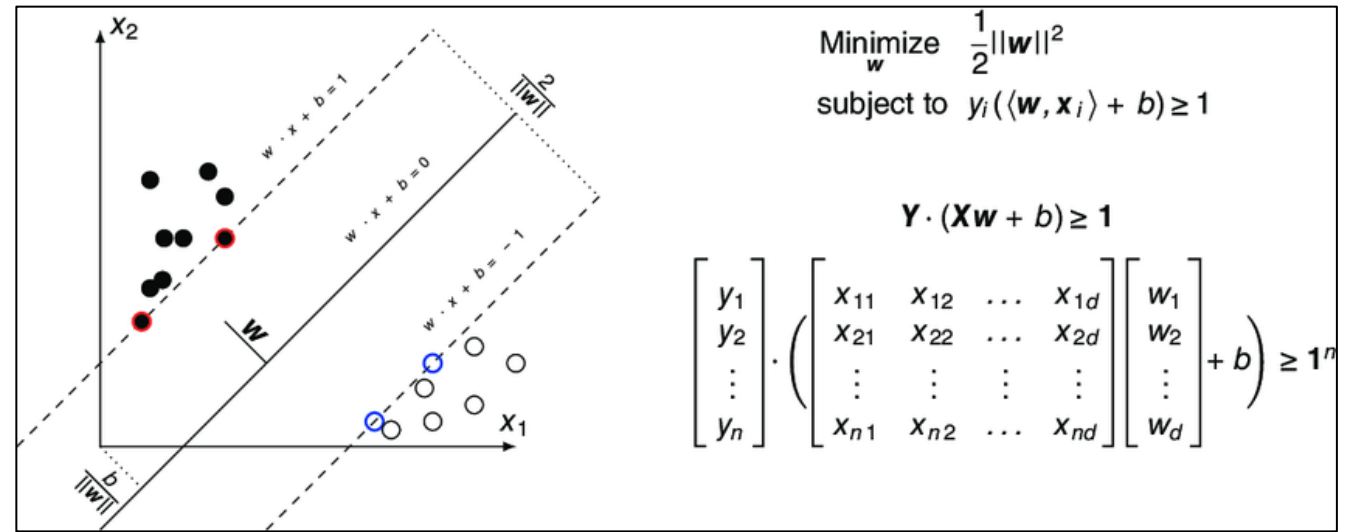
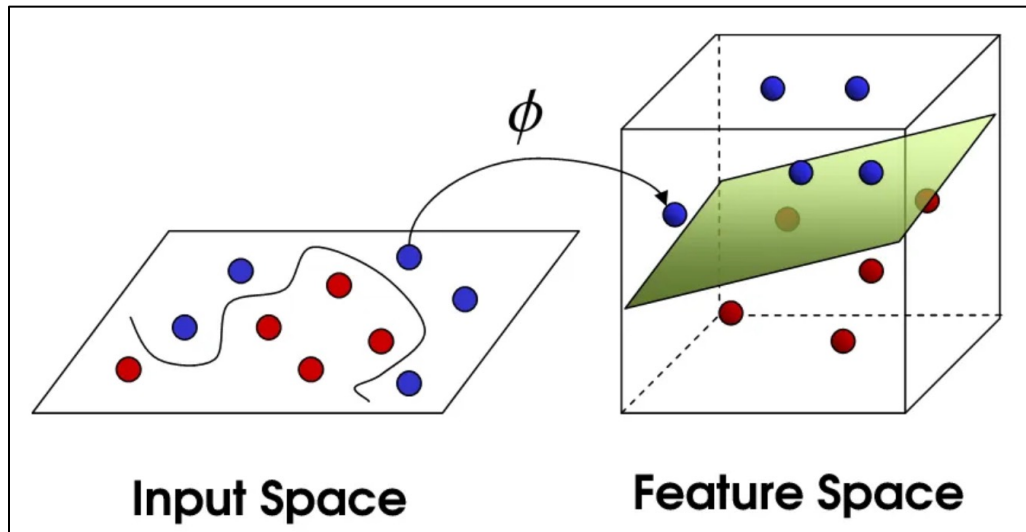


Support Vector Machine (First Look)



Outline

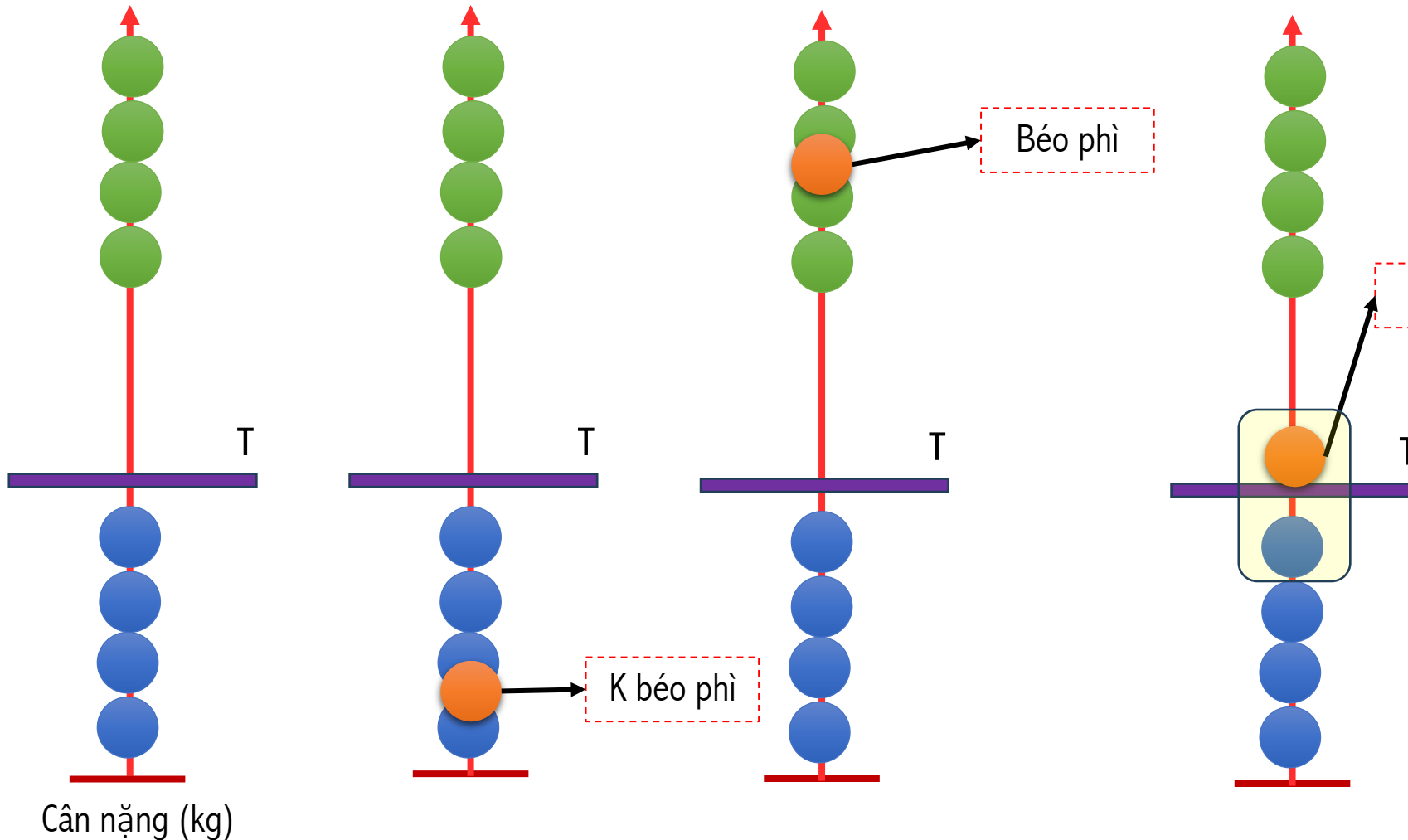
- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

Outline

- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

SVM Motivation

Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)



- Béo phì
- Không Béo phì
- Dữ liệu mới

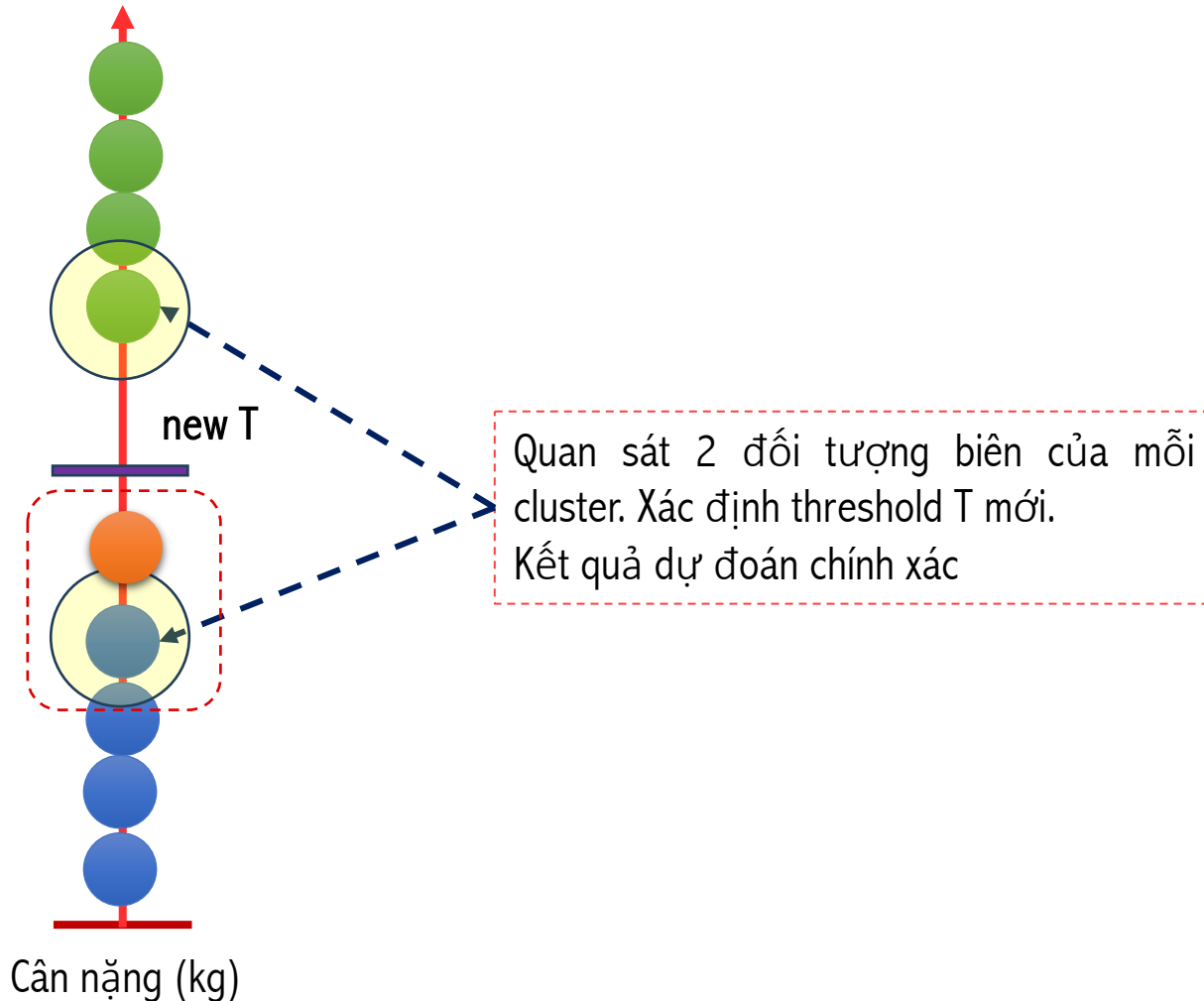
Có chính xác không?

Threshold T không chính xác

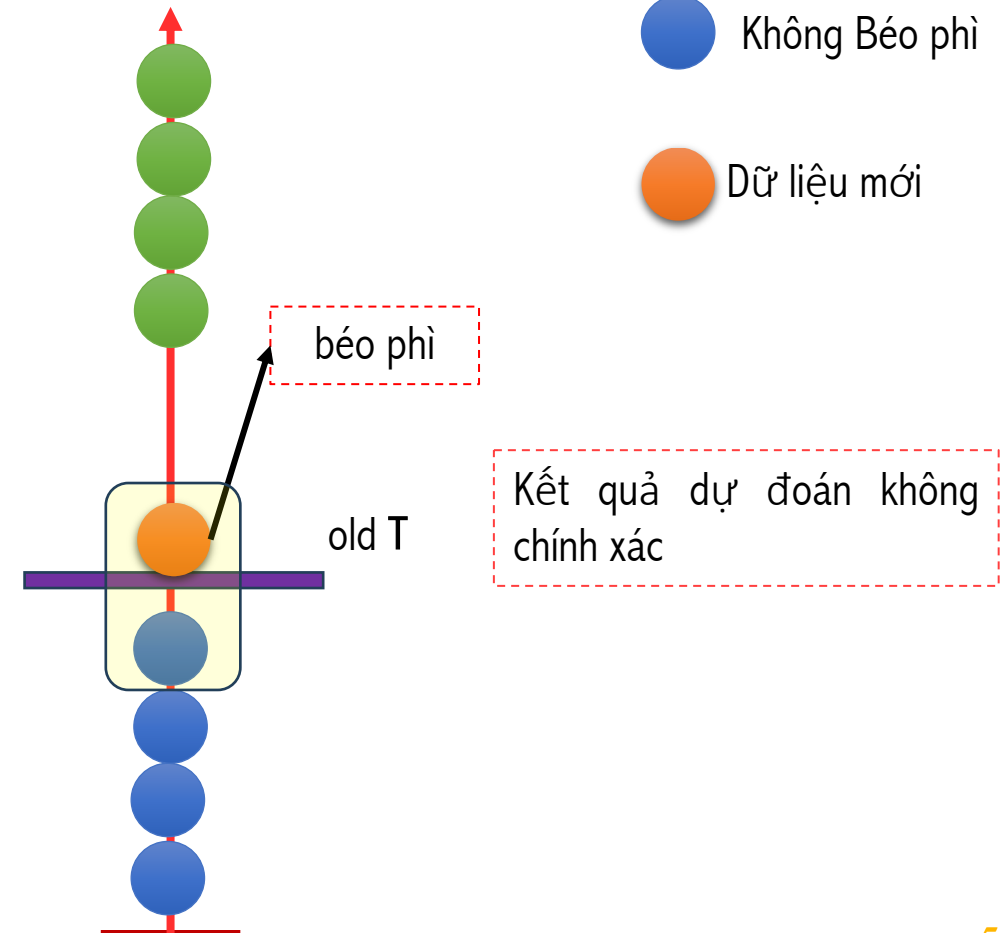
How to improve T ⁴

SVM Motivation

Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)



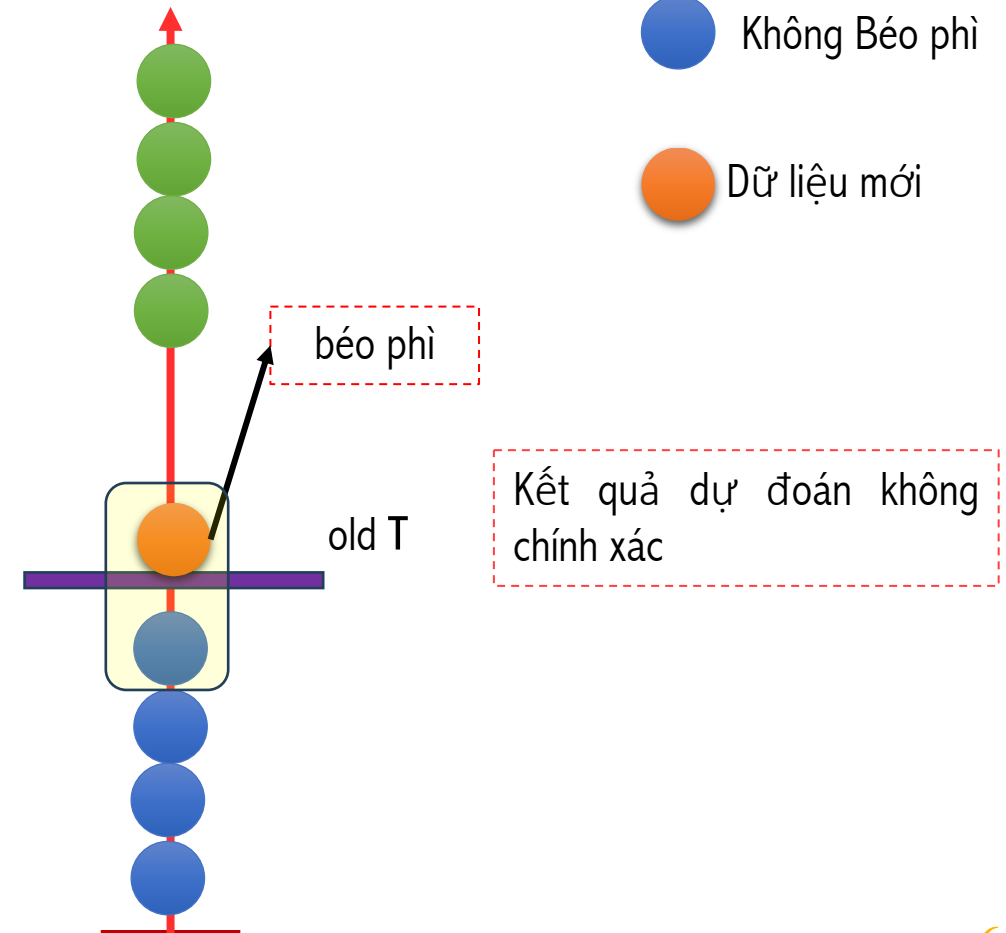
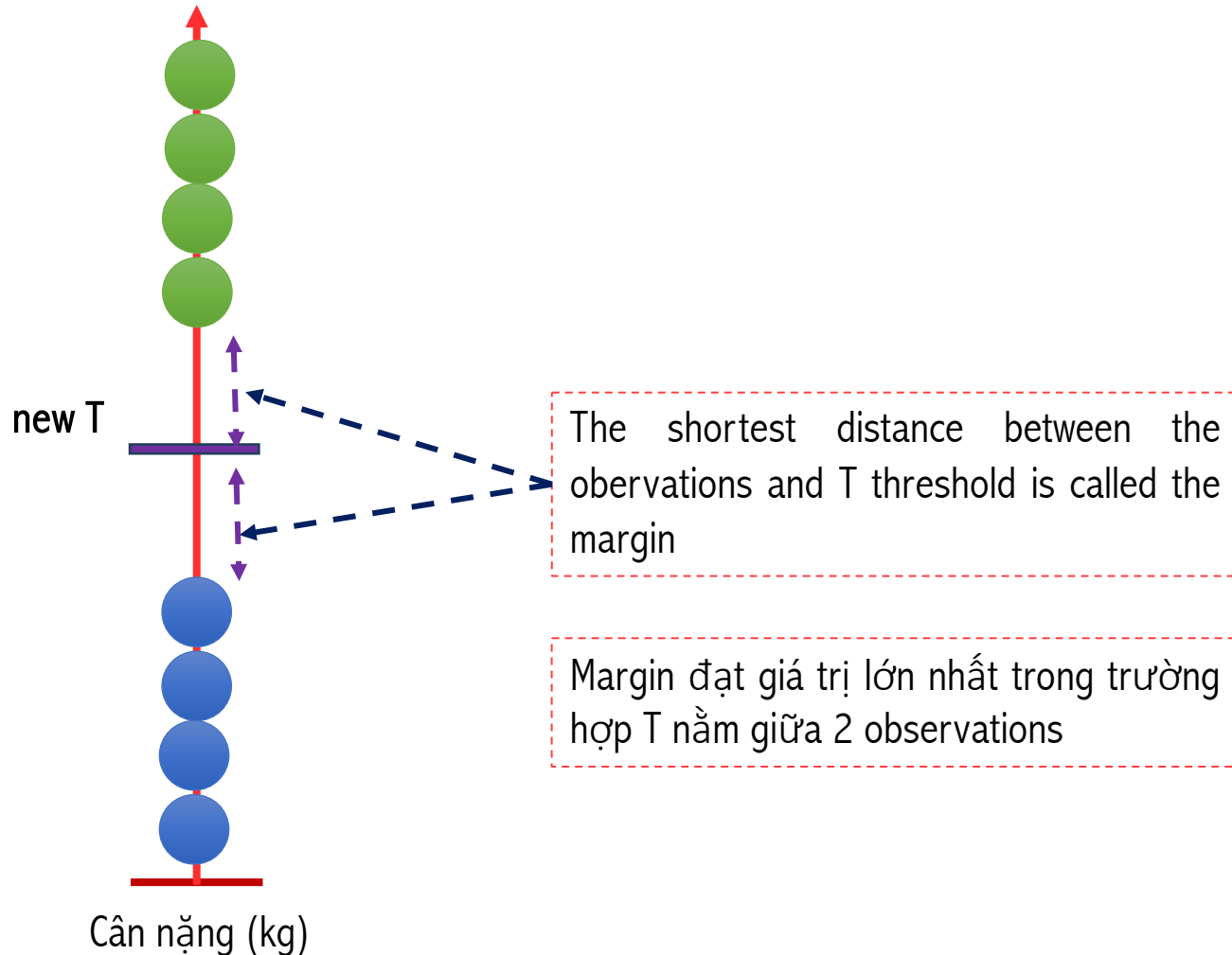
- Béo phì
- Không Béo phì
- Dữ liệu mới



