

# Support Vector Machine

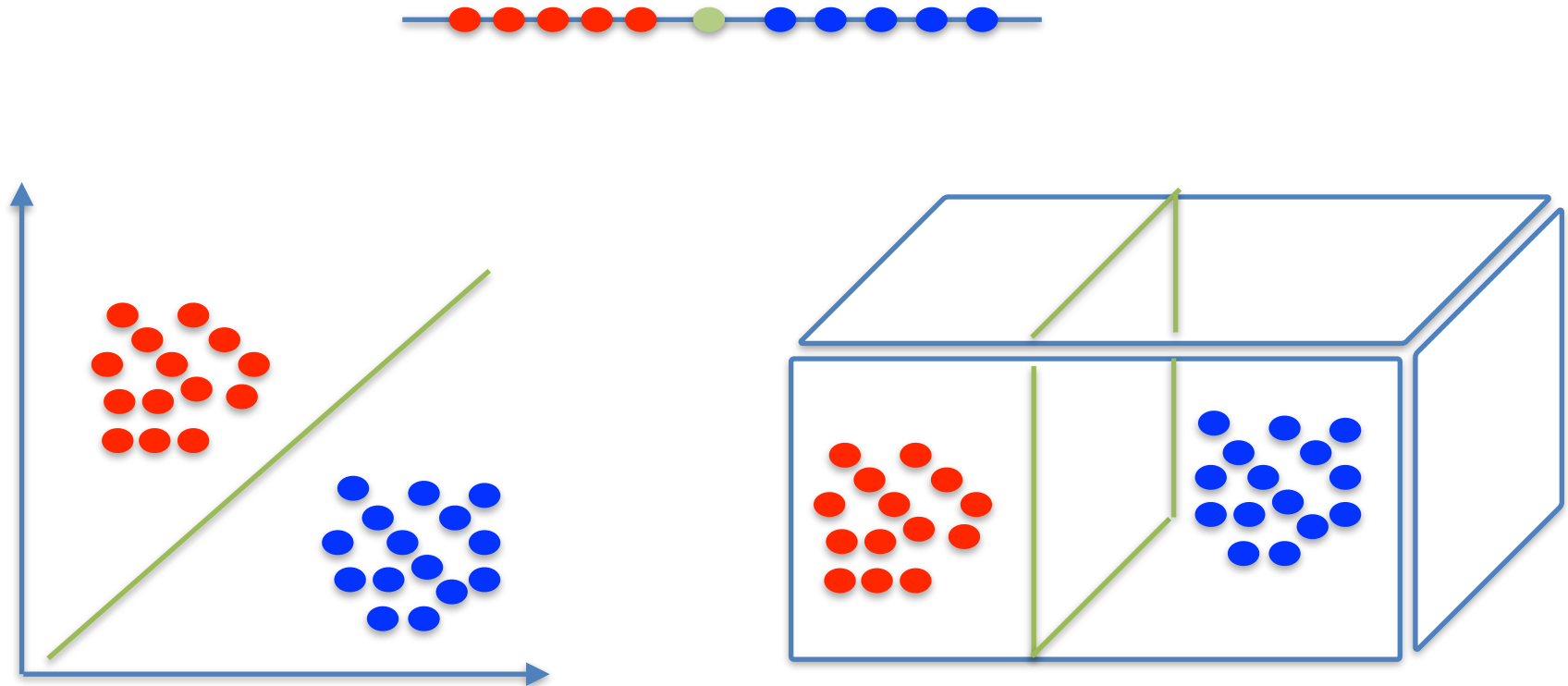
Raman Khurana

PhD in Physics

Researcher at Centre for Deep Learning

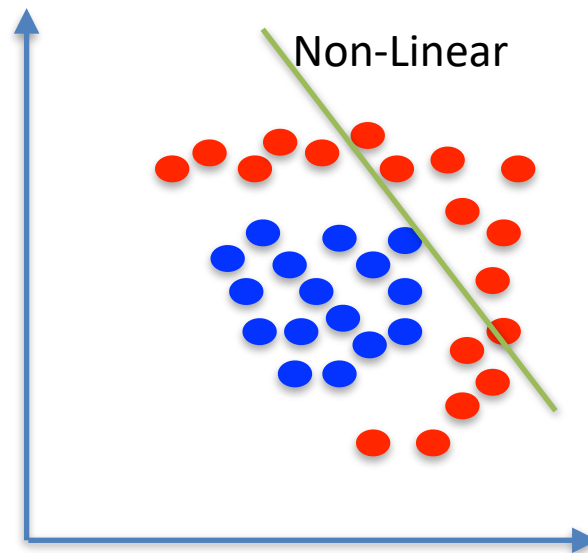
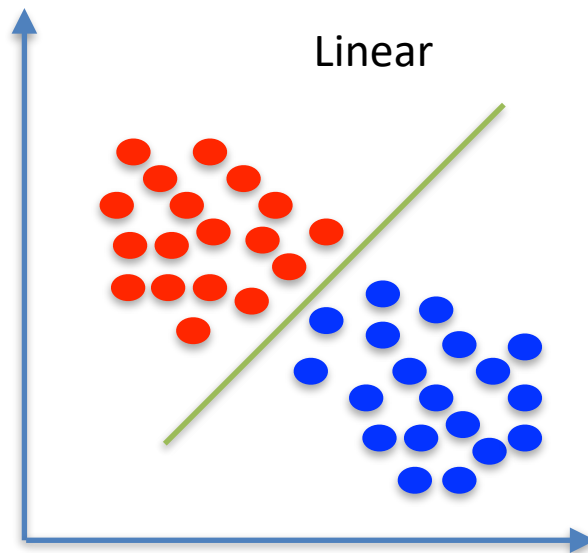
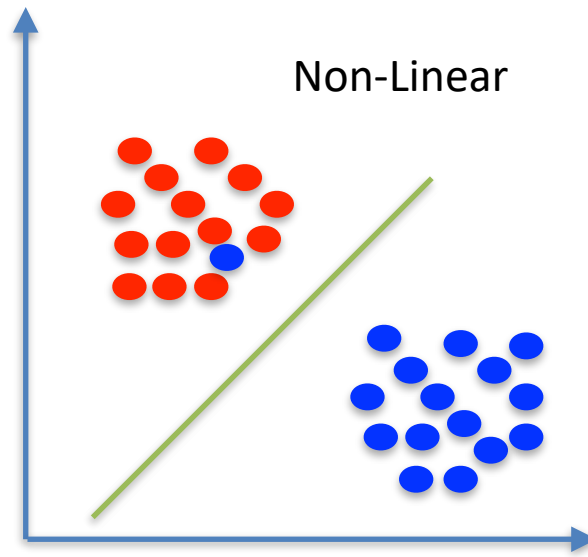
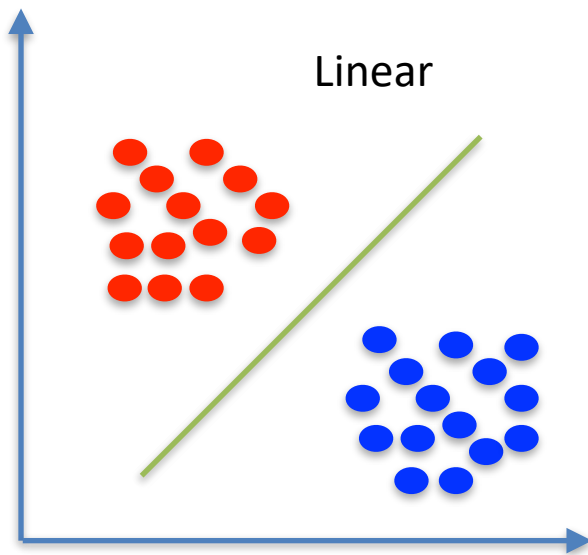
Northwestern

