x_{in}) \geq M(1 - \epsilon_i) \quad \forall i = 1, \dots, m \\ & \sum_{i=1}^n w_i^2 = 1, \quad \sum_{i=1}^m \epsilon_i = C \end{aligned}$$

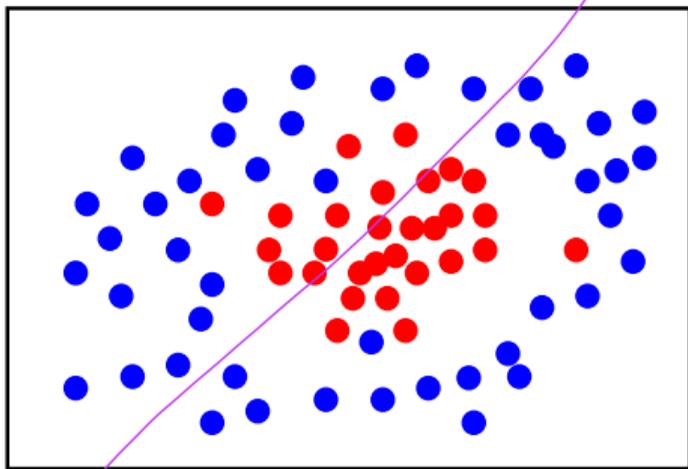
- C is non-negative tuning parameter, ϵ_i - slack variable
- Classification of test observation remains the same



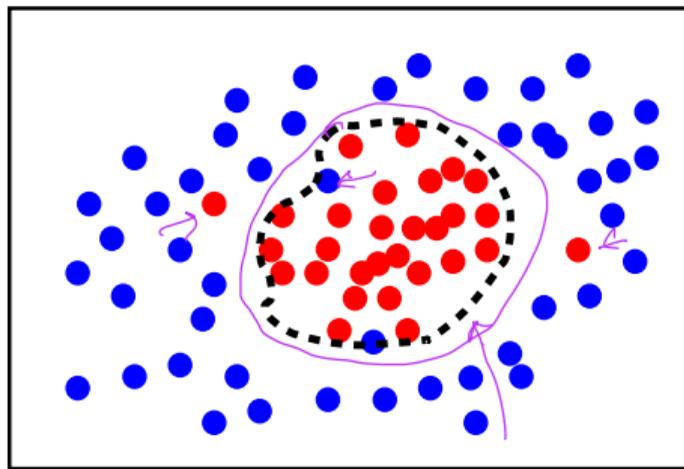
Observations

- $\epsilon_i = 0$ — i th observation is on the correct side of margin
- $\epsilon_i > 0$ — i th observation is on the wrong side of margin
- $\epsilon_i > 1$ — i th observation is on the wrong side of hyperplane
- C — budget for the amount that the margin can be violated by m observations
 - $C = 0$ — No violation, ie. maximal margin classifier ✓
 - $C > 0$ — No more than C observation can be on the wrong side of hyperplane
 - C is small — Narrow margin, highly fit to data, low bias and high variance
 - C is large — Fitting data is less hard, more bias and may have less variance

Classification with non-linear boundaries



bc



α_1^2 , $\alpha_1^2 \alpha_2$

Classification with non-linear boundaries

- Performance of linear regression can suffer for non-linear data
- Feature space can be enlarged using function of predictors
 - For example, instead of fitting with x_1, x_2, \dots, x_n features we could use $x_1, x_1^2, x_2, x_2^2, \dots, x_n, x_n^2$ as features
- Optimization problem becomes

$$\begin{aligned} & \max_{w_0, w_{11}, w_{12}, \dots, w_{n1}, w_{n2}, \epsilon_i, M} \quad M \\ & \text{subject to} \\ & \left(y_i \left(w_0 + \sum_{j=1}^n w_{j1} x_{ij} + \sum_{j=1}^n w_{j2} x_{ij}^2 \right) \right) \geq M(1 - \epsilon_i) \quad \forall i = 1, \dots, m \\ & \sum_{i=1}^n \sum_{j=1}^2 w_{ij}^2 = 1, \quad \sum_{i=1}^m \epsilon_i \leq C, \quad \epsilon_i \geq 0 \end{aligned}$$

Support Vector Machine

- Inner product is replaced with kernel, K or $K(x_i, x_{i'})$
- Kernel quantifies similarity between observations $K(x_i, x_{i'}) = \sum_{j=1}^n x_{ij}x_{i'j}$
 - Above one is Linear kernel ie. Pearson correlation
- Polynomial kernel $K(x_i, x_{i'}) = \left(1 + \sum_{j=1}^n x_{ij}x_{i'j}\right)^d$ where d is positive integer > 1
- Support vector classifier with non-linear kernel is known as support vector machine and the function will look

$$f(x) = w_0 + \sum_{i \in S} \alpha_i K(x, x_i)$$

- Radial kernel: $K(x_i, x_{i'}) = \exp(-\gamma \sum_{j=1}^n (x_{ij} - x_{i'j})^2)$ where $\gamma > 0$