SVM Motivation

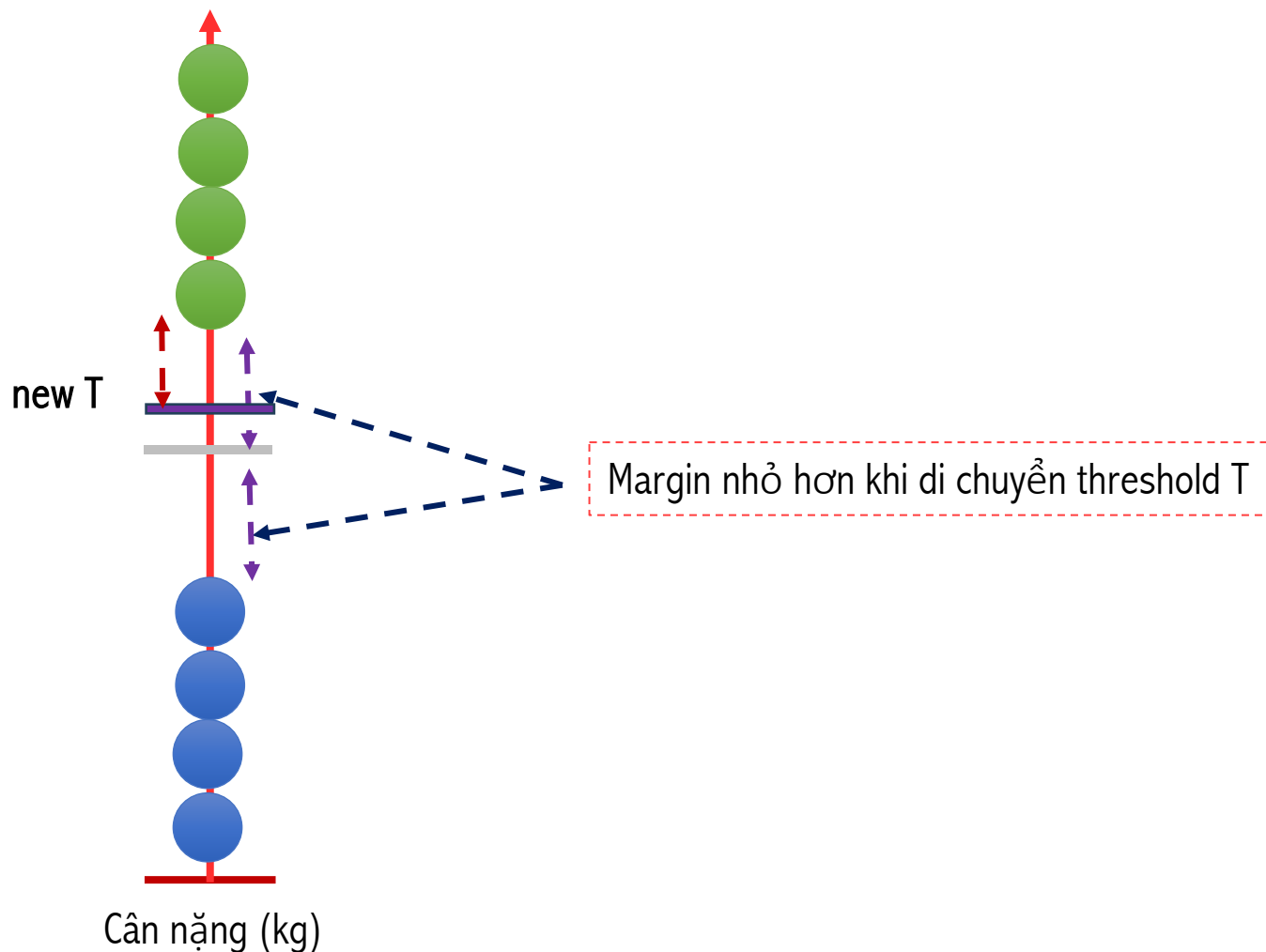
Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)

- Béo phì
- Không Béo phì
- Dữ liệu mới

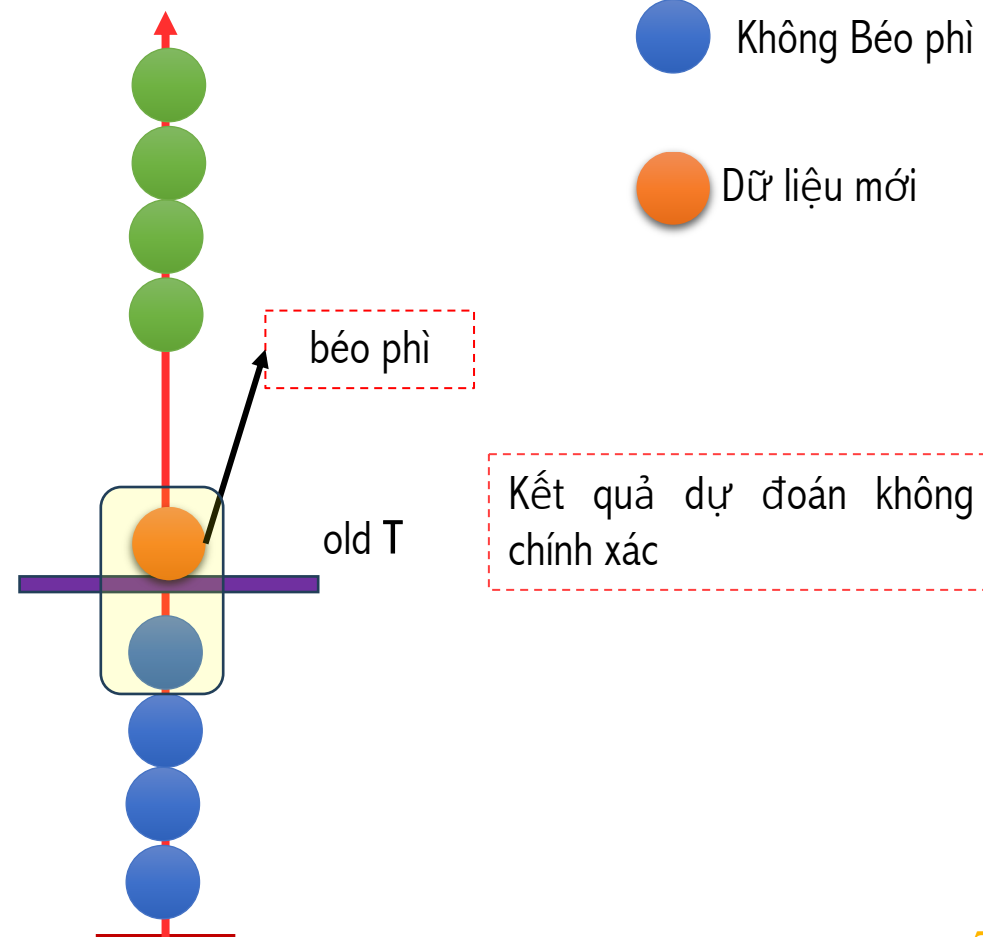


SVM Motivation

Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)

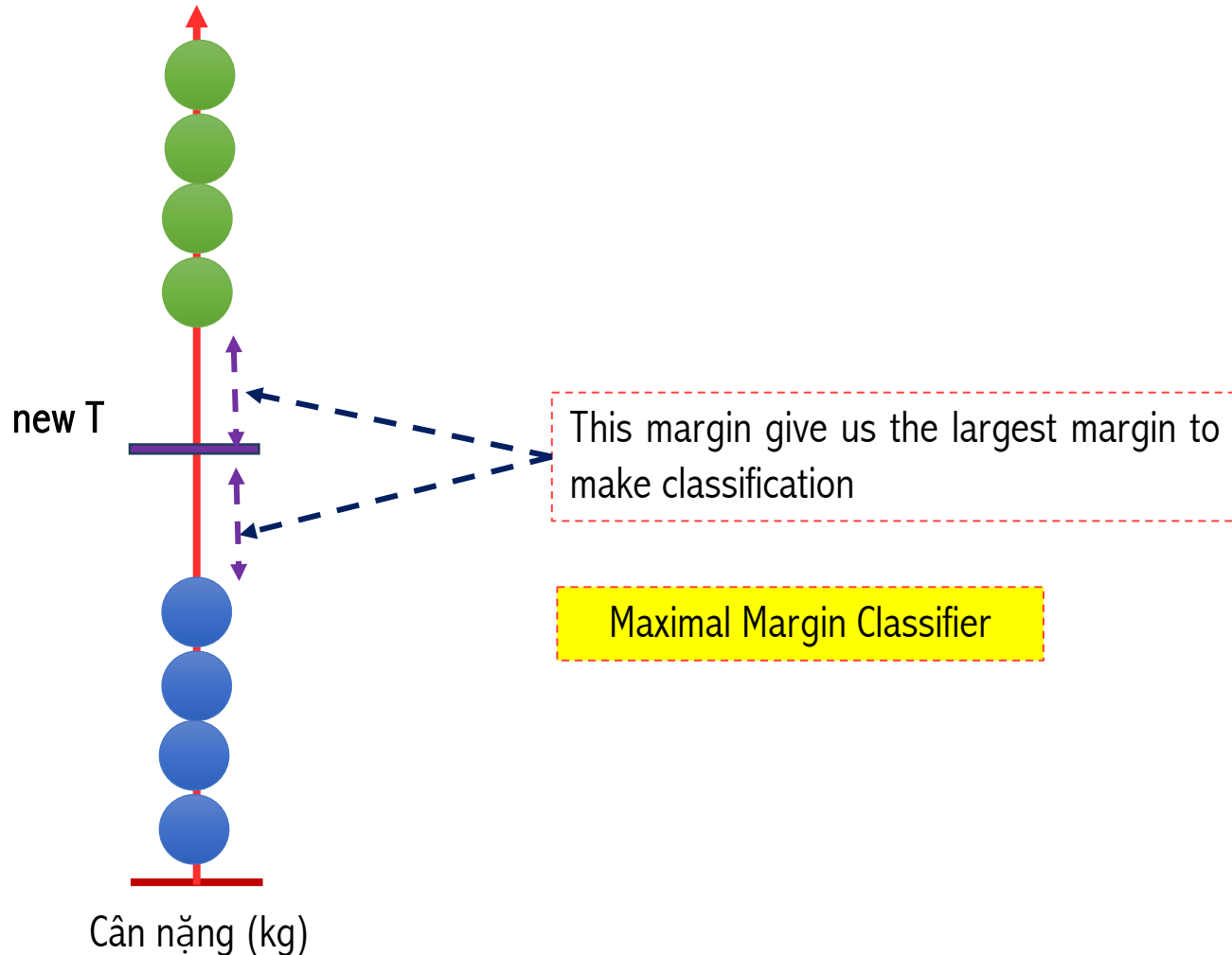


- Béo phì
- Không Béo phì
- Dữ liệu mới

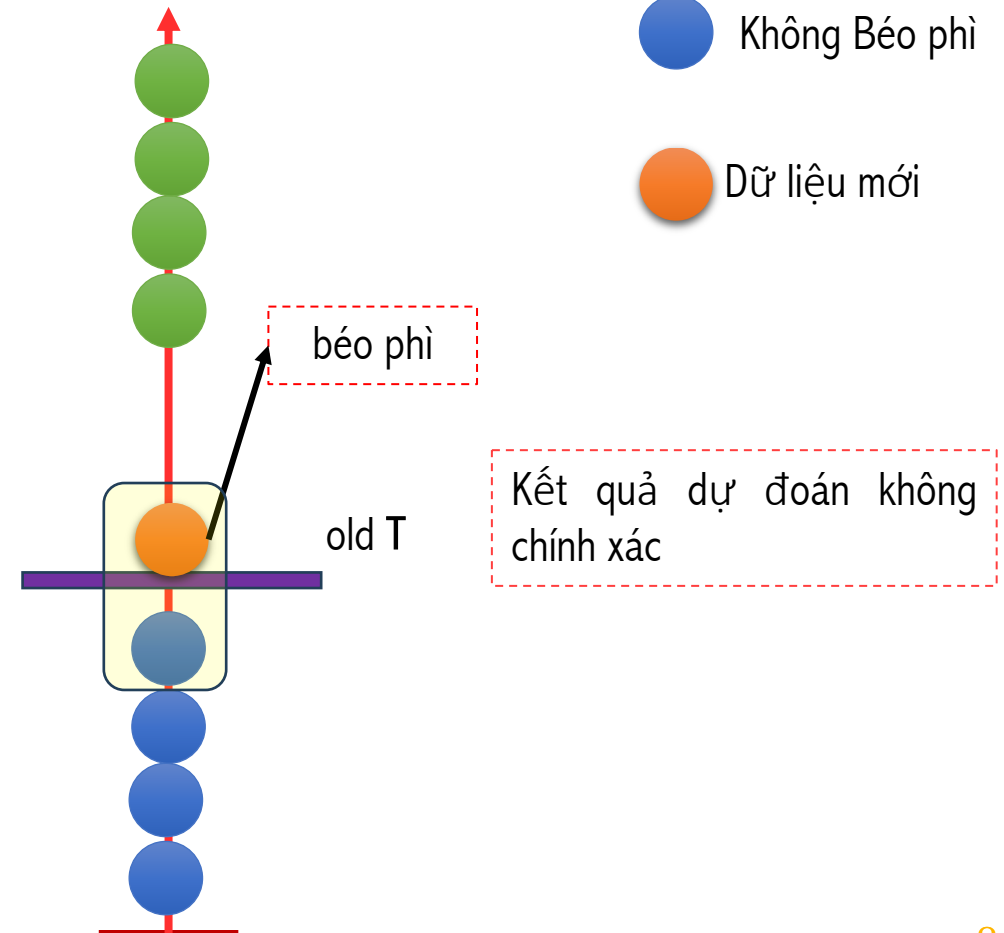


SVM Motivation

Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)



- Béo phì
- Không Béo phì
- Dữ liệu mới



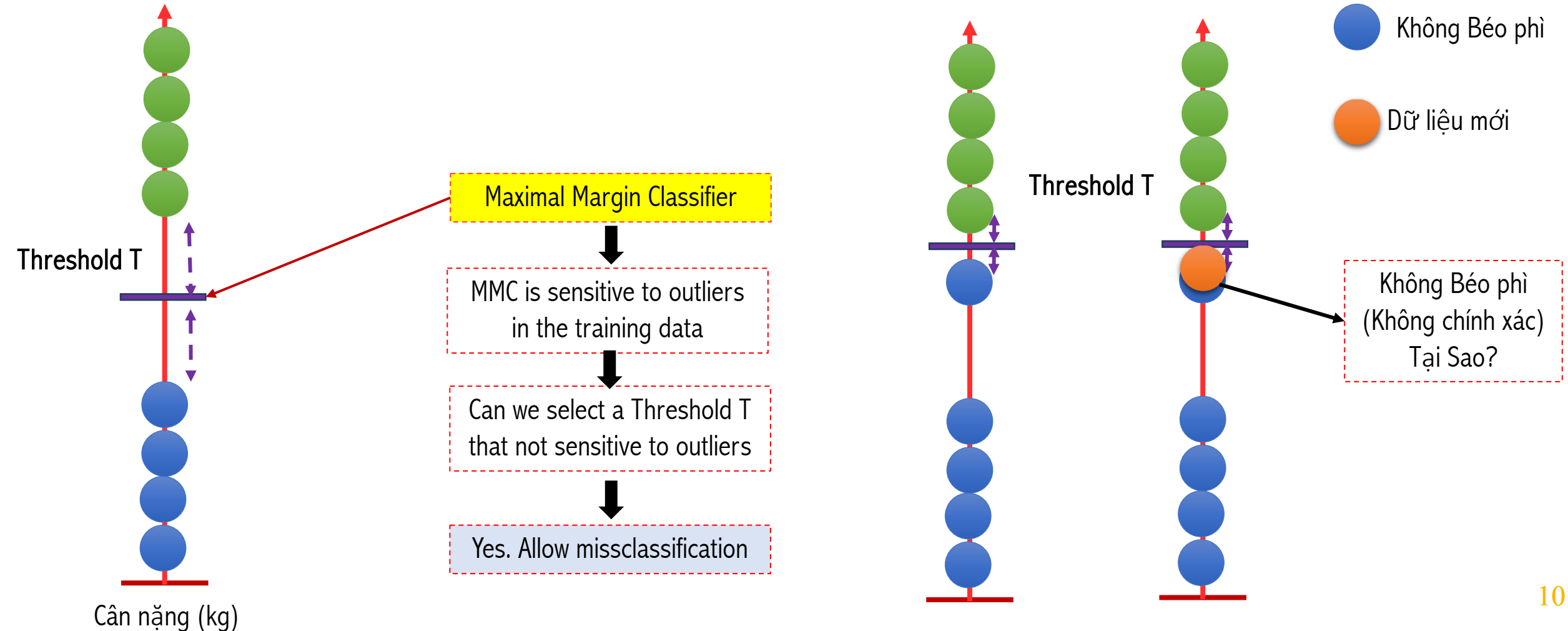
Outline

- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

SVM Motivation

Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)