# Classification Problem in n-dimension

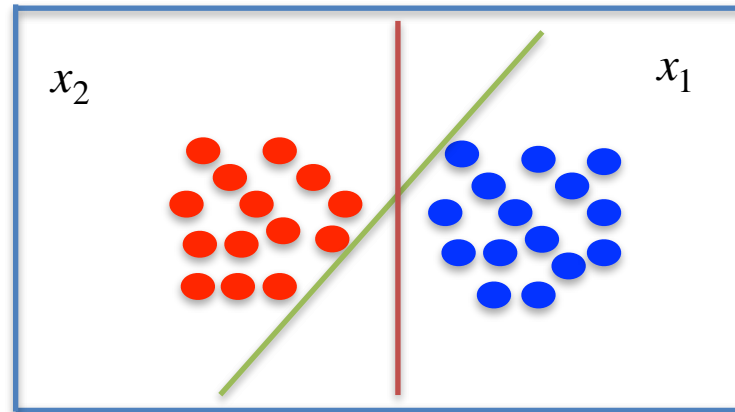


For  $n$ -dimensional feature space one can find a  $n-1$  dimensional hyperplane to classify data.

# Linear separability



# Find the best line

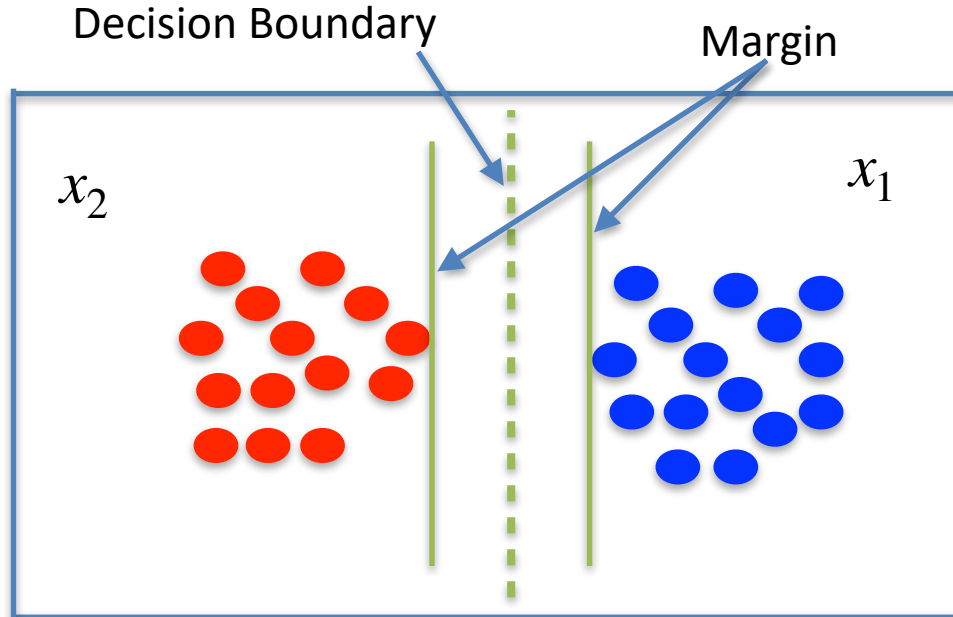


One can draw an infinite straight lines that will classify this data.

But How do I find the best one?

# Support vector machine (SVM): Idea

- Vladimir Vapnik wanted to separate two classes using a straight line.
- Used the widest street approach.
- Place the decision boundary in such a way that the separation b/w classes is the widest.

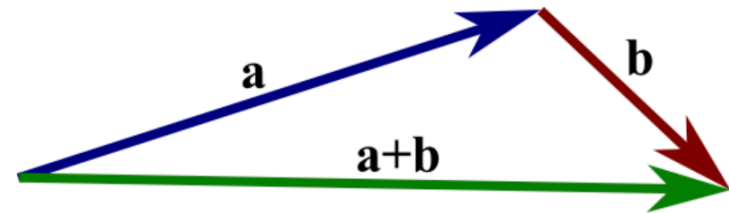


# SVM: Prerequisite

What is a vector?



Triangle law of vector addition.

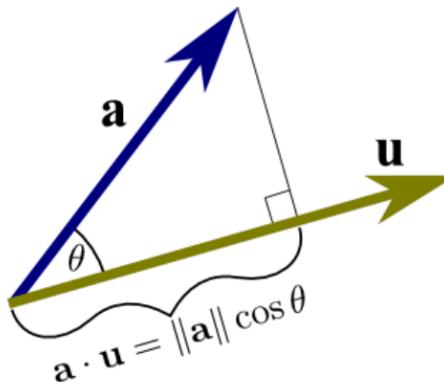


Unit vector.

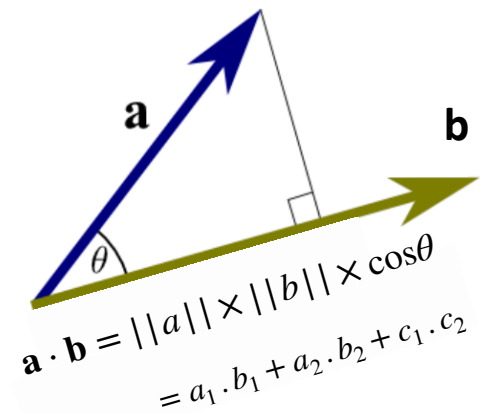
$||b|| = \text{magnitude of } \vec{b}$

$$\hat{b} = \frac{\vec{b}}{||b||} = \mathbf{u}$$

Projection on Unit vector

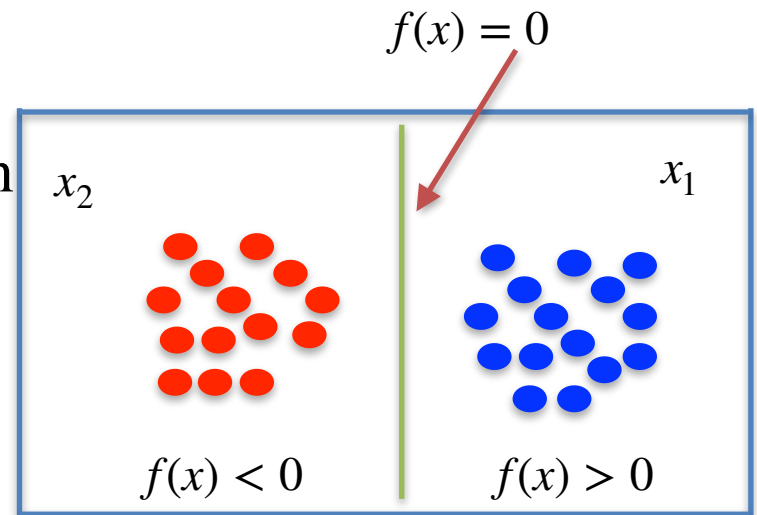


Scalar or dot or inner product.