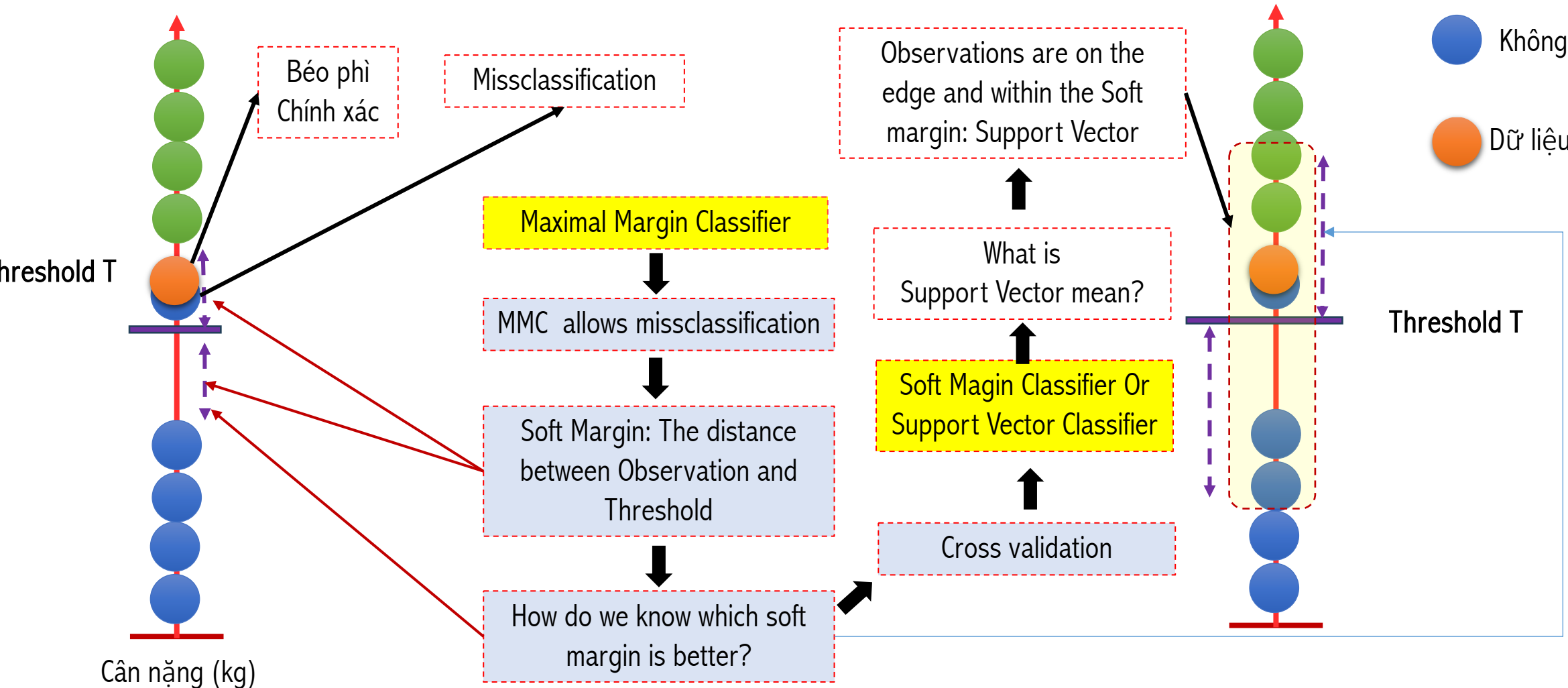
- Béo phì
- Không Béo phì
- Dữ liệu mới



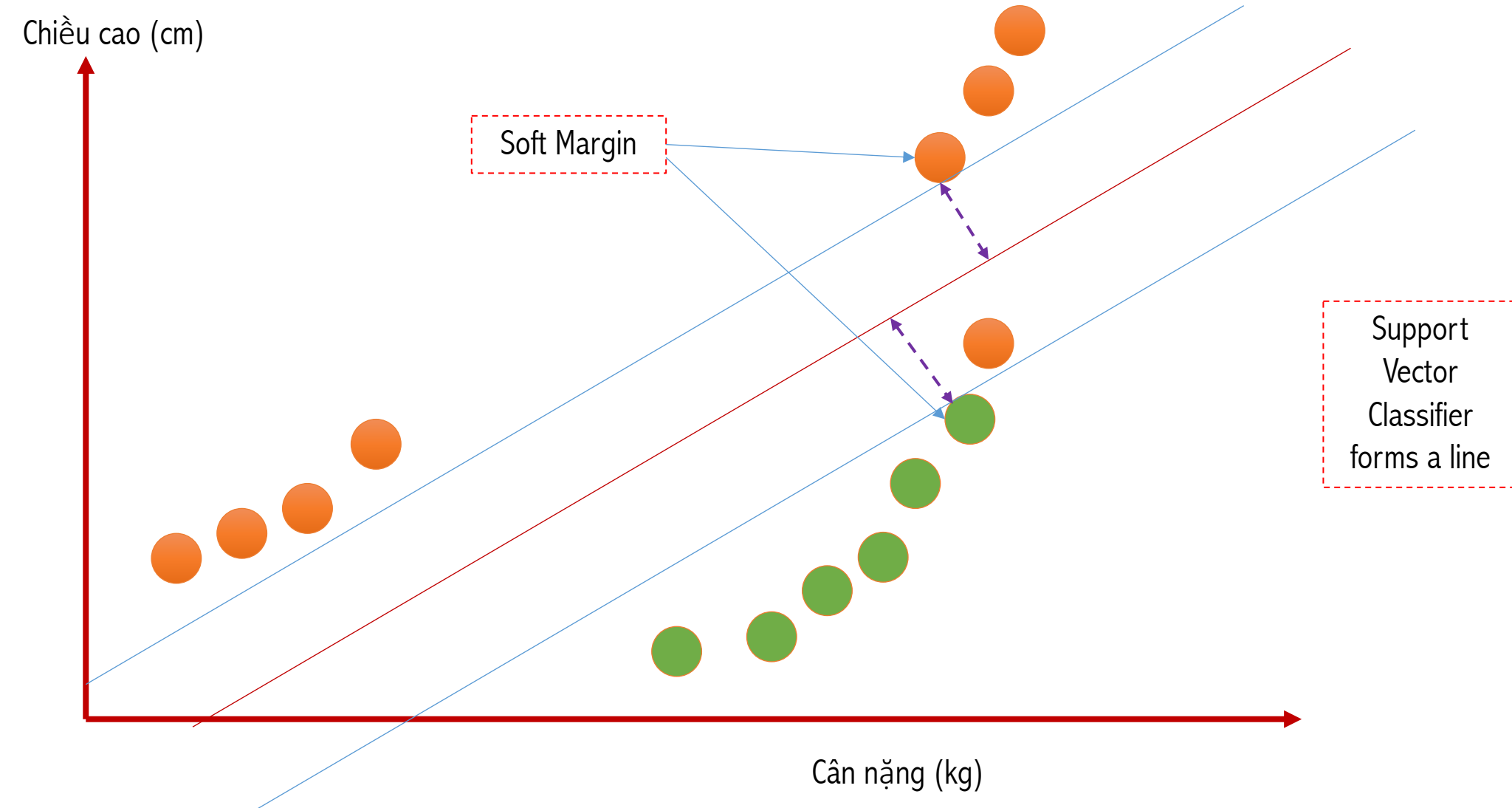
SVM Motivation

Phát triển chương trình dự đoán Lợn có khả năng bị béo phì hay không dựa vào cân nặng (kg)

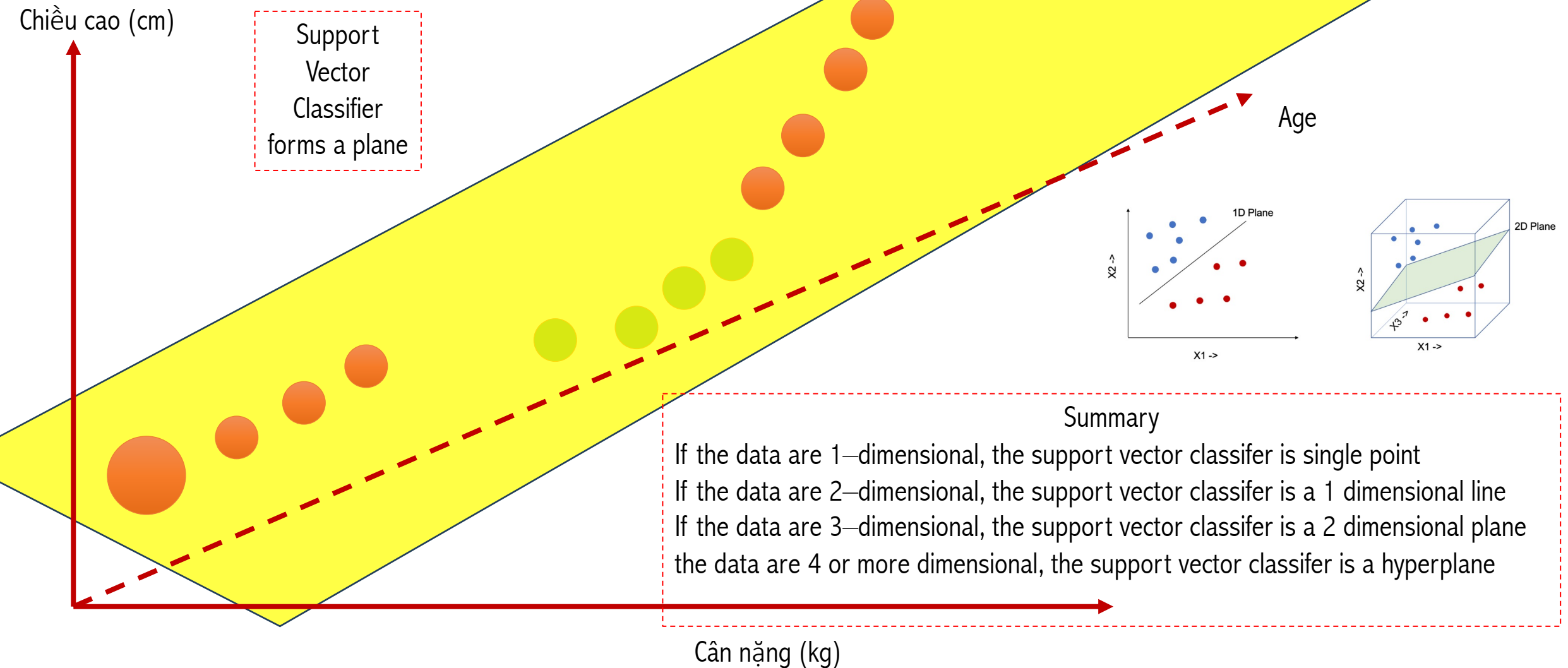
- Béo phì
- Không Béo phì
- Dữ liệu mới



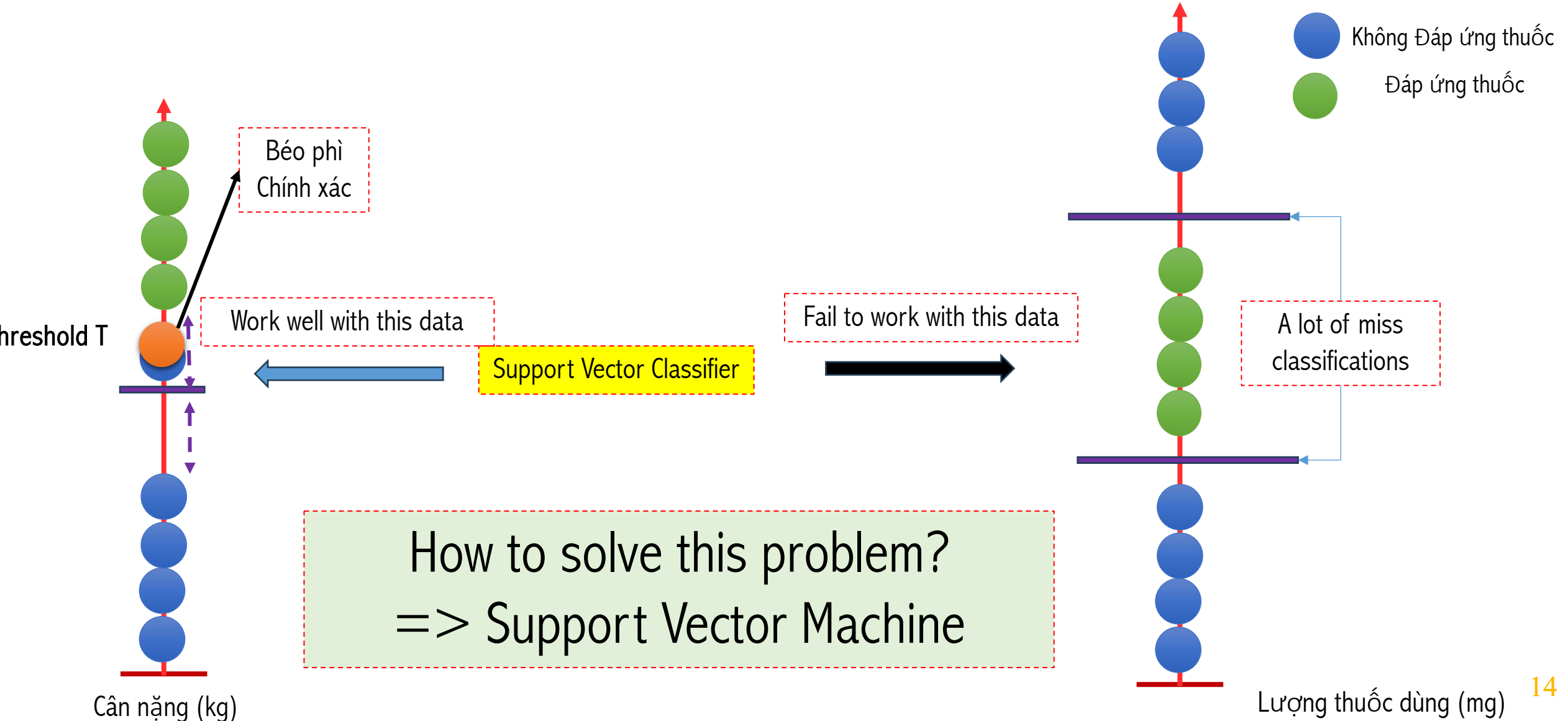
SVM Motivation: 2 Dimensional



SVM Motivation: 3 Dimensional



Support Vector Classifier: Limitation

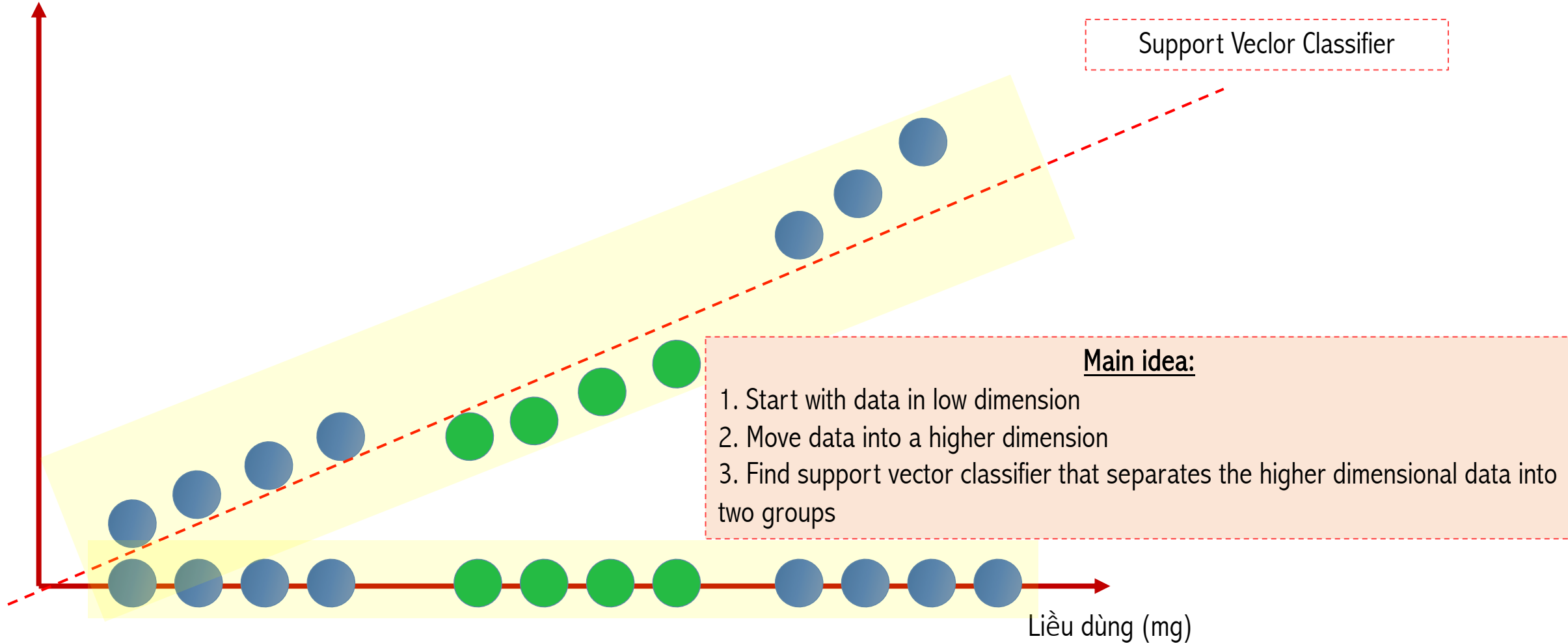


Outline

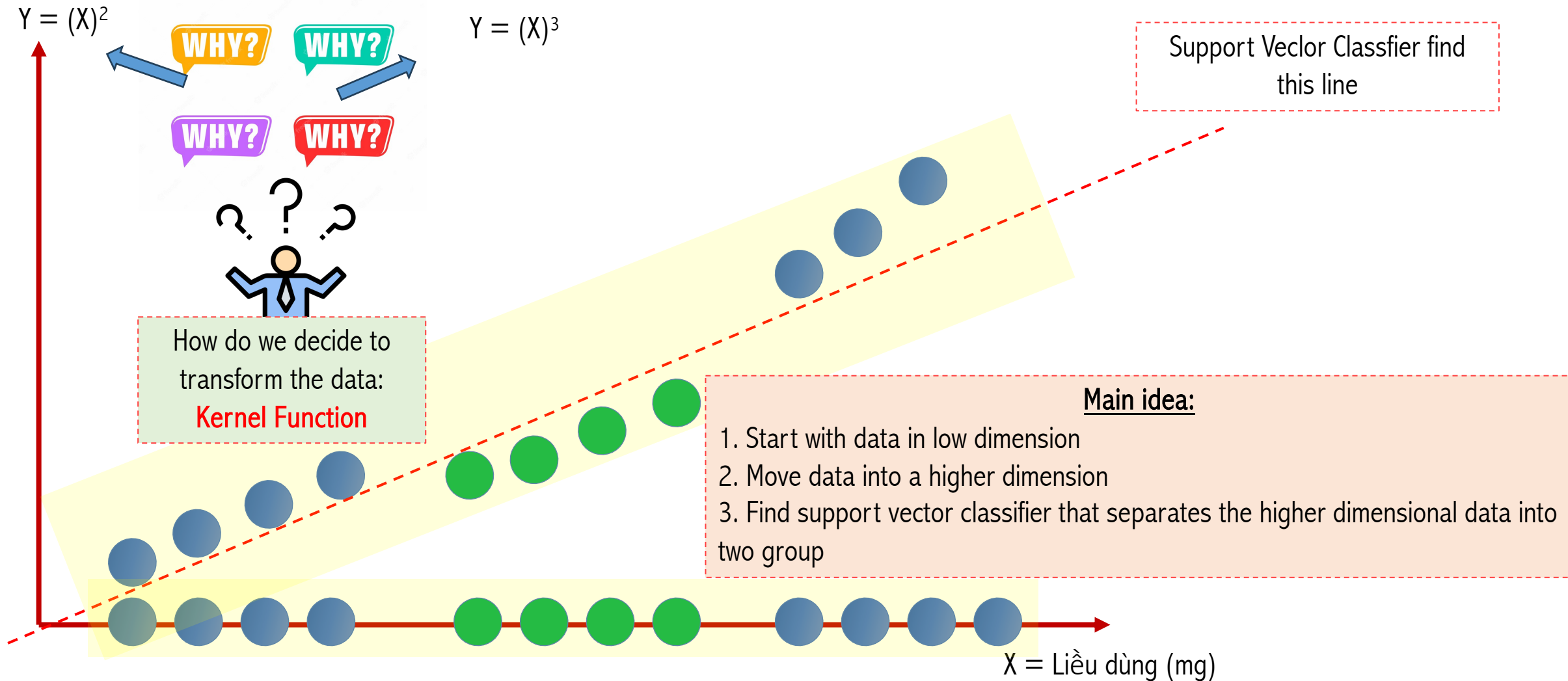
- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

Support Vector Machine: Main Idea

$Y = (\text{liều dùng})^2$



Support Vector Machine: Kernel Function



Kernel functions: None-linear functions that help us to transform data from lower dimension to higher dimension

Outline

- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

Polynomial Kernel

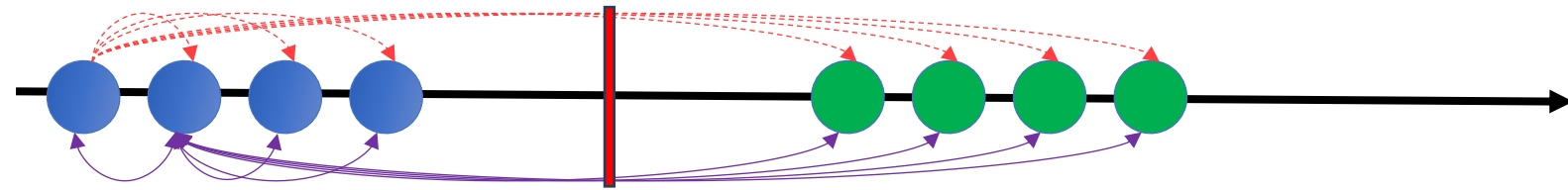
Polynomial Kernel: the degree of the polynomial d

$d = 1$. Compute the relationship between each pair of observations in 1-Dimensional to find SVC

$d = 2$. Compute 2-Dimensional relationship between each pair of observations find SVC

$d = n$. Compute n-Dimensional relationship between each pair of observations. Those relationship are used to find SVC

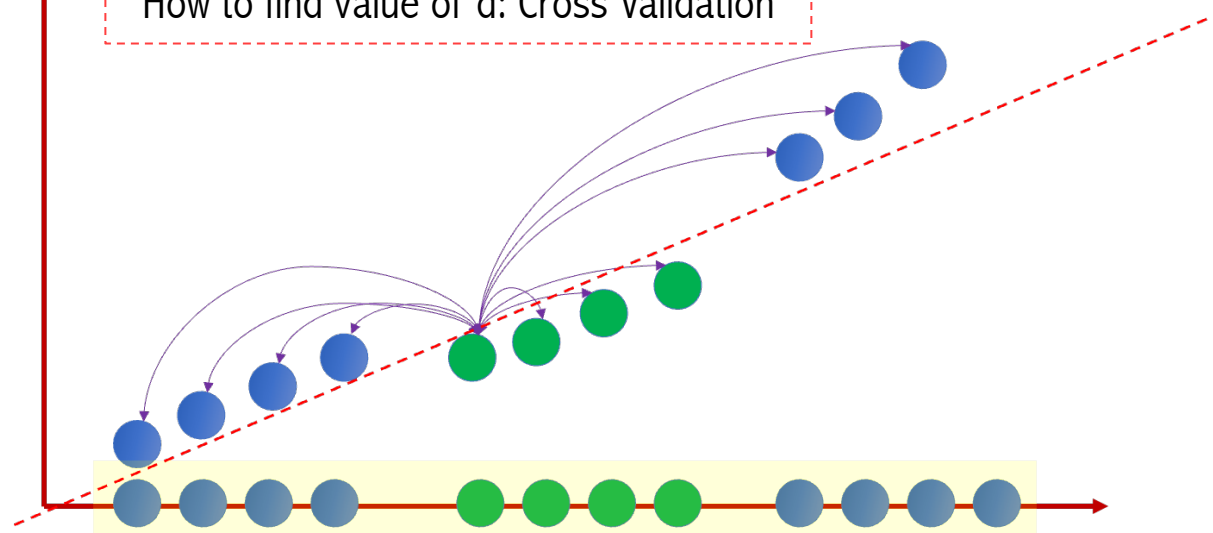
Other kernel: Radial Basic Function Kernel (RBF)



$$Y = (\text{liều dùng})^2$$

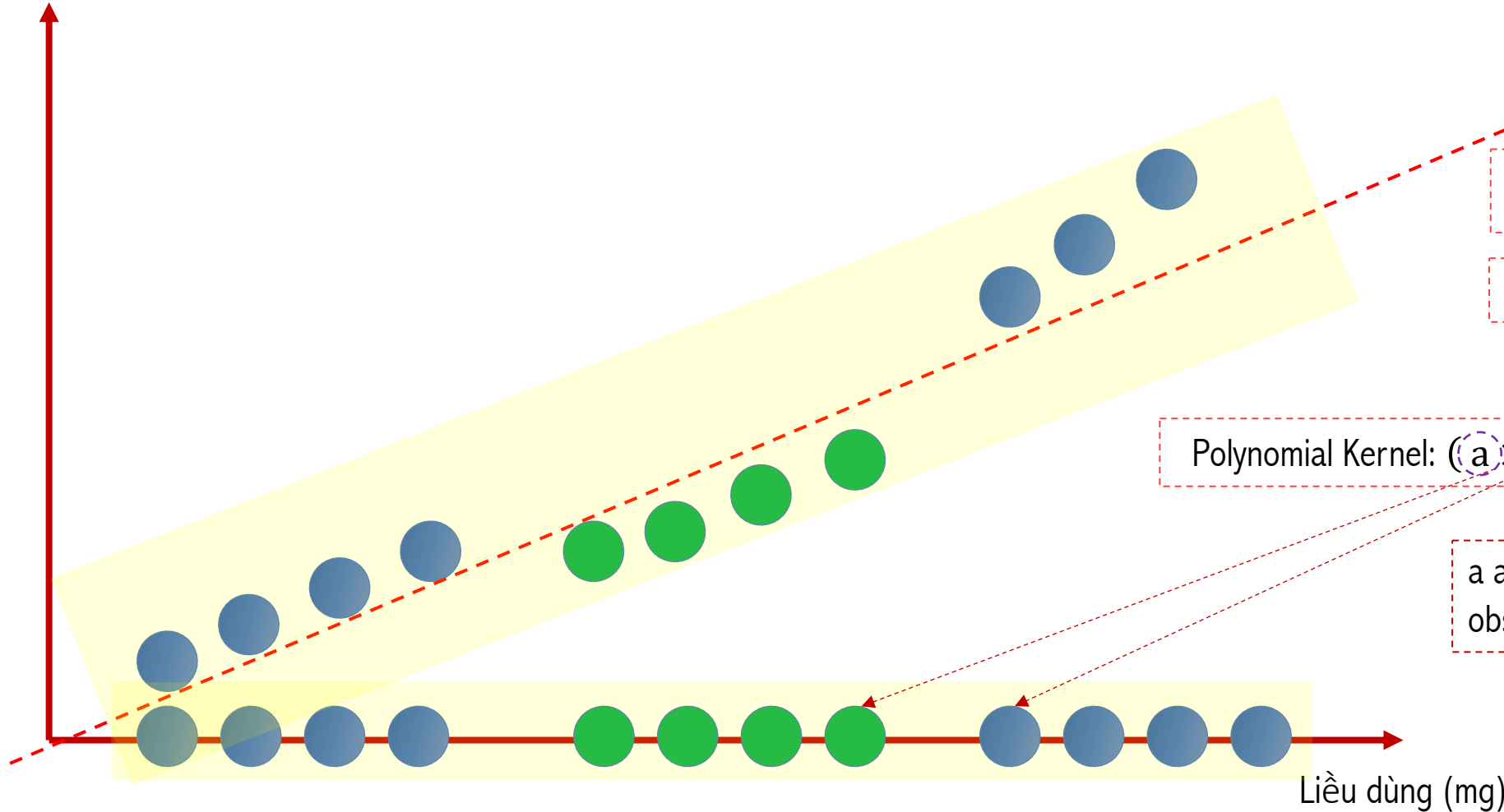
How to find value of d : Cross Validation

Support Vector Classifier



Polynomial Kernel

$Y = (\text{liều dùng})^2$



Example: $r = \frac{1}{2}$, $d = 2$

Coefficient of the polynomial kernel

Degree of polyno

Polynomial Kernel: $(a \times b + r)^d$

a and b refer to two different observations in the dataset

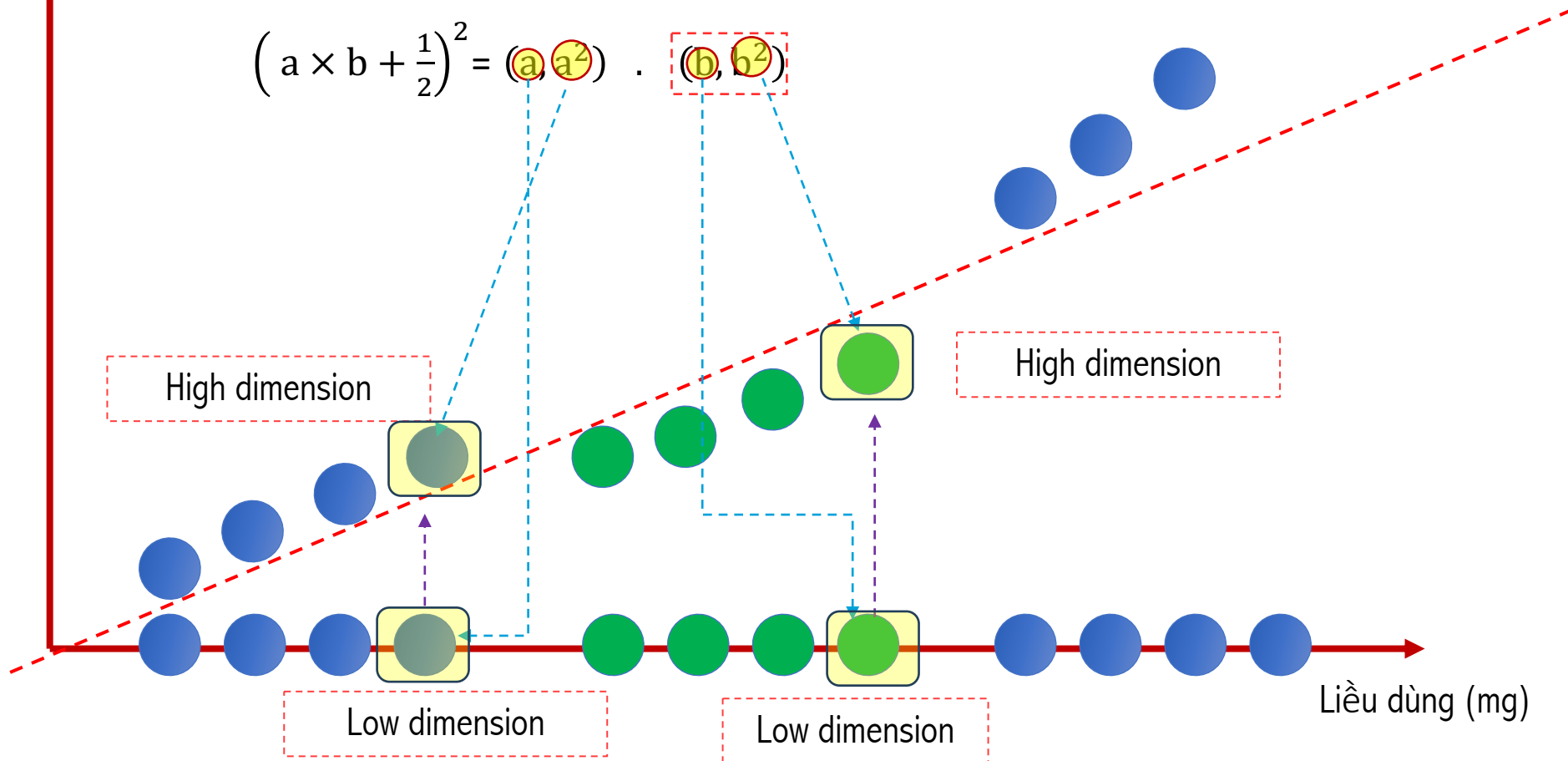
Polynomial Kernel

$Y = (\text{liều dùng})^2$

$$\left(a \times b + \frac{1}{2}\right)^2 = \left(a \times b + \frac{1}{2}\right) \left(a \times b + \frac{1}{2}\right) = a^2b^2 + ab + \frac{1}{4} = \left(a, a^2, \frac{1}{2}\right) \cdot \left(b, b^2, \frac{1}{2}\right)$$

$r = \frac{1}{2}, d = 2$

$$\left(a \times b + \frac{1}{2}\right)^2 = (a, a^2) \cdot (b, b^2)$$



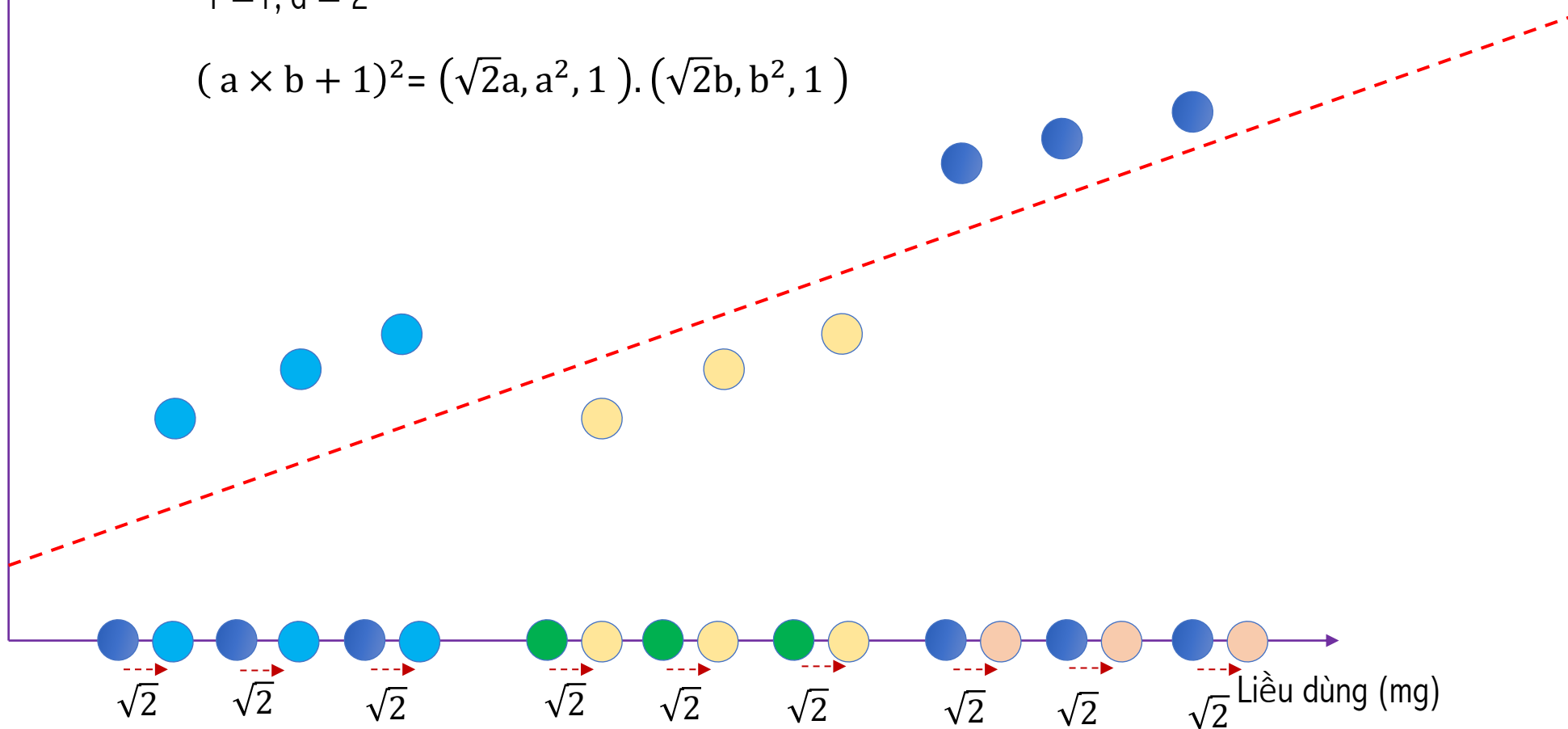
Polynomial Kernel

$Y = (\text{liều dùng})^2$

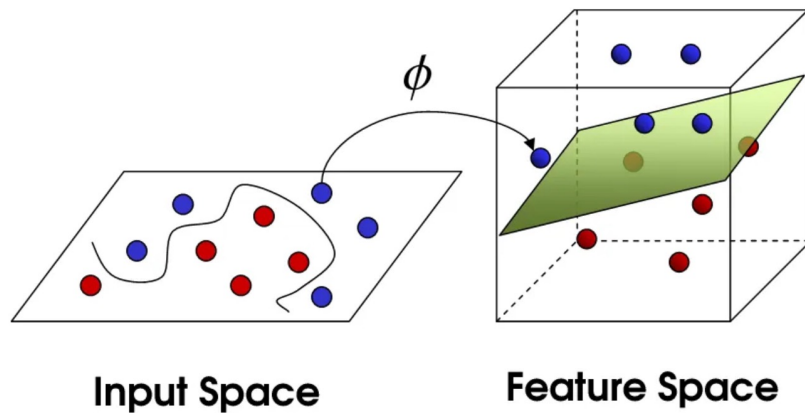
$$(a \times b + 1)^2 = (a \times b + 1) (a \times b + 1) = a^2b^2 + 2ab + 1 = (\sqrt{2}a, a^2, 1) \cdot (\sqrt{2}b, b^2, 1)$$