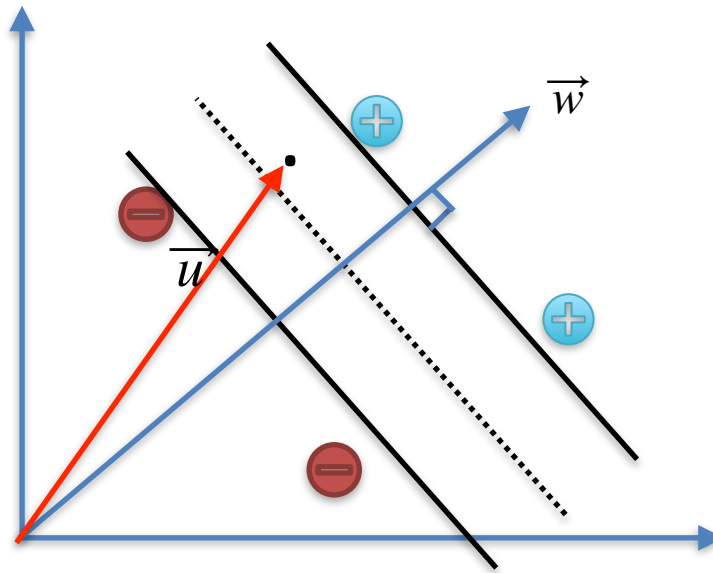


# Linear classifiers

- A Linear classifier in 2d is straight line of the form:  $y = mx + c$ ; eqn of line
  - $\implies ax + by + c = 0$ ; standard eqn
  - If  $w = (a, b)^T$  and  $\mathbf{x} = (x, y)$ ;
  - $w^T \cdot \mathbf{x} + c = 0$
- In 3d space, classifier is a 2d plane with a generic equation  $ax + by + cz + d = 0$ .
  - $w = (a, b, c)^T$  and  $\mathbf{x} = (x, y, z)$ ;
  - $w^T \cdot \mathbf{x} + d = 0$
- In n-dimensional space,
  - $f(x) = w^T \cdot \mathbf{x} + b$ ; b is bias term



# SVM: Decision Rule



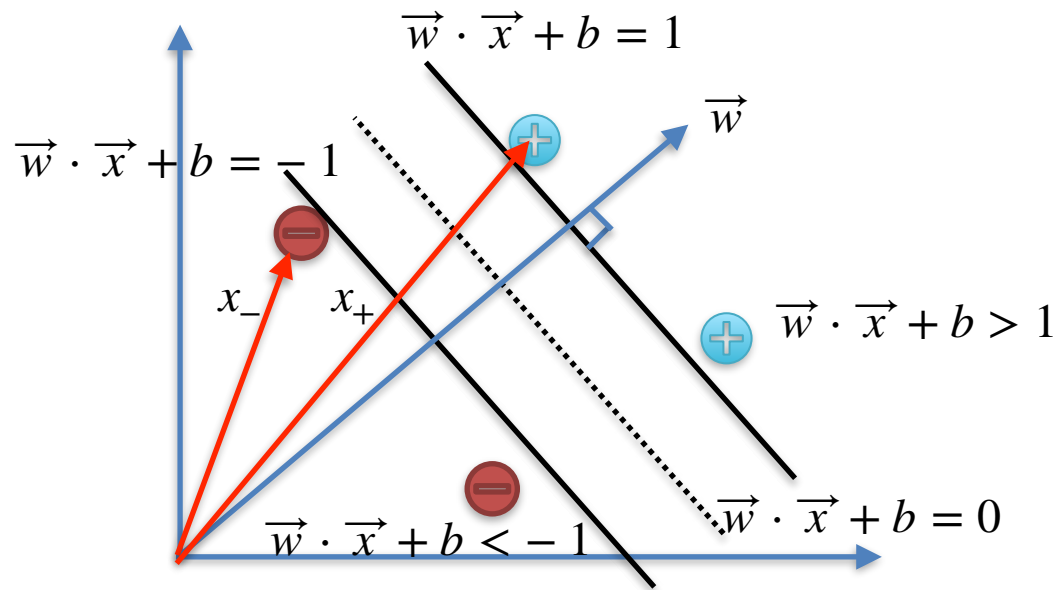
$$\Rightarrow \vec{w} \cdot \vec{u} + b \geq 0, \text{ then it is } \oplus$$

Decision Rule

We don't know what are  $\vec{w}$  and  $b$



# SVM: Constraints



$$\vec{w} \cdot \vec{x}_+ + b \geq +1,$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1$$

For mathematical convenience assume that the separation between  $x_-$  and  $x_+$  is from -1 to +1.

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq +1, \text{ where } x_i \text{ is } x_+, x_-$$

$y_i = +1$  for +ve samples

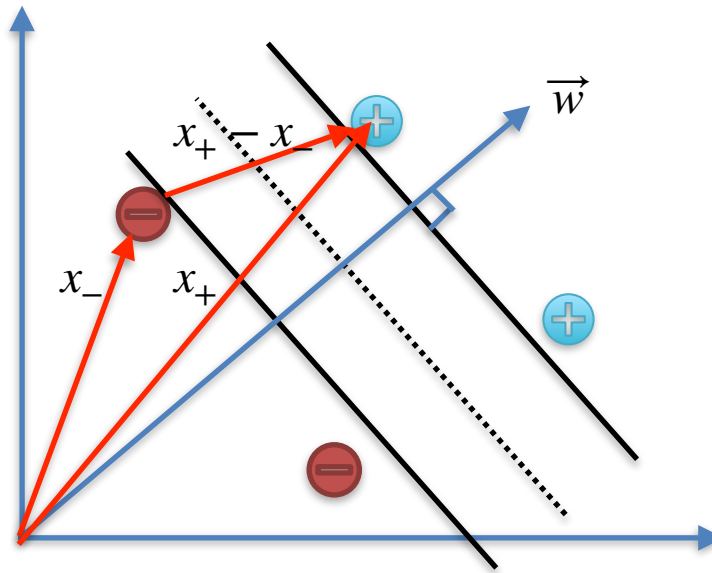
$y_i = -1$  for -ve samples

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 = 0$$

For  $x_i$  to be exactly on the margin line.