$$r=1, d=2$$

$$(a \times b + 1)^2 = (\sqrt{2}a, a^2, 1) \cdot (\sqrt{2}b, b^2, 1)$$



Kernel Trick



Low Dimension Data

Kernel Function

High Dimension Data

Test Data

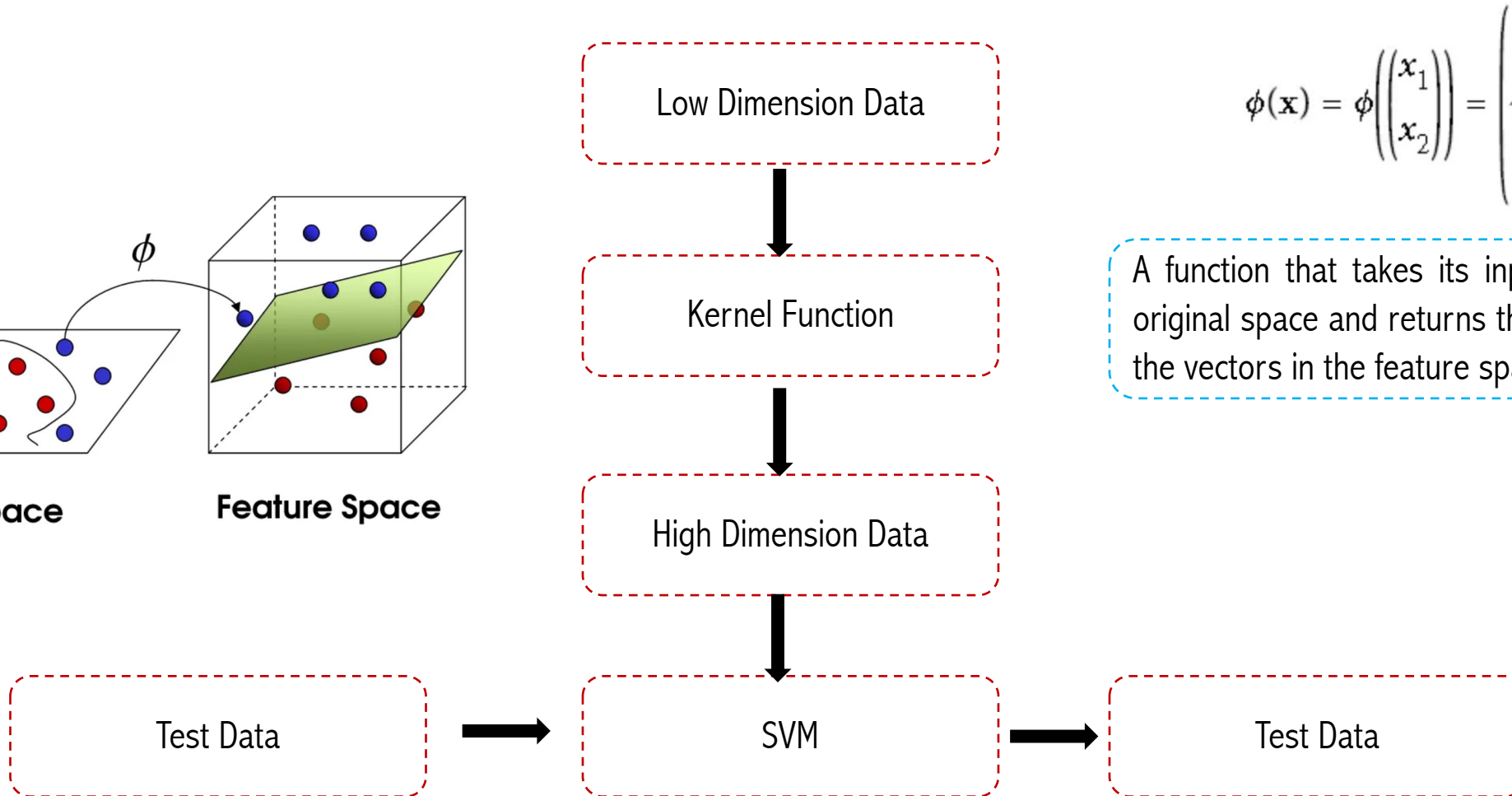
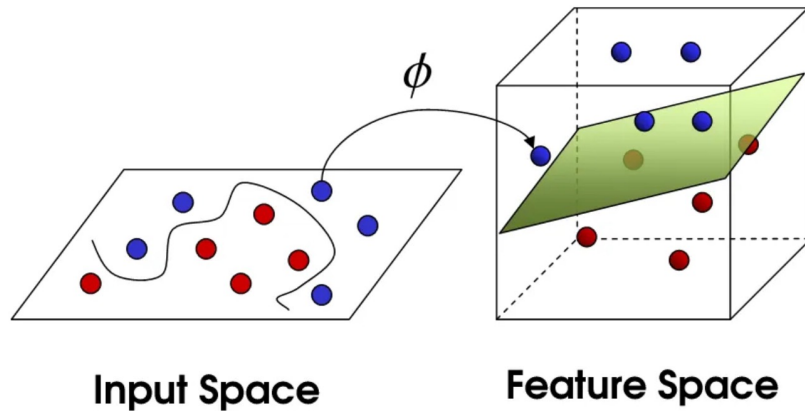
SVM

Test Data

$$\phi(\mathbf{x}) = \phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

A function that takes its input vector in the original space and returns the dot product of the vectors in the feature space

Kernel Trick

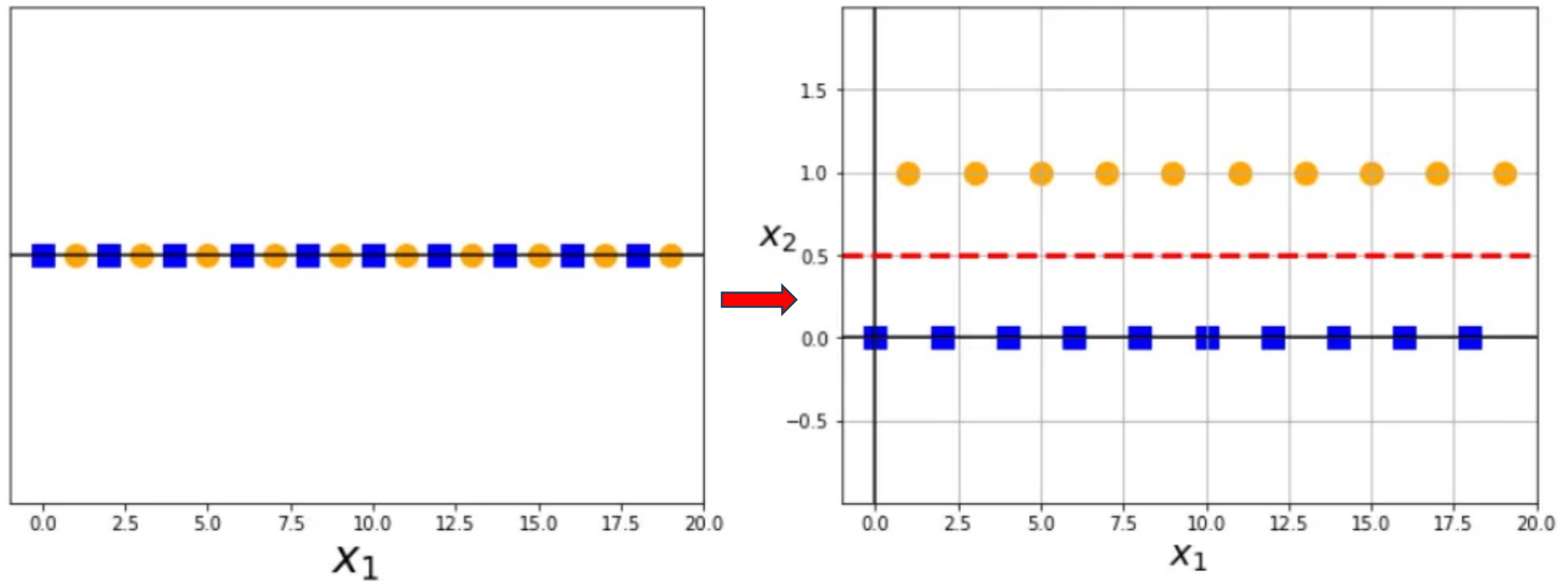


$$\phi(\mathbf{x}) = \phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

A function that takes its input vector in the original space and returns the dot product of the vectors in the feature space

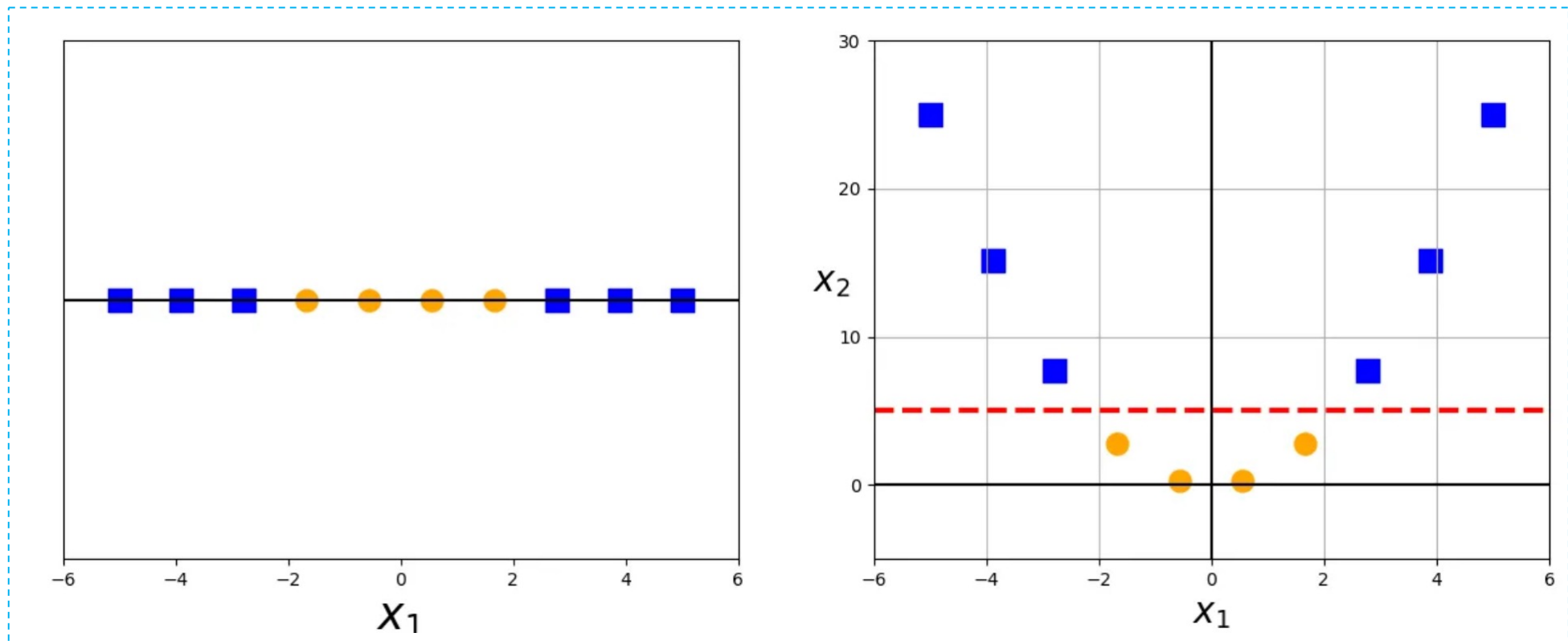
Kernel Trick

Apply the transformation $\phi(x) = x \bmod 2$



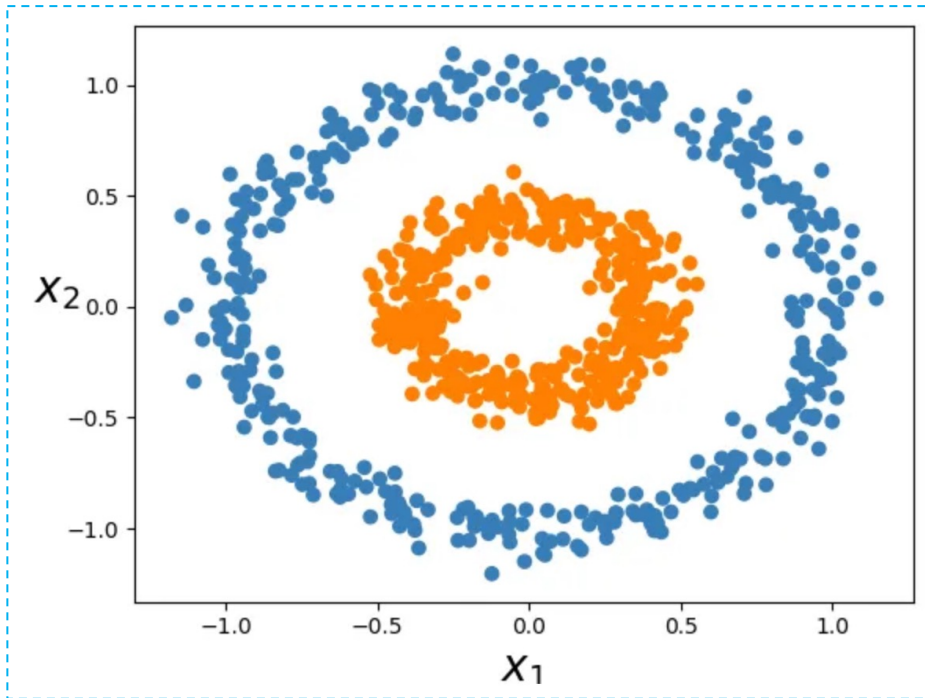
Kernel Trick

Apply the transformation $\phi(x) = x^2$




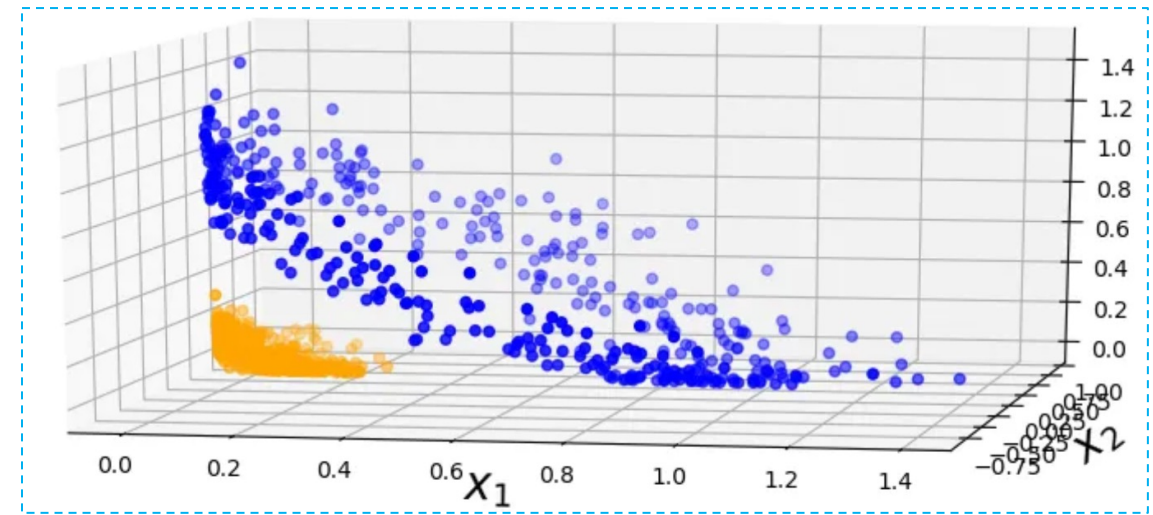
Kernel Trick

The **kernel trick** provides a solution to this problem. It allows us to operate in the original feature space without computing the coordinates of the data in a higher dimensional space.



Second-degree polynomial mapping

$$\phi(\mathbf{x}) = \phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$




We have seen how higher dimensional transformations can allow us to separate data in order to make classification predictions. It seems that in order to train a support vector classifier and optimize our objective function, we would have to perform operations with the higher dimensional vectors in the transformed feature space => **extremely high and impractical computational costs**

Kernel Trick

The “trick” is that kernel methods represent the data only through a set of pairwise similarity comparisons between the original data observations \mathbf{x} (with the original coordinates in the lower dimensional space), instead of explicitly applying the transformations $\phi(\mathbf{x})$ and representing the data by these transformed coordinates in the higher dimensional feature space.