# SVM: Width of the street

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 = 0$$



$$\text{Width} = (x_+ - x_-) \cdot \frac{\vec{w}}{||w||}$$

$\swarrow$   $\searrow$   
 $1 - b$   $1 + b$

$$\text{Width} = \frac{2}{||w||} \leftarrow \text{Maximize}$$

$$\implies \text{Minimize } \frac{||w||}{2}$$

$$\implies \text{Minimize } \frac{||w||^2}{2}$$

# Revisit

$$\Rightarrow \vec{w} \cdot \vec{u} + b \geq 0, \text{ where } b = -c \text{ then it is } \oplus$$

Decision Rule

.....Equation 1

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 \geq 0$$

For  $x_i$  to be correctly classified.

.....Equation 2

$$\Rightarrow \text{Minimize } \frac{||w||^2}{2}$$

.....Equation 3

Aim: minimize Equation 3 with constraint in Equation 2.

# Lagrange Multiplier

- Find extremum of a function with constraints
- Formulate Lagrangian using objective function and all the constraints.

Plan of action

$$L = \text{function to maximize (or minimize)} - \sum \alpha_i \cdot R_i$$



Take partial derivative w.r.t each variable and set to ZERO



Rearrange expression of L using all the new information

# Simplifying the optimization problem

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

$\alpha$ 's are non-zero for samples on the margin. For all other samples  $\alpha$  are zero. These are called support vectors.

$$\frac{\delta L}{\delta \vec{w}} = \vec{w} - \sum \alpha_i y_i \vec{x}_i = 0$$

$$\Rightarrow \vec{w} = \sum \alpha_i y_i \vec{x}_i$$

$$\frac{\delta L}{\delta b} = \sum \alpha_i y_i = 0$$

$$L = \frac{1}{2} (\sum \alpha_i y_i \vec{x}_i) \cdot (\sum \alpha_j y_j \vec{x}_j) - (\sum \alpha_i y_i \vec{x}_i) \cdot (\sum \alpha_j y_j \vec{x}_j) - \sum \alpha_i y_i + \sum \alpha_i$$

$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

Exercise: Find expression for bias term, b

$$b = y_i - \sum_j \alpha_j y_j (\mathbf{x}_j \cdot \mathbf{x}_i)$$

# What do we have so far?

$$\Rightarrow \vec{w} \cdot \vec{u} + b \geq 0, \text{ where } b = -c \text{ then it is } \oplus$$

Decision Rule

.....Equation 1

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 \geq 0$$

For  $x_i$  to be correctly classified.

.....Equation 2

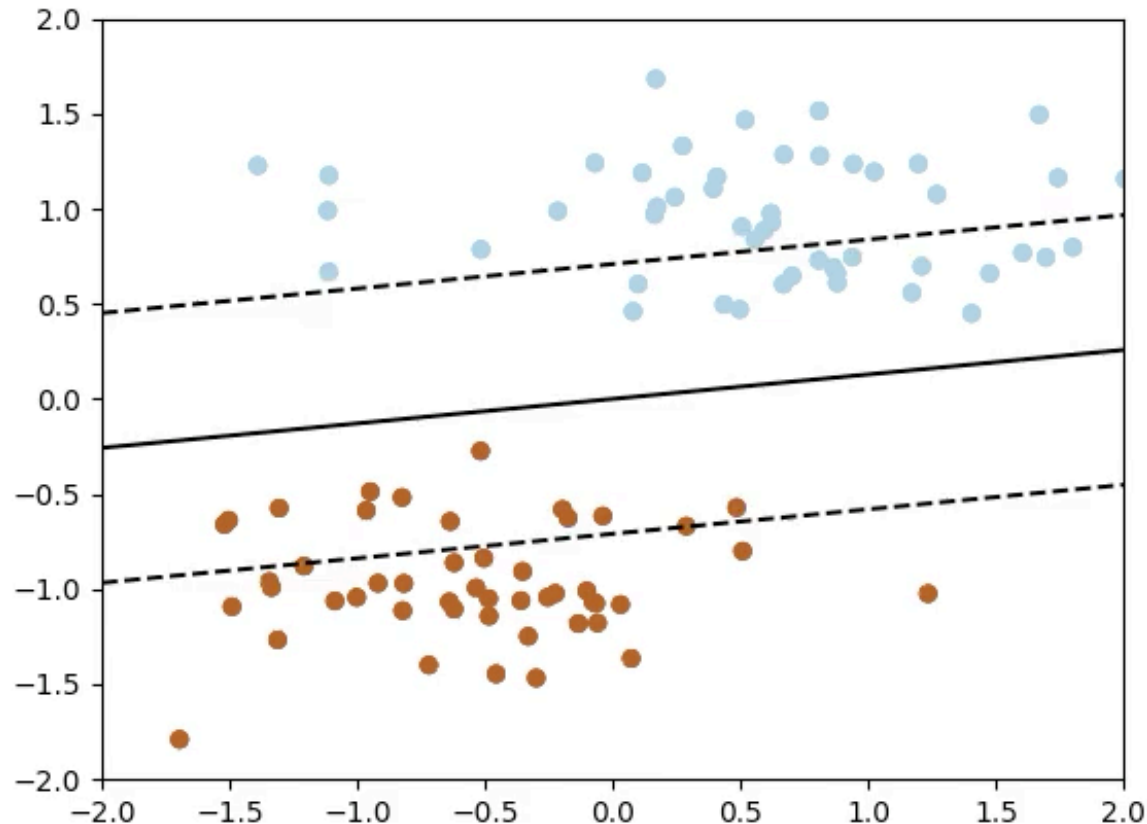
$$\Rightarrow \text{Minimize } \frac{||w||^2}{2}$$

.....Equation 3

$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j \quad \dots \text{Equation 4}$$