Kernel Function:

More formally, if we have data $\mathbf{x}, \mathbf{z} \in X$ and a map $\phi: X \rightarrow \mathbb{R}^N$ then

$$k(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$$

is a kernel function

The ultimate benefit of the kernel trick is that the objective function we are optimizing to fit the higher dimensional decision boundary only includes the dot product of the transformed feature vectors. Therefore, we can just substitute these **dot product terms** with the kernel function, and we don't even use $\phi(\mathbf{x})$.

$$\begin{aligned} \phi(\mathbf{a})^T \cdot \phi(\mathbf{b}) &= \begin{pmatrix} a_1^2 \\ \sqrt{2} a_1 a_2 \\ a_2^2 \end{pmatrix}^T \cdot \begin{pmatrix} b_1^2 \\ \sqrt{2} b_1 b_2 \\ b_2^2 \end{pmatrix} = a_1^2 b_1^2 + 2a_1 b_1 a_2 b_2 + a_2^2 b_2^2 \\ \phi(\mathbf{x}) &= \phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix} \\ &= (a_1 b_1 + a_2 b_2)^2 = \left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}^T \cdot \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right)^2 = (\mathbf{a}^T \cdot \mathbf{b})^2 \end{aligned}$$

The kernel function here is the polynomial kernel

$$k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \cdot \mathbf{b})^2$$

Kernel Trick

$Y = (\text{liều dùng})^2$

$$\left(a \times b + \frac{1}{2}\right)^2 = \left(a \times b + \frac{1}{2}\right) \left(a \times b + \frac{1}{2}\right) = a^2b^2 + ab + \frac{1}{4} = \left(a, a^2, \frac{1}{2}\right) \cdot \left(b, b^2, \frac{1}{2}\right)$$

$$r = \frac{1}{2}, d = 2$$

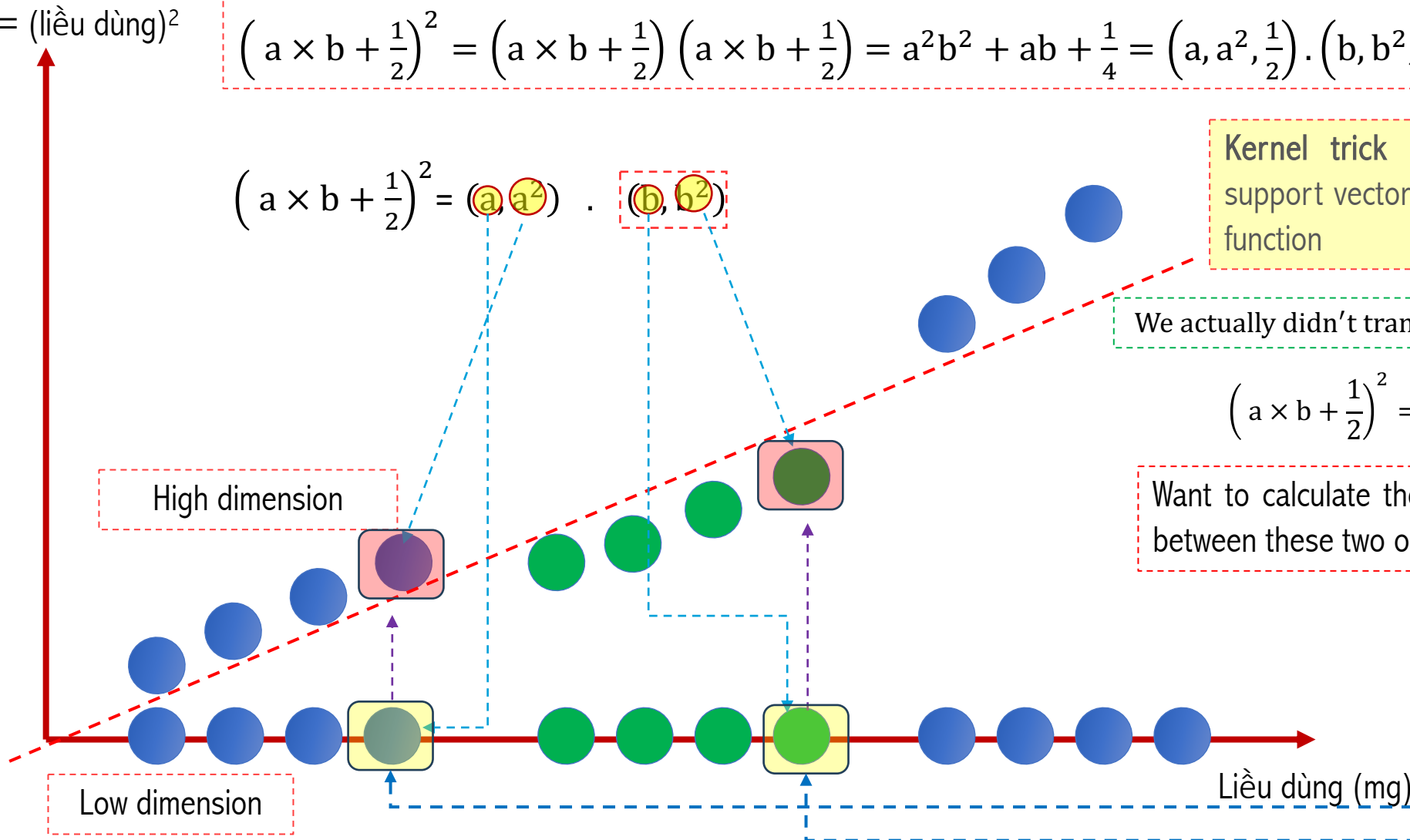
$$\left(a \times b + \frac{1}{2}\right)^2 = (a, a^2) \cdot (b, b^2)$$

Kernel trick is to convert dot product of support vectors to the dot product of mapping function

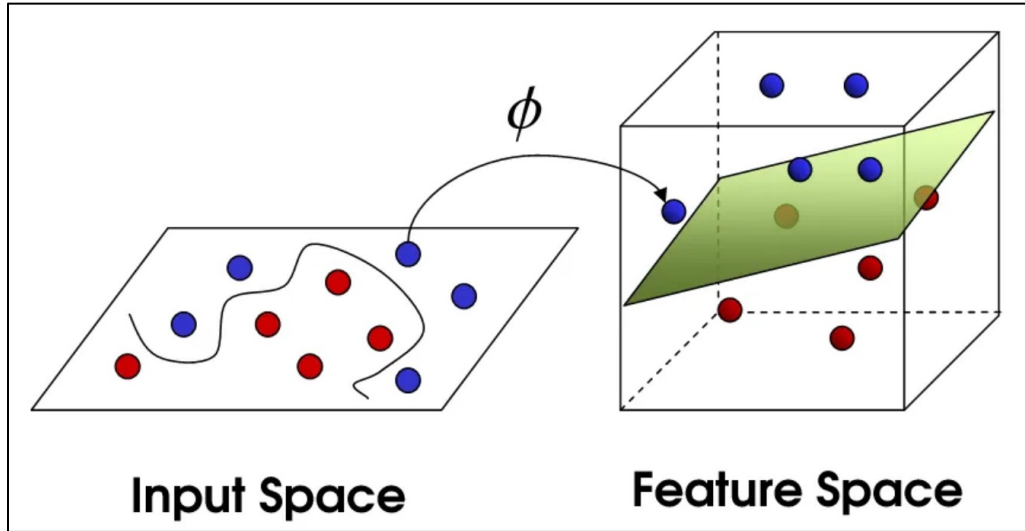
We actually didn't transform the data to 2 - Dimensions.

$$\left(a \times b + \frac{1}{2}\right)^2 = \left(5 \times 10 + \frac{1}{2}\right)^2 = 110.25$$

Want to calculate the high dimensional relationship between these two observations (samples)



Kernel Trick



Input space $x = (x_1, x_2)$

Feature space $\phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$

$$\begin{aligned}\phi(x) \cdot \phi(z) &= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1z_2) \\ &= 1 + 2x_1z_1 + 2x_2z_2 + x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (1 + x_1z_1 + x_2z_2)^2 \\ &= (1 + x \cdot z)^2\end{aligned}$$

Did you notice the magic of mathematics ? We can represent the dot product $\phi(x) \cdot \phi(z)$ in feature space just by using a simple formula $(1 + x \cdot z)^2$ in input space.

So we do not have to perform any complex transformations or store the feature space in memory, if the dot product of feature space can be represented using dot product of input space.

Kernel Trick

If the feature space is **abstract vector space** then the Kernel is represented using $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ and when the feature space is **vector space** then the **transpose** operator can be used. So don't be confused seeing these two different representations, they are essentially the same.

$$\phi(x_i)^T \phi(x_j) \equiv \langle \phi(x_i), \phi(x_j) \rangle$$

Any Machine Learning algorithm can use Kernel Method, however it's mainly used in **SVM** and **Clustering**.

Radial Kernel

Find support vector classifier in infinite dimensions

Radial kernel: $e^{-\gamma(a-b)^2}$

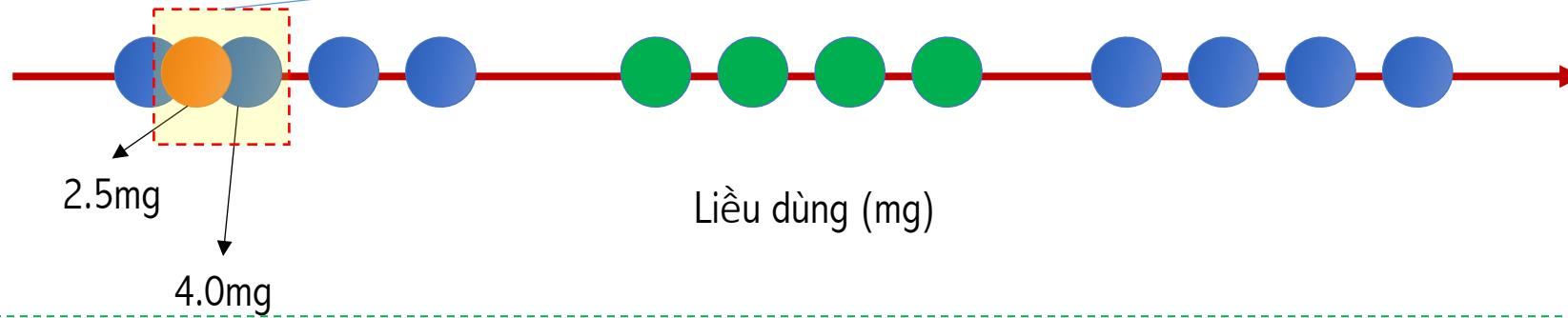
$$\gamma = 1$$

$$e^{-(2.5-4)^2} = 0.11$$

$$\gamma = 2$$

$$e^{-(2.5-4)^2} = 0.01$$

γ scales the amount of influence two points have each other



The nearest neighbors have a lot of influence on how we classify the new observation.

Radial kernel determines how much influence each observation in the Training Dataset has on classifying new observation

Outline

- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

Radial Kernel

Find support vector classifier in infinite dimensions

Radial kernel: $e^{-\gamma(a-b)^2}$

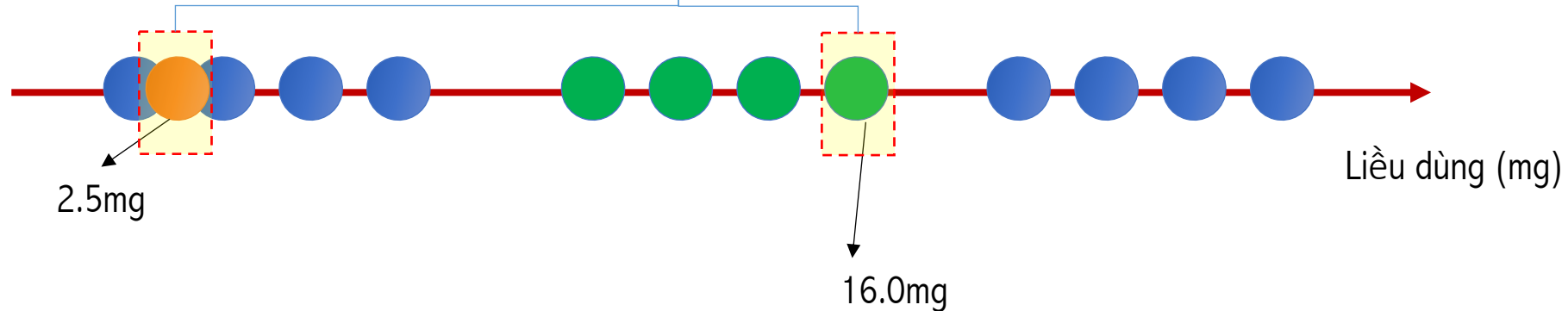
$$\gamma = 1$$

$$e^{-(2.5-16)^2} \approx 0$$

$$\gamma = 2$$

$$e^{-(2.5-16)^2} \approx 0$$

γ scales the amount of influence two points have each other



The further two observations are from each other, the less influence they have on each other

Radial Kernel: Intuition

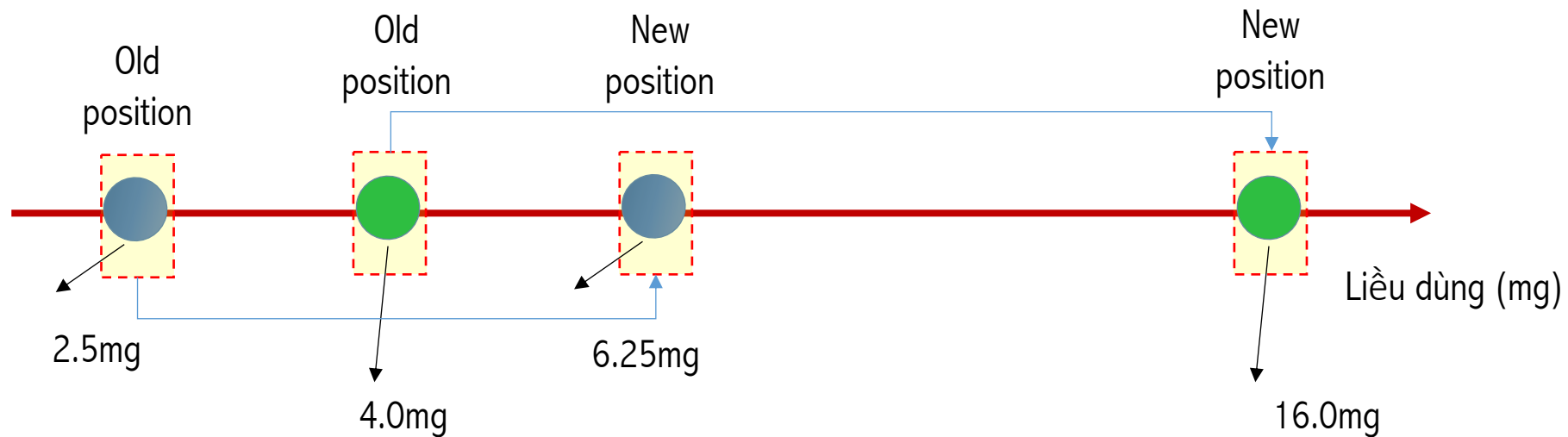
Polynomial Kernel: $(a \times b + r)^d$ with $r = 0$

$$(a \times b + r)^d = (a \times b)^d = (a^d) (b^d)$$

$$d = 2 \Rightarrow (a \times b + r)^2 = (a^2) (b^2)$$

$$d = 2 \Rightarrow (a \times b + r)^2 = (2.5^2) (4^2)$$

This dot product only has one coordinate. The new coordinate is square of the original measurement on the original axis



The original data are shift with $r = 0$ and $d = 2$

Radial Kernel: Intuition

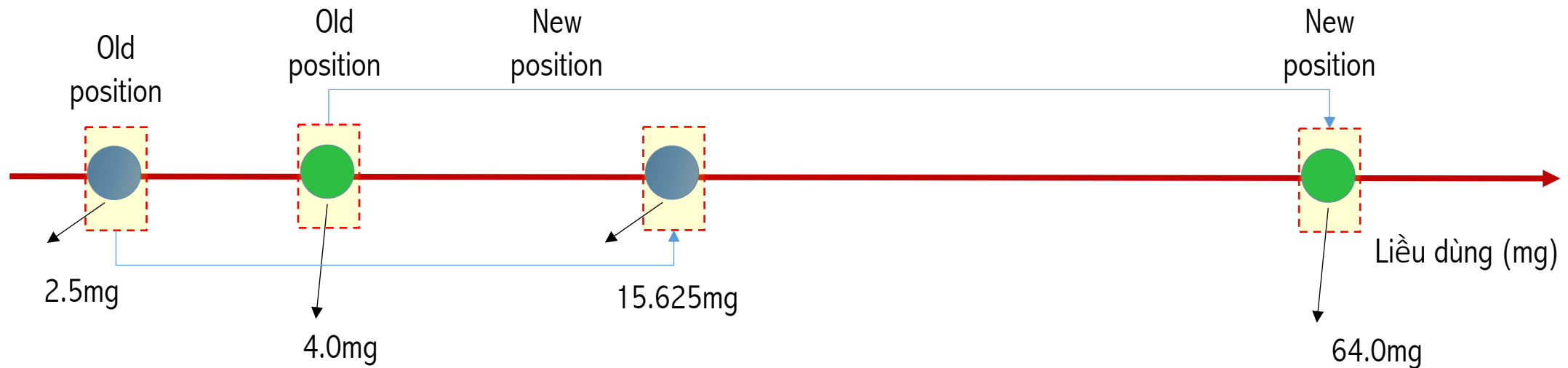
Polynomial Kernel: $(a \times b + r)^d$ with $r = 0$

$$(a \times b + r)^d = (a \times b)^d = (a^d) (b^d)$$

$$d = 3 \Rightarrow (a \times b + r)^3 = (a^3) (b^3)$$

$$d = 3 \Rightarrow (a \times b + r)^3 = (2.5^3) (4^3)$$

This dot product only has one coordinate. The new coordinate is square of the original measurement on the original axis



The original data are shift further with $r = 0$ and $d = 3$

Radial Kernel: Intuition

Polynomial Kernel: $(a \times b + r)^d$ with $r = 0$

$$(a \times b + r)^d = (a \times b)^d = (a^d) (b^d)$$

$$d = 3 \Rightarrow (a \times b + r)^3 = (a^1) (b^1)$$

$$d = 3 \Rightarrow (a \times b + r)^3 = (2.5^1) (4^1)$$

This dot product only has one coordinate. The new coordinate is square of the original measurement on the original axis



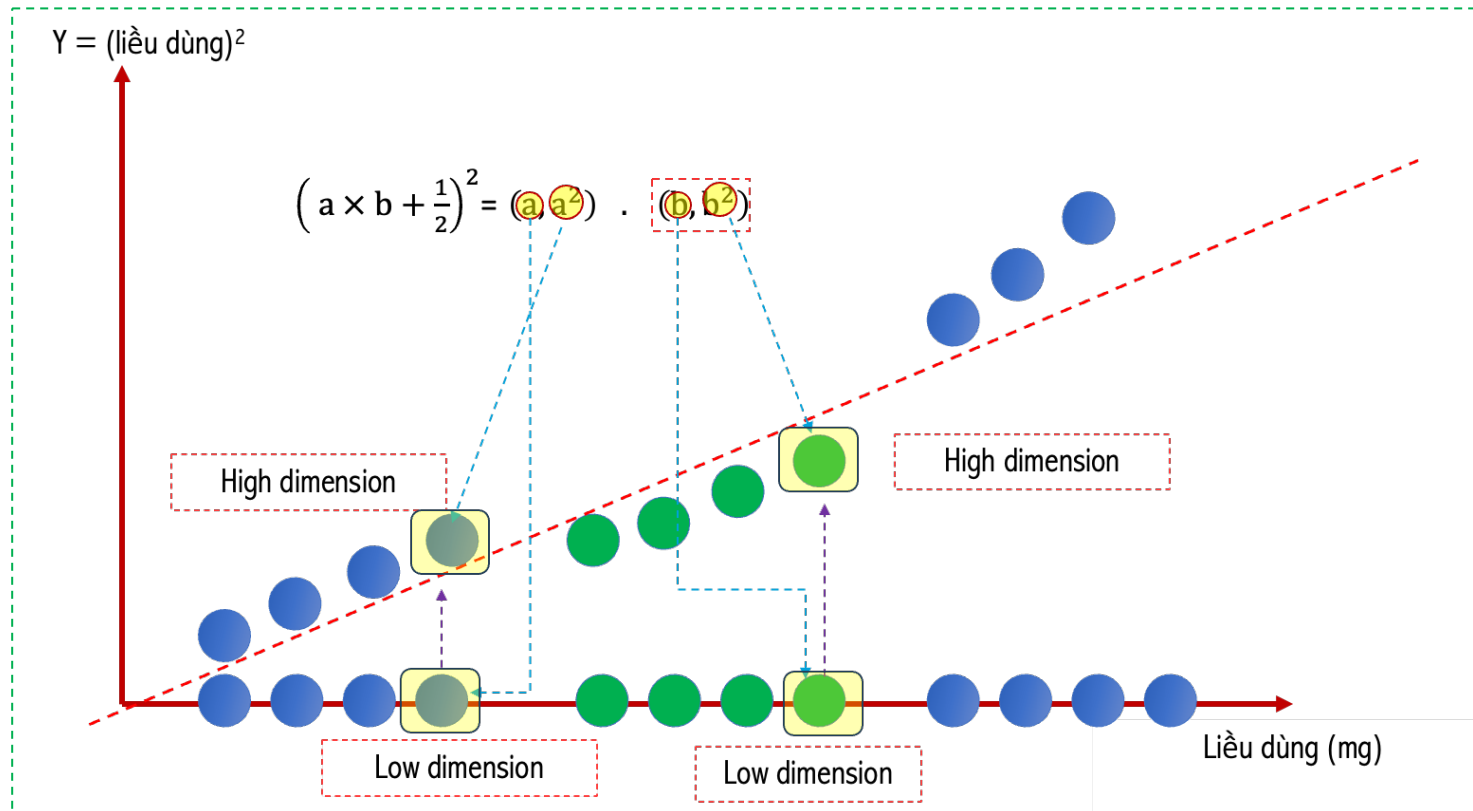
The original data stays in its original position with $r = 0$ and $d = 1$.
The data stays on the same 1-dimensional line regardless the value of d

Radial Kernel: Intuition

Polynomial Kernel: $(a \times b)^1 = (a^1) (b^1)$ with $r = 0$ and $d = 1$

Polynomial Kernel: $(a \times b)^2 = (a^2) (b^2)$ with $r = 0$ and $d = 2$

$$(a^1) (b^1) + (a^2) (b^2) = (a, a^2) \cdot (b, b^2)$$



We do not actually do the transformation, we just solve for Dot Product the get high dimensional relationship!