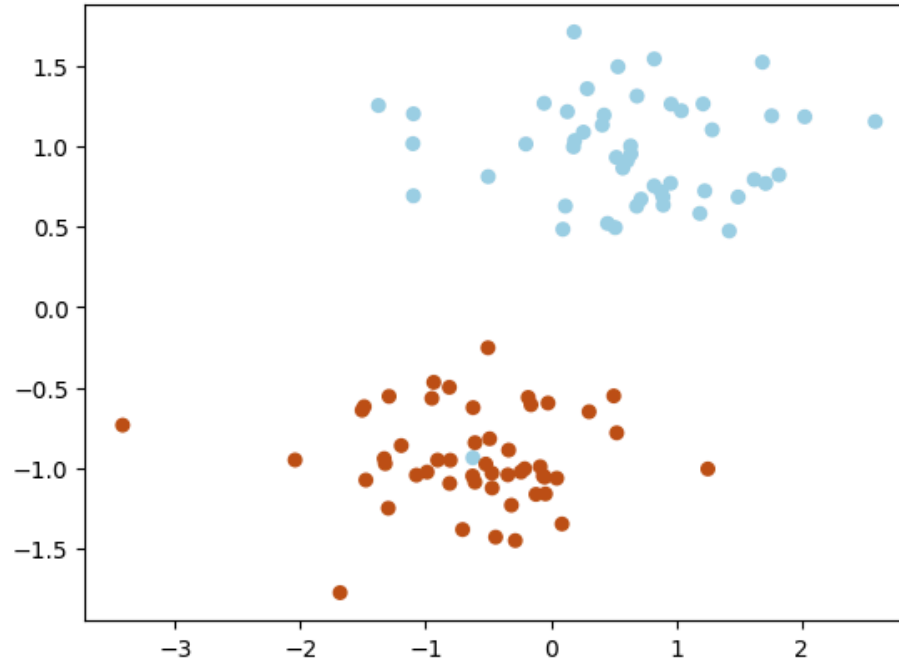
The problem is now simplified to find maximum of equation 4, i.e. optimization depends only on the dot product of pair of samples

# SVM at action



Support vector are the data points/samples which lie on the boundary of margin.

# Noisy/Non-linear data



- Hard margin will never converge due to shifting of data point.
- Model will never find a minima and hence no decision boundary to completely separate two classes.

Practical solution: Ignore this one data point, by allowing a certain degree of misclassification.



# Soft Margin

$\Rightarrow \vec{w} \cdot \vec{u} + b \geq 0$ , where  $b = -c$  then it is  $\oplus$

Decision Rule

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 \geq 0$$

$$y_i(\vec{w} \cdot \vec{x}_i + b) - (1 - \zeta_i) \geq 0$$
$$\zeta_i > 0 \forall i$$

$$\Rightarrow \text{Minimize } \frac{||w||^2}{2}$$

$$\Rightarrow \text{Minimize } \frac{||w||^2}{2} + C \sum_{i=1}^n \zeta_i$$

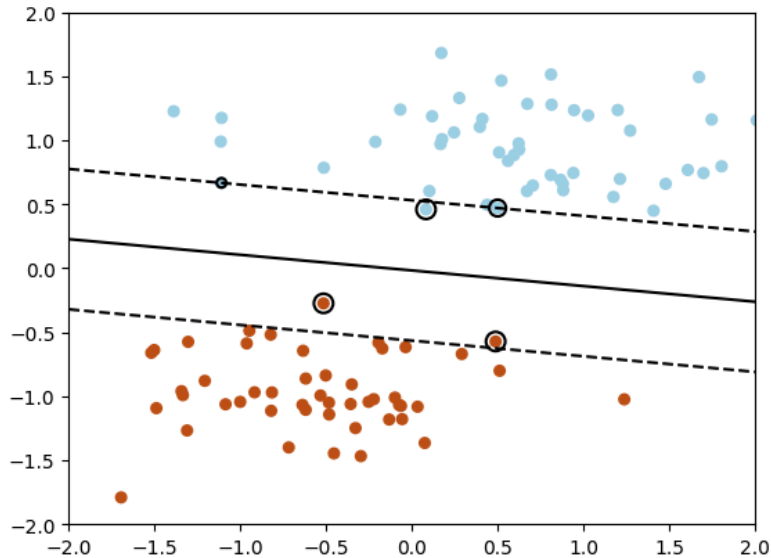
$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

# Soft Margin

- Two scenarios

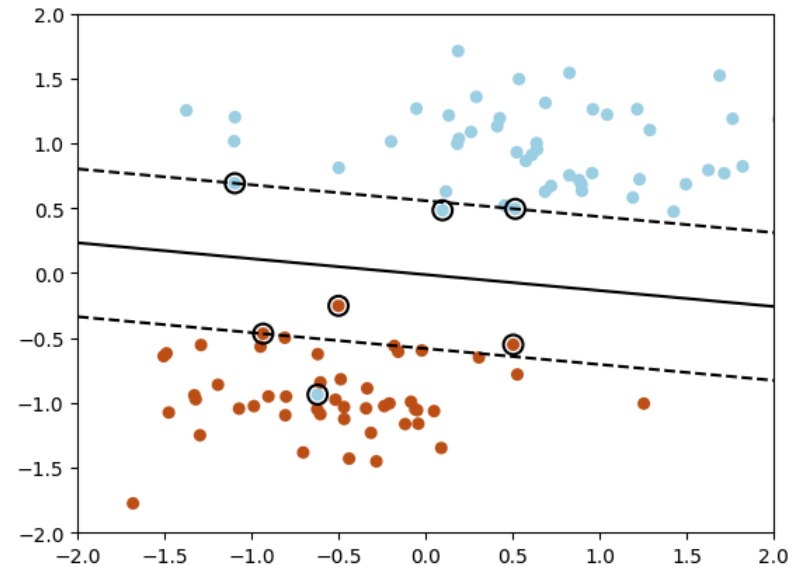
- Correctly classified data point within the margin

- $0 < \zeta_i < 1$



- Misclassified data point on wrong side of margin

- $\zeta_i \geq 1$



Support vector are the data points/samples which lie **on the decision boundary or in the margin or mis-classified.**

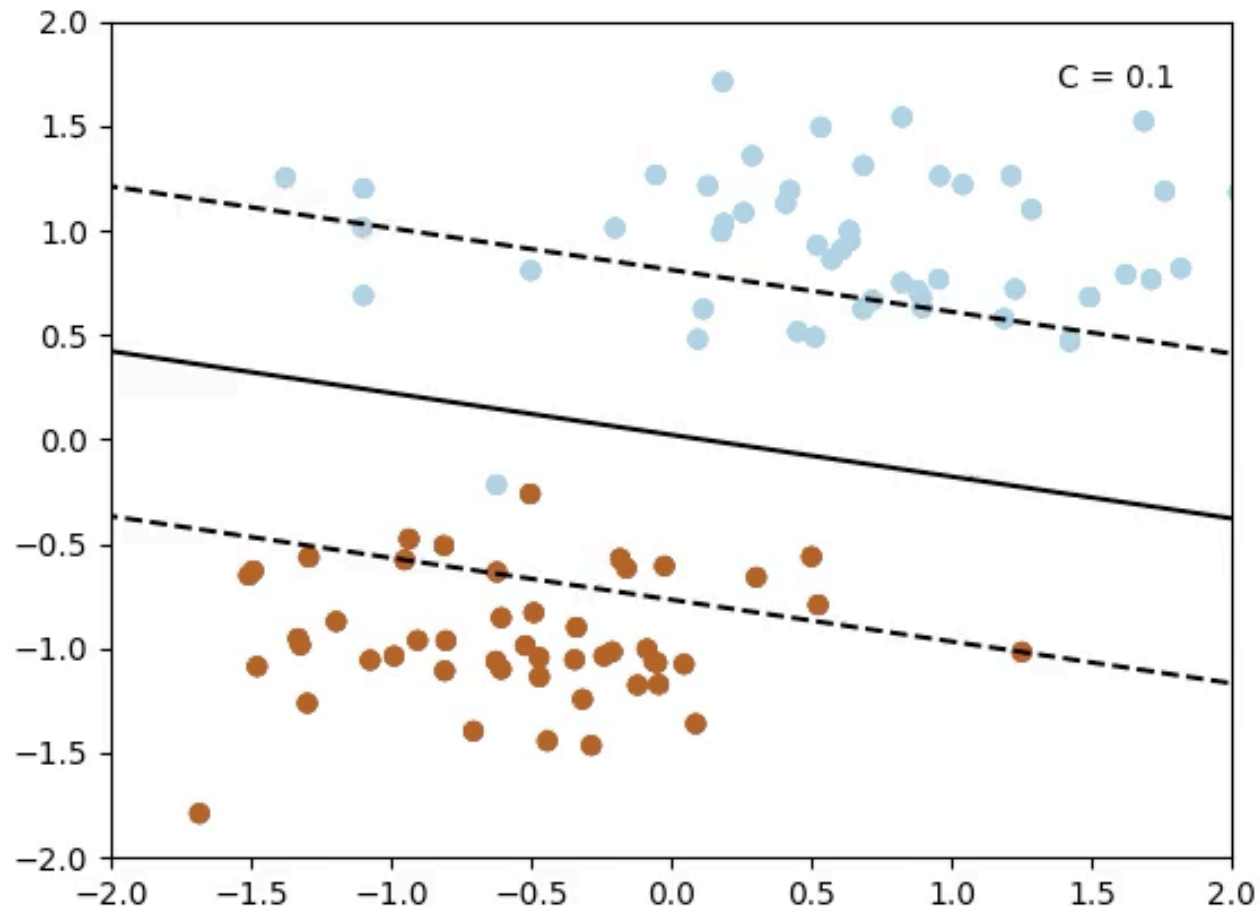
# Importance of C

- C is the hyper parameter of model which controls tradeoff b/w margin width and classification error.

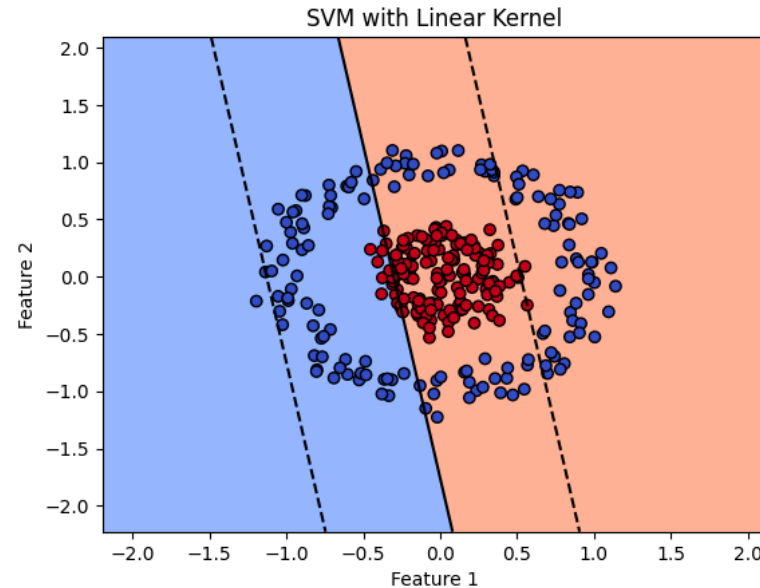
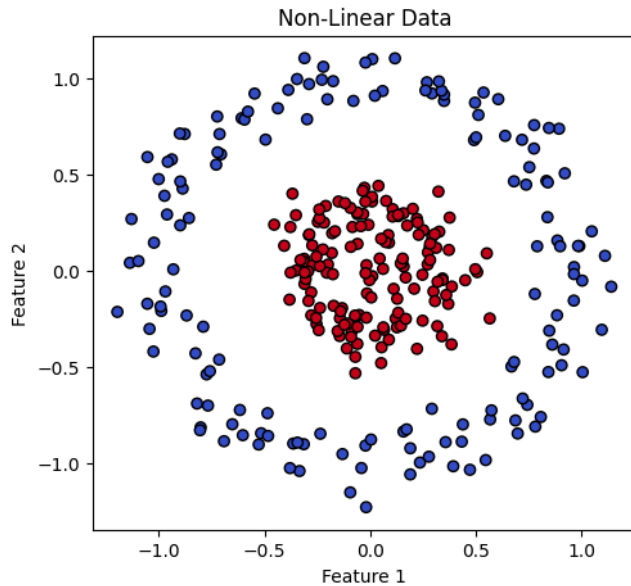
$$\text{Minimize } \frac{||w||^2}{2} + C \sum_{i=1}^n \zeta_i$$

- After introducing slack variable in the Lagrangian it is feasible to find optimized solution even if some data points are misclassified.
- Large C: classification error has more weightage  $\implies$  model tries to minimize mis-classification of classes  $\implies$  the margin width can be smaller  $\implies$  sensitive to small variation in the (unseen) data  $\implies$  can not be generalized on unseen data [**Overfitting**]
- Small C: allows misclassification of class  $\implies$  model will tolerate errors for larger margin width  $\implies$  better generalization [**less likely to overfit, but more likely to underfit**]

# Importance of C



# Kernel Trick



A linear classifier will not be generalized enough to deal with non-linear data.

$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

$$\vec{x}_i \rightarrow \phi(\vec{x}_i) \text{ and } \vec{x}_j \rightarrow \phi(\vec{x}_j)$$

$$\vec{x}_i \cdot \vec{x}_j \rightarrow \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

$$\vec{x}_i \cdot \vec{x}_j \rightarrow K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

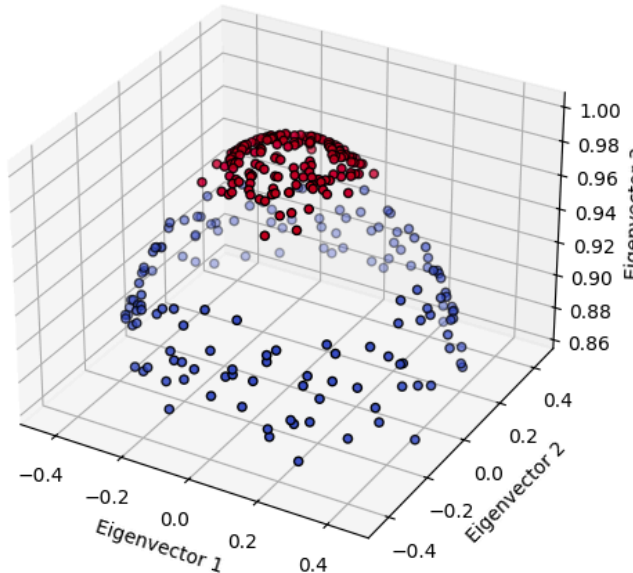
# Kernel Trick

$$\vec{x}_i \cdot \vec{x}_j \rightarrow K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