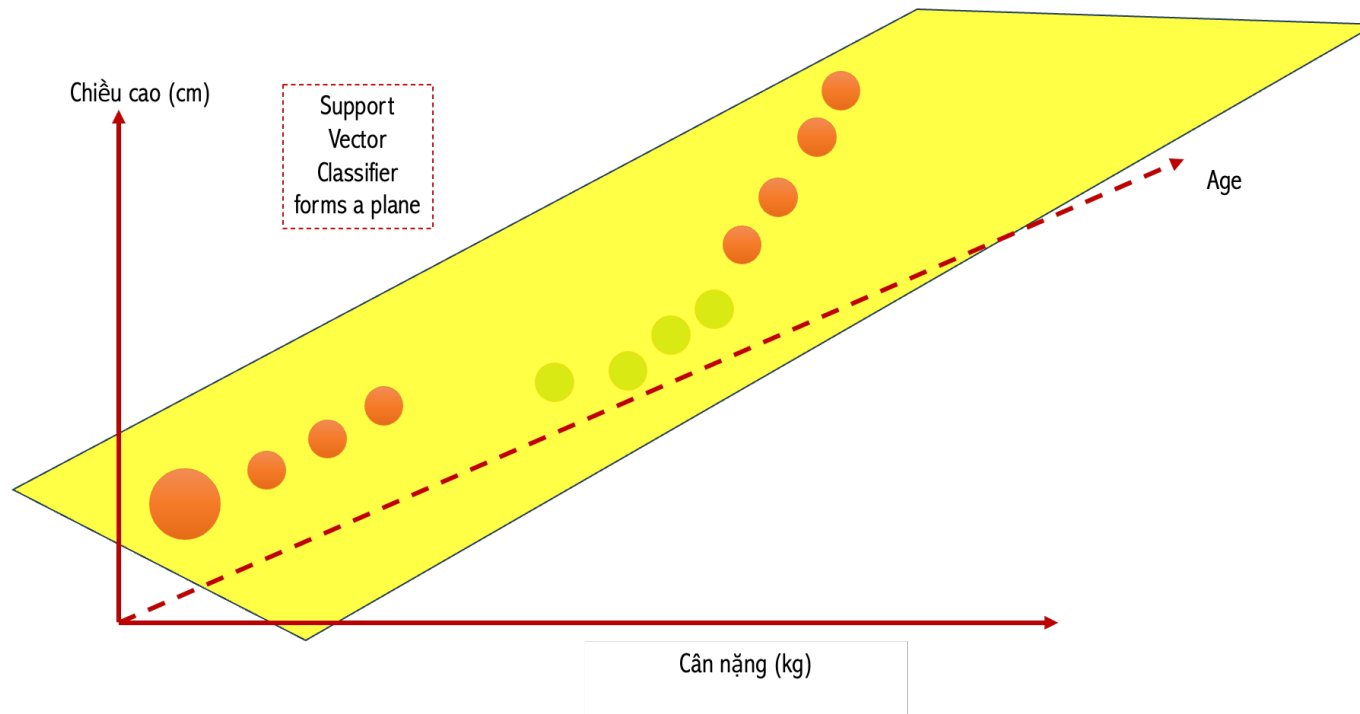
Radial Kernel: Intuition

Polynomial Kernel: $(a \times b)^1 = (a^1) (b^1)$ with $r = 0$ and $d = 1$

Polynomial Kernel: $(a \times b)^2 = (a^2) (b^2)$ with $r = 0$ and $d = 2$

Polynomial Kernel: $(a \times b)^3 = (a^3) (b^3)$ with $r = 0$ and $d = 3$

$$(a^1) (b^1) + (a^2) (b^2) + (a^3) (b^3) = (a, a^2, a^3) \cdot (b, b^2, b^3)$$



We do not actually do the transformation, we just solve for Dot Product to get high dimensional relationship!

Radial Kernel: Intuition

Polynomial Kernel: $(a \times b)^1 = (a^1) (b^1)$ with $r = 0$ and $d = 1$

Polynomial Kernel: $(a \times b)^2 = (a^2) (b^2)$ with $r = 0$ and $d = 2$

Polynomial Kernel: $(a \times b)^3 = (a^3) (b^3)$ with $r = 0$ and $d = 3$

Polynomial Kernel: $(a \times b)^{\dots} = (a^{\dots}) (b^{\dots})$ with $r = 0$ and $d = \dots$

Polynomial Kernel: $(a \times b)^{\infty} = (a^{\infty}) (b^{\infty})$ with $r = 0$ and $d = \infty$

$$(a^1) (b^1) + (a^2) (b^2) + (a^3) (b^3) + \dots + (a^{\infty}) (b^{\infty}) = (a, a^2, a^3, \dots, a^{\infty},) \cdot (b, b^2, b^3, \dots, b^{\infty})$$

$$\text{Radial kernel: } e^{-r(a-b)^2} = e^{-\frac{1}{2}(a^2-2ab+b^2)} = e^{-\frac{1}{2}(a^2+b^2)} e^{ab}$$

Taylor Series:

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots + \frac{f^{\infty}(a)}{\infty!}(x-a)^{\infty}$$

$$e^x = e^a + \frac{e^a}{1!}(x-a) + \frac{e^a}{2!}(x-a)^2 + \frac{e^a}{3!}(x-a)^3 + \dots + \frac{e^a}{\infty!}(x-a)^{\infty}$$

$$e^x = e^0 + \frac{e^0}{1!}(x-a) + \frac{e^0}{2!}(x-0)^2 + \frac{e^0}{3!}(x-0)^3 + \dots + \frac{e^0}{\infty!}(x-0)^{\infty}$$

Radial Kernel: Intuition

Polynomial kernel

$$(a^0)(b^0) + (a^1)(b^1) + (a^2)(b^2) + (a^3)(b^3) + \dots + (a^\infty)(b^\infty) = (a, a^2, a^3, \dots, a^\infty) \cdot (b, b^2, b^3, \dots, b^\infty)$$

Radial kernel

$$\text{Radial kernel: } e^{-\gamma(a-b)^2} = e^{-\frac{1}{2}(a^2-2ab+b^2)} = e^{-\frac{1}{2}(a^2+b^2)} e^{ab}$$

$$e^{ab} = e^0 + \frac{e^0}{1!} ab + \frac{e^0}{2!} (ab)^2 + \frac{e^0}{3!} (ab)^3 + \dots + \frac{e^0}{\infty!} (ab)^\infty$$

$$e^{ab} = 1 + \frac{1}{1!} ab + \frac{1}{2!} (ab)^2 + \frac{1}{3!} (ab)^3 + \dots + \frac{1}{\infty!} (ab)^\infty$$

$$(a^0)(b^0) + (a^1)(b^1) + (a^2)(b^2) + (a^3)(b^3) + \dots + (a^\infty)(b^\infty) = (1, a, a^2, a^3, \dots, a^\infty) \cdot (1, b, b^2, b^3, \dots, b^\infty)$$

$$e^{ab} = \left(1, \sqrt{\frac{1}{1!}} a, \sqrt{\frac{1}{2!}} a^2, \sqrt{\frac{1}{3!}} a^3, \dots, \sqrt{\frac{1}{\infty!}} a^\infty, \right) \cdot \left(1, \sqrt{\frac{1}{1!}} b, \sqrt{\frac{1}{2!}} b^2, \sqrt{\frac{1}{3!}} b^3, \dots, \sqrt{\frac{1}{\infty!}} b^\infty, \right)$$

Radial Kernel: Intuition

$$\text{Radial kernel: } e^{-\gamma(a-b)^2} = e^{-\frac{1}{2}(a^2-2ab+b^2)} = e^{-\frac{1}{2}(a^2+b^2)} e^{ab}$$

$$e^{ab} = \left(1, \sqrt{\frac{1}{1!}}a, \sqrt{\frac{1}{2!}}a^2, \sqrt{\frac{1}{3!}}a^3, \dots, \sqrt{\frac{1}{\infty!}}a^\infty, \right) \cdot \left(1, \sqrt{\frac{1}{1!}}b, \sqrt{\frac{1}{2!}}b^2, \sqrt{\frac{1}{3!}}b^3, \dots, \sqrt{\frac{1}{\infty!}}b^\infty, \right)$$

$$e^{-\frac{1}{2}(a^2+b^2)} e^{ab} = e^{-\frac{1}{2}(a^2+b^2)} \left(1, \sqrt{\frac{1}{1!}}a, \sqrt{\frac{1}{2!}}a^2, \sqrt{\frac{1}{3!}}a^3, \dots, \sqrt{\frac{1}{\infty!}}a^\infty, \right) \cdot \left(1, \sqrt{\frac{1}{1!}}b, \sqrt{\frac{1}{2!}}b^2, \sqrt{\frac{1}{3!}}b^3, \dots, \sqrt{\frac{1}{\infty!}}b^\infty, \right)$$

$$e^{-\frac{1}{2}(a^2+b^2)} e^{ab} = \left(\delta, \delta \sqrt{\frac{1}{1!}}a, \delta \sqrt{\frac{1}{2!}}a^2, \delta \sqrt{\frac{1}{3!}}a^3, \dots, \delta \sqrt{\frac{1}{\infty!}}a^\infty, \right) \cdot \left(\delta, \delta \sqrt{\frac{1}{1!}}b, \delta \sqrt{\frac{1}{2!}}b^2, \delta \sqrt{\frac{1}{3!}}b^3, \dots, \delta \sqrt{\frac{1}{\infty!}}b^\infty, \right)$$

$$\delta = \sqrt{e^{-\frac{1}{2}(a^2+b^2)}}$$

Radial kernel is equal to a dot product that has coordinates for infinite number of dimensions

Radial Kernel: Example

The relationship between two points in infinite – dimensions

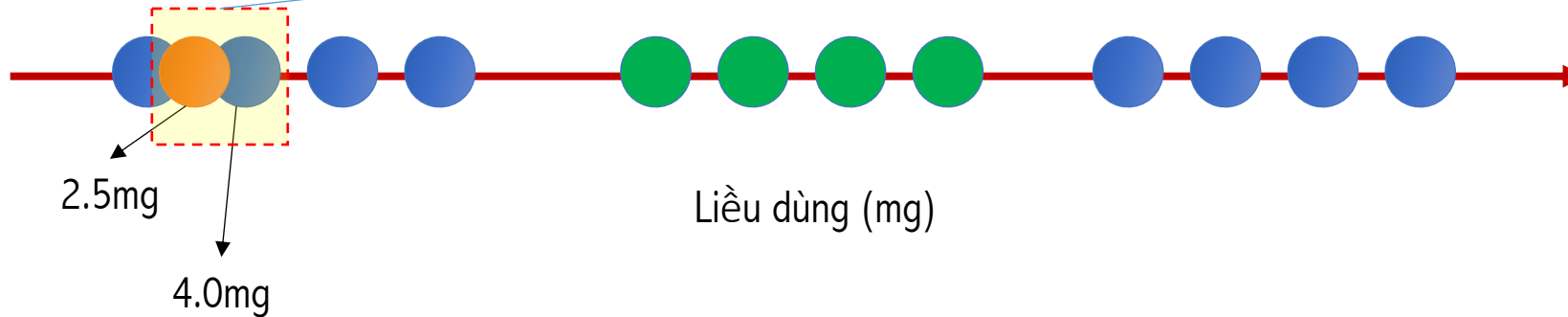
$$\gamma = 1$$

$$e^{-(2.5-4)^2} = 0.11$$

$$\gamma = 2$$

$$e^{-(2.5-4)^2} = 0.01$$

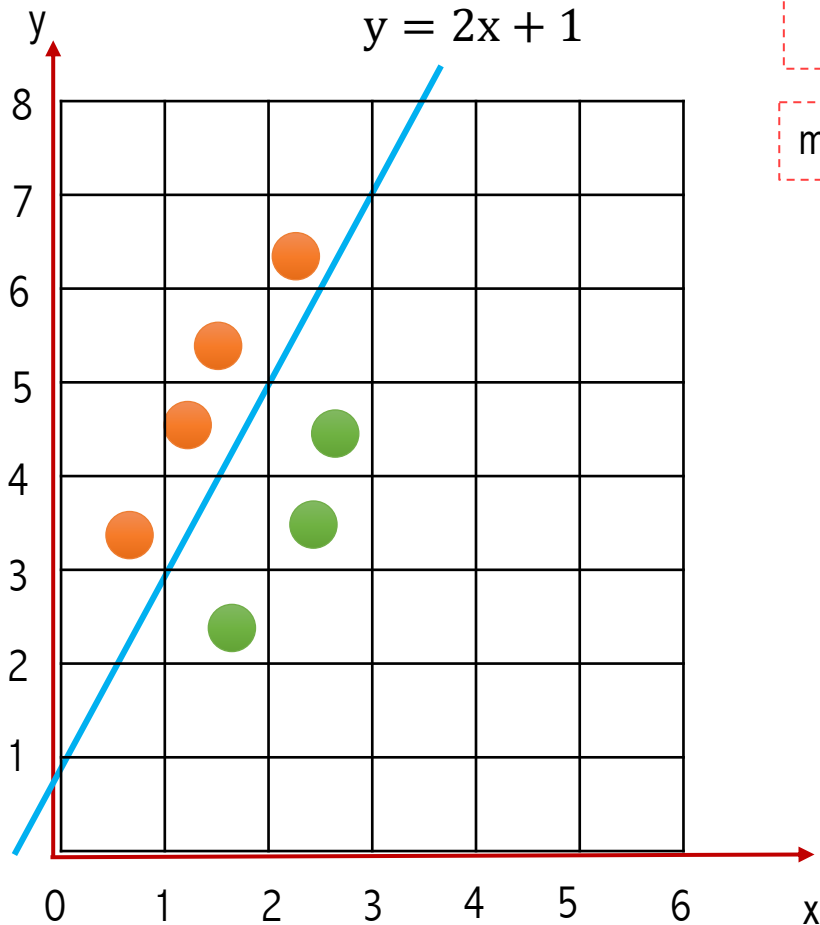
γ scales the amount of influence two points have each other



Outline

- Maximal Margin Classifier
- Support Vector Classifier
- Support Vector Machine
- Polynomial Kernel
- Radial Basic Function Kernel (RBF)
- Example