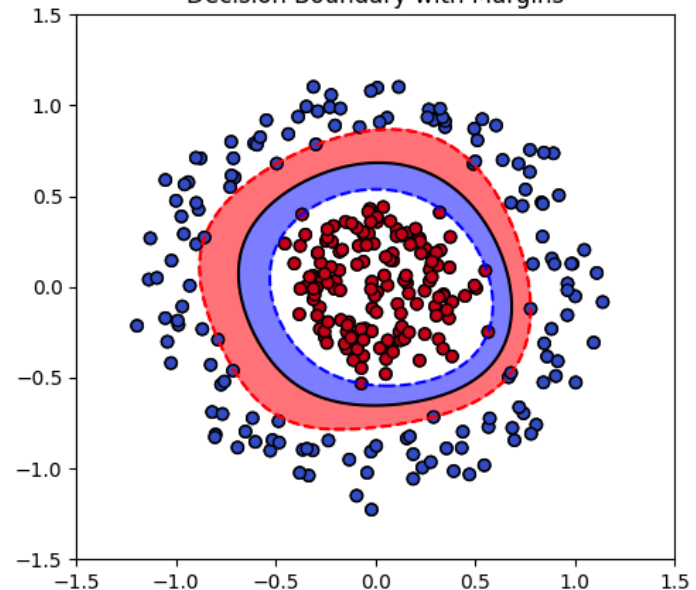
Kernel function

- **Kernel Trick:** Kernel Function computes the dot product of two vectors in a potentially higher dimensional feature space without explicitly performing the transformation into that higher dimensional space.

RBF Kernel Transformation



Decision Boundary with Margins



# Common Kernels

- **Linear:**

$$K(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

- The simplest kernel to be used when data is linearly separable.
- # of features > # of data points.

- **Polynomial:**

$$K(x_i, x_j) = (x_i \cdot x_j + c)^d; c \text{ is constant and } d \text{ is degree of polynomial.}$$

- Use for non-linear case; complexity of model controlled by parameter  $d$

- **Radial Basis Function (RBF) or Gaussian Function:**

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

- Another non-linear kernel transforming data to infinite dimensional space.
- Transformation to Infinite dimensional space allows the modeling of complex data.

# Key points

- Pros
  - Dataset size is small
  - # of features  $>$  # of samples
  - Robust (not immune) to overfitting; if used carefully
  - Memory efficient; uses only a subset of samples
  - Kernel trick; handles non-linear data.
  - Effective for binary classification
  - Better interpretability
- Cons
  - When dataset is very large
  - # of feature  $\gg$  # of samples
  - Choice of right Kernel
  - Computational expensive for large dataset or large number of features.



# Learnings and take aways

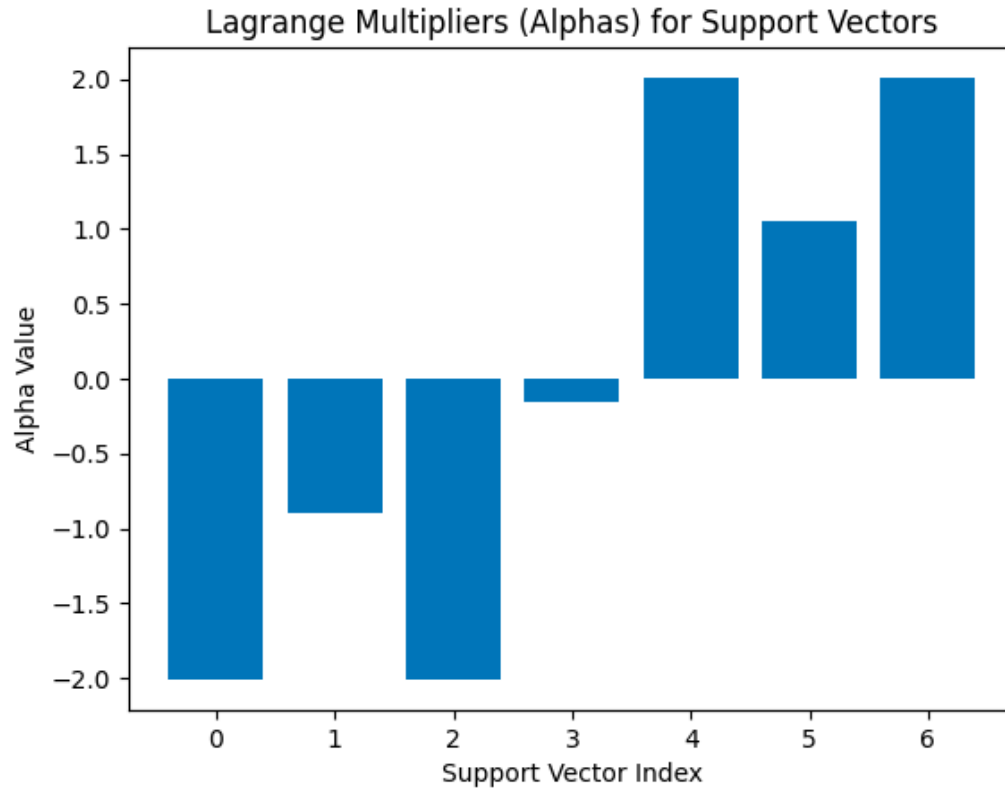
- A powerful supervised learning algorithm used primarily for classification task.
- Finds the best hyperplane that best separates the classes.
  - A decision boundary separates classes in the feature space.
  - For 2-d data, it's a line; for 3-d data it's a plane; and for higher dimension it's a hyperplane.
- The Lagrangian/Objective function depends on the dot product of pair of vectors.
- Data points closest to the hyperplane on either side constitute the Support vector.
  - Drives the position and orientation of hyperplane.
- The margin of the hyperplane is defined as the distance between the hyperplane and the closest support vector.
  - SVM aims to maximize this margin width ensuring better generalization on unseen data.
- Kernel trick transform the non-linearly separable data into higher dimensional space where it become linearly separable.
- A slack variable ( $\zeta$ ) allows for limited misclassification, making the model more robust against noisy data or a few non-linearly separable data points.
- Regularization parameter (C) controls the tradeoff between margin width and classification error.

# Further reading

- [Support Vector Networks](#), Corinna Cortes & Vladimir Vapnik (1995)
- [The Nature of Statistical Learning Theory](#) :
- [SVM in sklearn](#)
- [A tutorial on Support Vector Regression:](#)  
[Alex J Smola](#)
- [Support Vector Machine: Hype or Hallelujah](#)



# Dual coefficients



```
1 dual_coefs
```

```
✓ 0.0s
```

```
array([-2.01      , -0.89680413, -2.01      , -0.15962996,  2.01      ,  
       1.05643409,  2.01      ])
```

```
1 np.sum(dual_coefs)
```

```
✓ 0.0s
```

```
4.440892098500626e-16
```

$$\frac{\delta L}{\delta b} = \sum \alpha_i y_i = 0$$

Class weighted sum of Lagrange Multipliers is zero.

# Kernel Trick



$$\vec{x}_i \cdot \vec{x}_j \rightarrow K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

- The kernel trick is to perform the dot product of pair of vector in a higher dimensional space without actually transforming the data points itself.
  - Saves computation when there is huge data to be transformed.
- Why do we need generic kernels?
  - Data is not always simple to transform using known/simple mathematical equations.
  - Generalized kernels help them separate in higher dimensional without knowing the exact mathematical formula.
    - However, you still need to know something about the data to make decision about the kernel.

# Dimensionality in Kernel space

- Polynomial:

- For  $n$  number of features and  $d$  as degree of polynomial the dimension of feature space after polynomial transformation is:  
$$\frac{(n + d)!}{d!}$$

- Gaussian:

- For any value of  $n$  and width of gaussian function the dimension of feature space after RBF transformation is INFINITE.
- Gaussian function can be expressed as infinite sum using Taylor Series.

$$e^{-x^2} = \sum_{n=0}^{\infty} \frac{(-1)^n (x^2)^n}{n!}$$

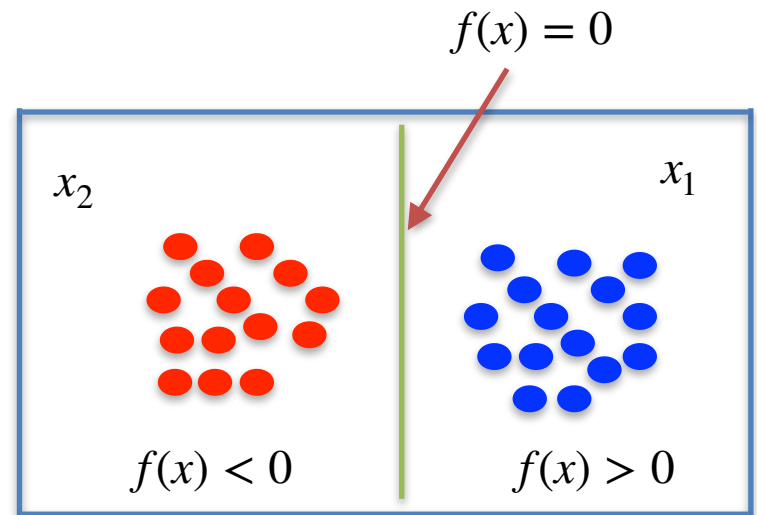
# Self questions

- Ask all type of What, When, Where, How questions and find answer to them?
- Is it possible to know the dimension of the higher dim kernel space where data is transformed and transformation of the dot product is performed.
- What are assumptions for SVM?
  - Model, kernel trick, and more assumptions associated with SVM itself.
- Ask ChatGPT: Be my interviewer who is an expert in Statistics and technical. Ask 30 question on SVM and keep a few important questions on nearby topics in ML. Start with basic, and move to advanced level in steps.
- What are other kernel model than SVM?
-

# Linear classifiers

- A Linear classifier has the form

$$f(x) = w^T x + b$$



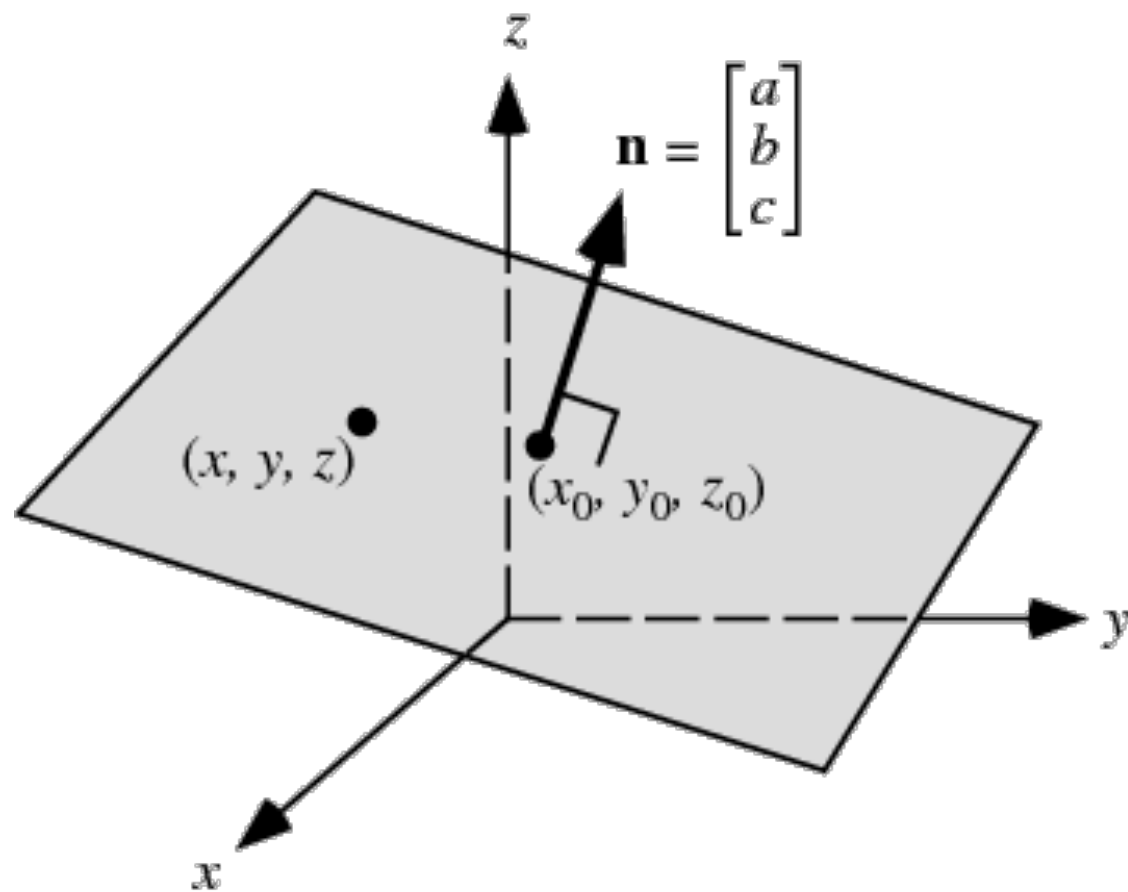
- In 2d the discriminant is a line.
- $w$  is the **normal** to the discriminant line, and  $b$  is **bias**.
- $w$  is also known as **weight vector**.



# Hinge loss

# SVM vs Logistic regression

# Multiple classes



# Find the best line