Line Equation Review



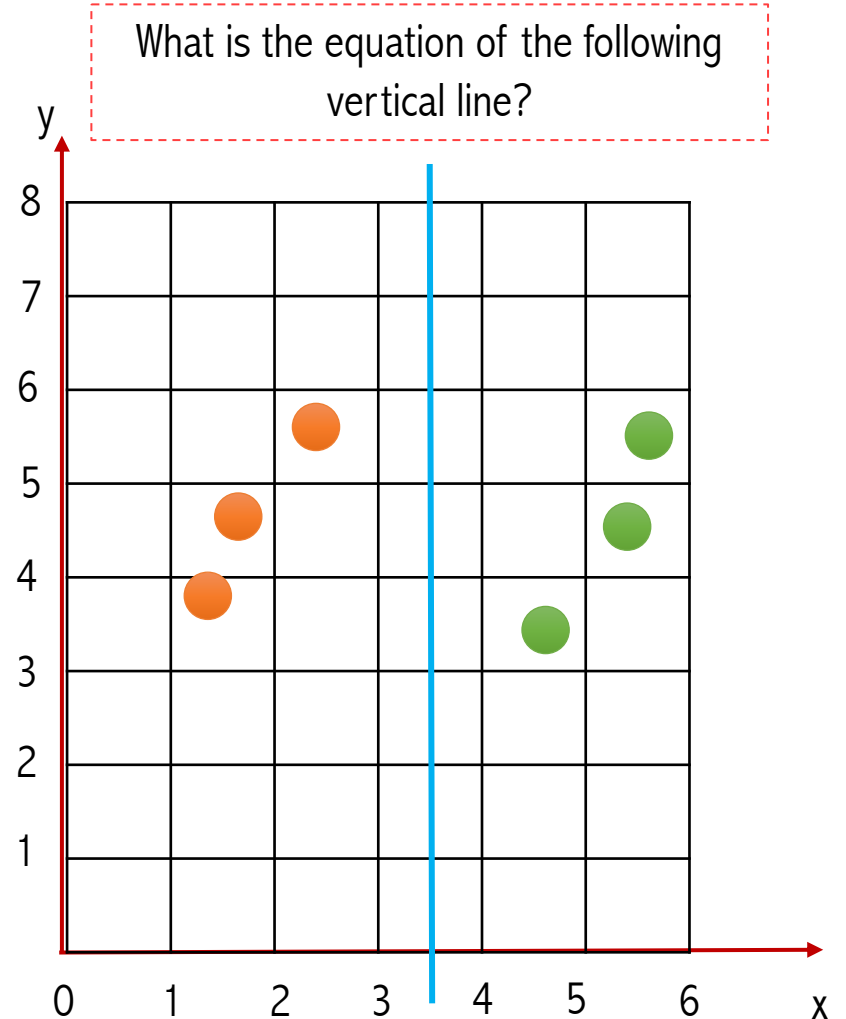
$$y = mx + b$$

m : slope, b: intercept

General form of the equation of the straight line

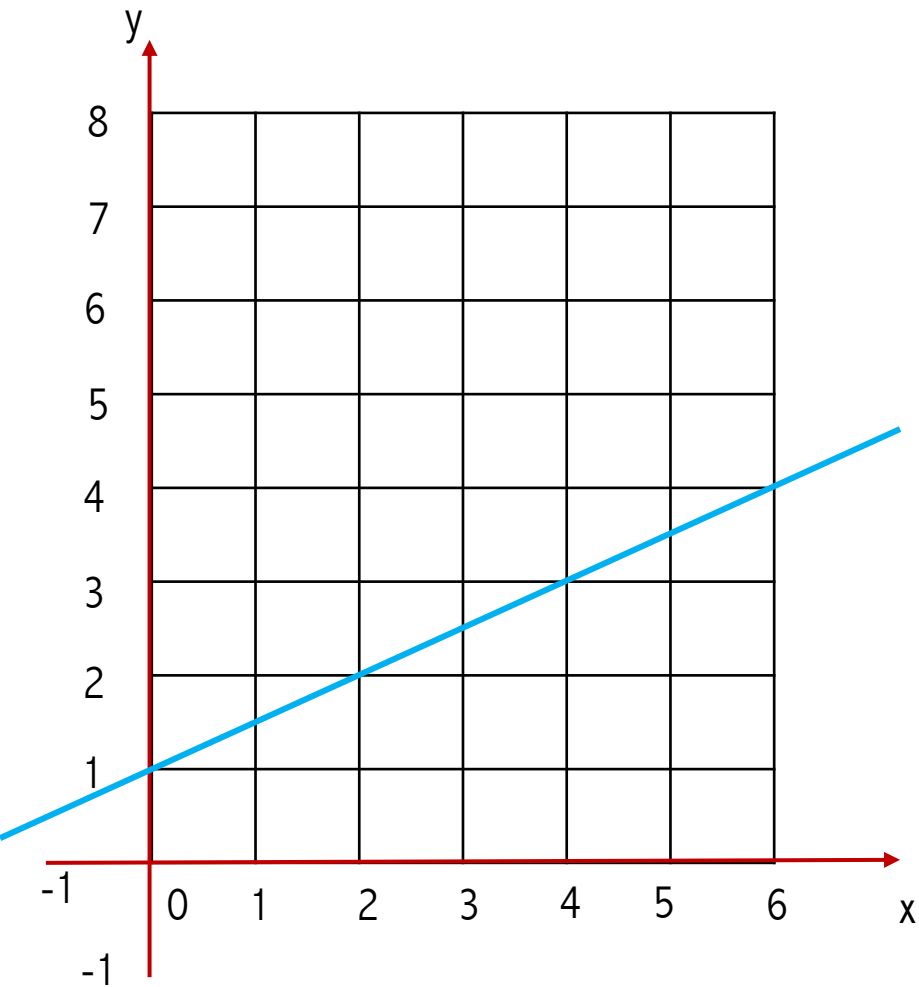
$$Ax + By + C = 0$$

$$y = -\frac{A}{B}x - \frac{C}{B}$$

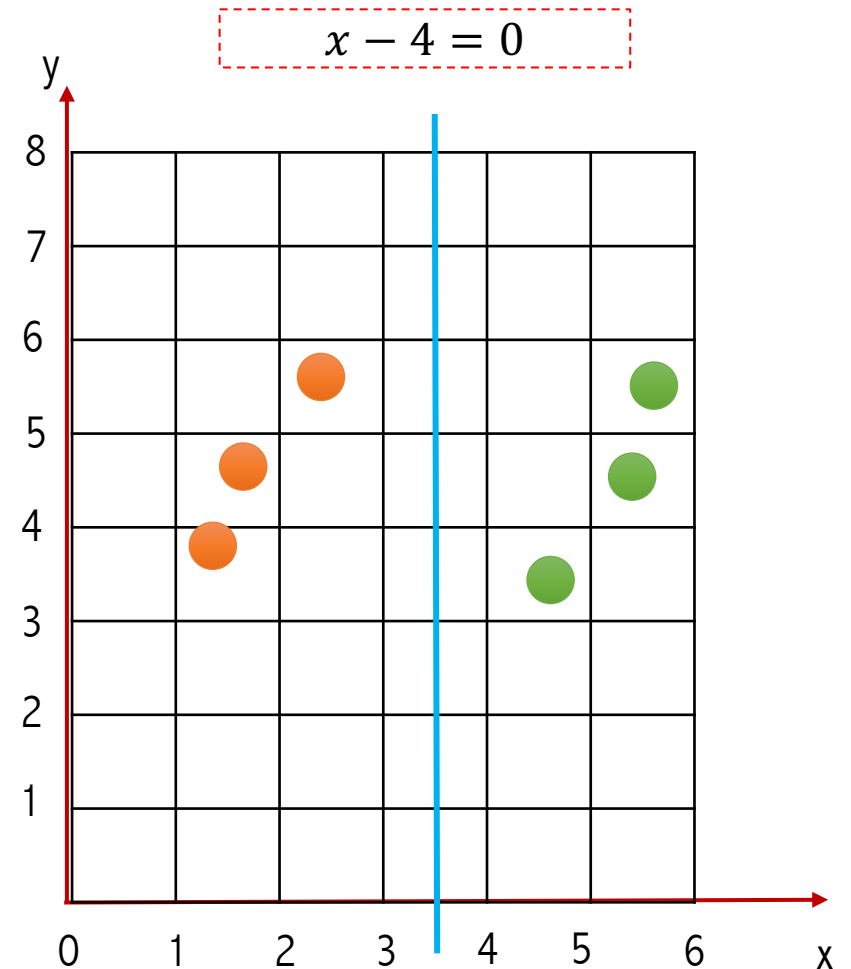


What is the equation of the following vertical line?

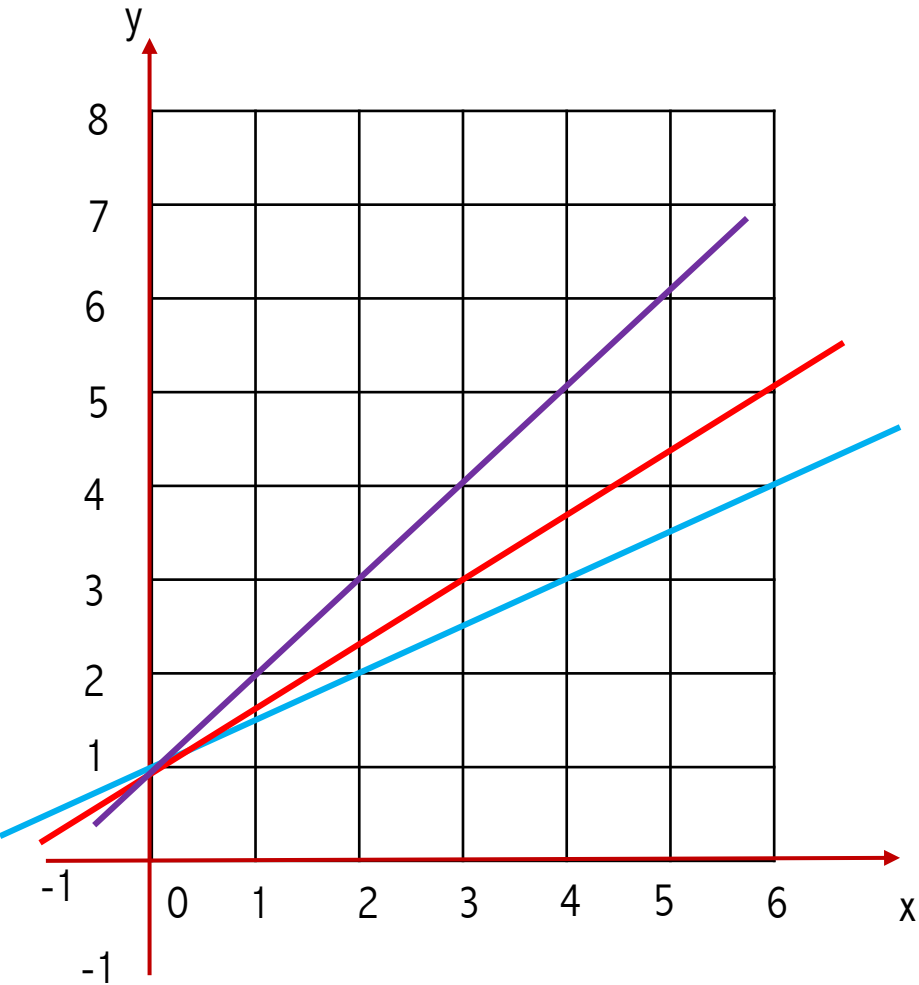
Line Equation Review



$$\begin{aligned} y &= 0.5x + 1 \\ \downarrow \\ -0.5x + y - 1 &= 0 \\ \downarrow \\ -2x + 4y - 4 &= 0 \\ \downarrow \\ 2x - 4y + 4 &= 0 \end{aligned}$$



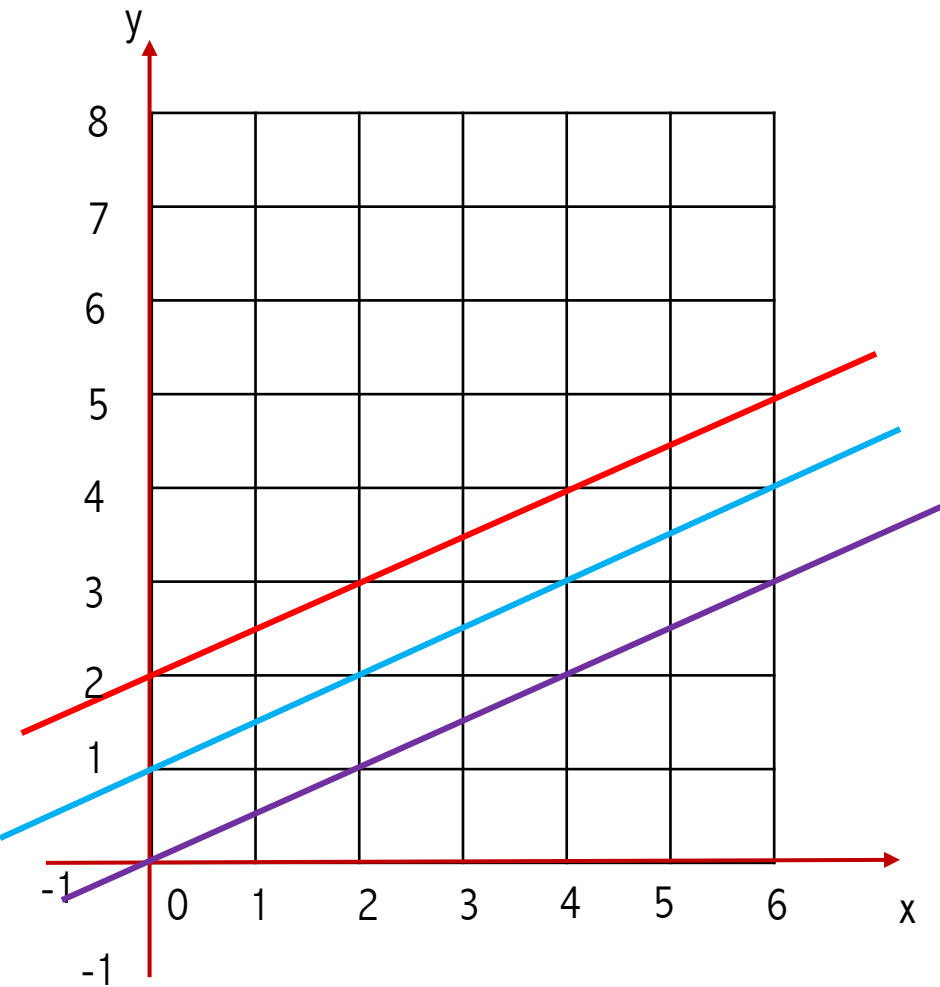
Line Equation Review: A, B, C Changes



$$\begin{array}{c} Ax + By + C = 0 \\ \text{---} -2x + 4y - 4 = 0 \text{---} \\ \swarrow \quad \searrow \\ \text{---} -x + 4y - 4 = 0 \text{---} \quad \text{---} -4x + 4y - 4 = 0 \text{---} \end{array}$$

When we change A, the Line is rotating around 1.

Line Equation Review: A, B, C Changes



$$y = -\frac{A}{B}x - \frac{C}{B}$$

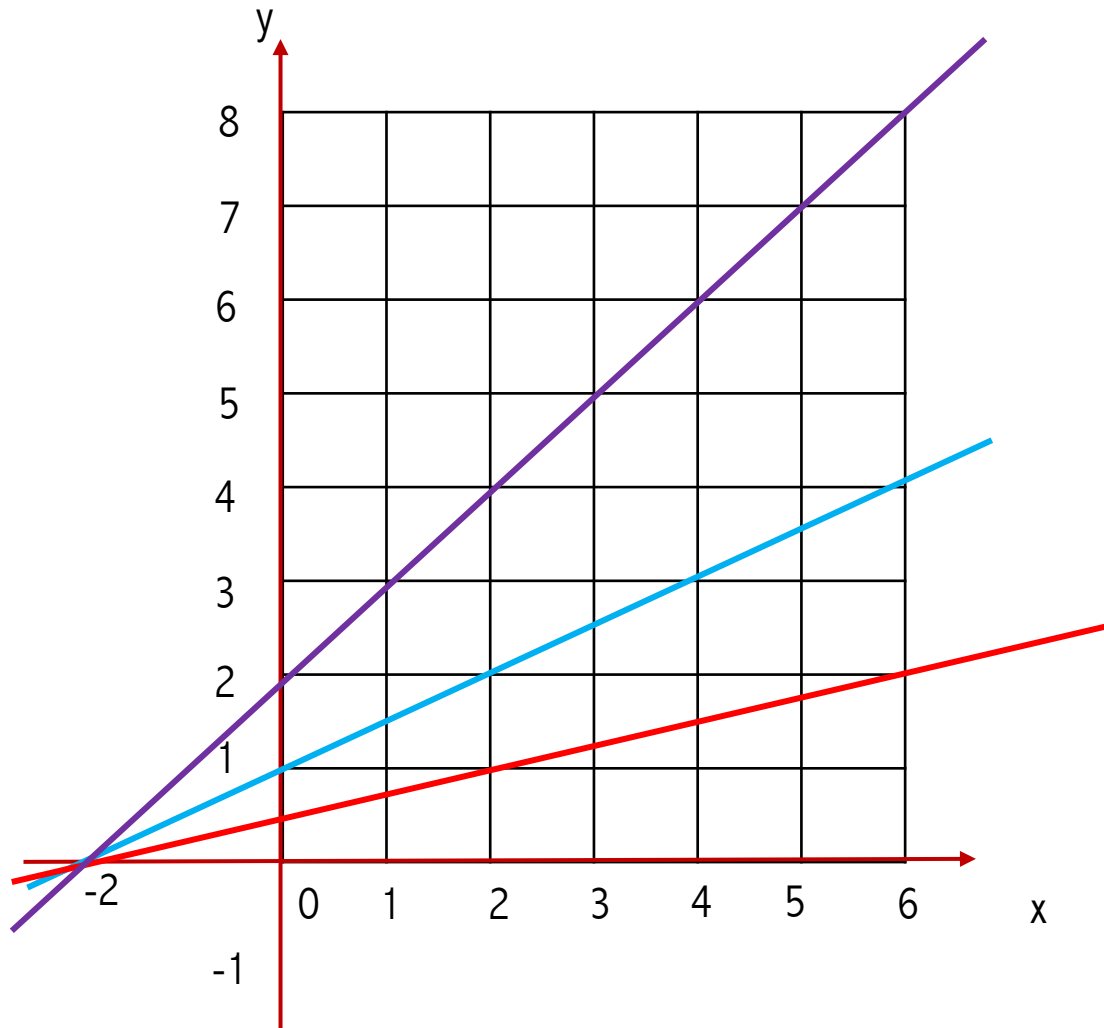
$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2x + 4y - 8 = 0$$

$$-2x + 4y - 0 = 0$$

Line Equation Review: A, B, C Changes



$$y = -\frac{A}{B}x - \frac{C}{B}$$

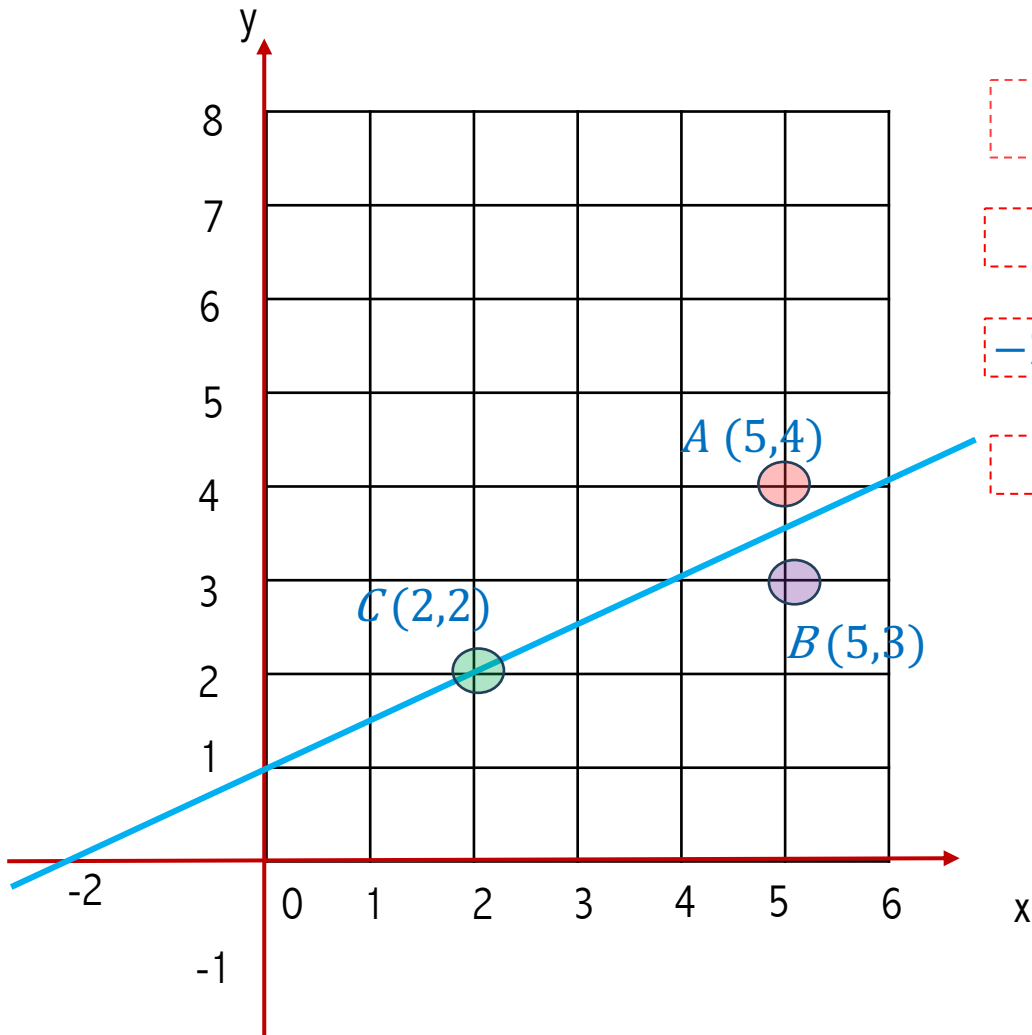
$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2x + 8y - 4 = 0$$

$$-4x + 2y - 4 = 0$$

General Form Line Equation



$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2 * 5 + 4 * 4 - 4 = 2$$

A is above the line

$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$\begin{aligned} -2 * 5 + 3 * 4 - 4 \\ = -2 \end{aligned}$$

B is below the line

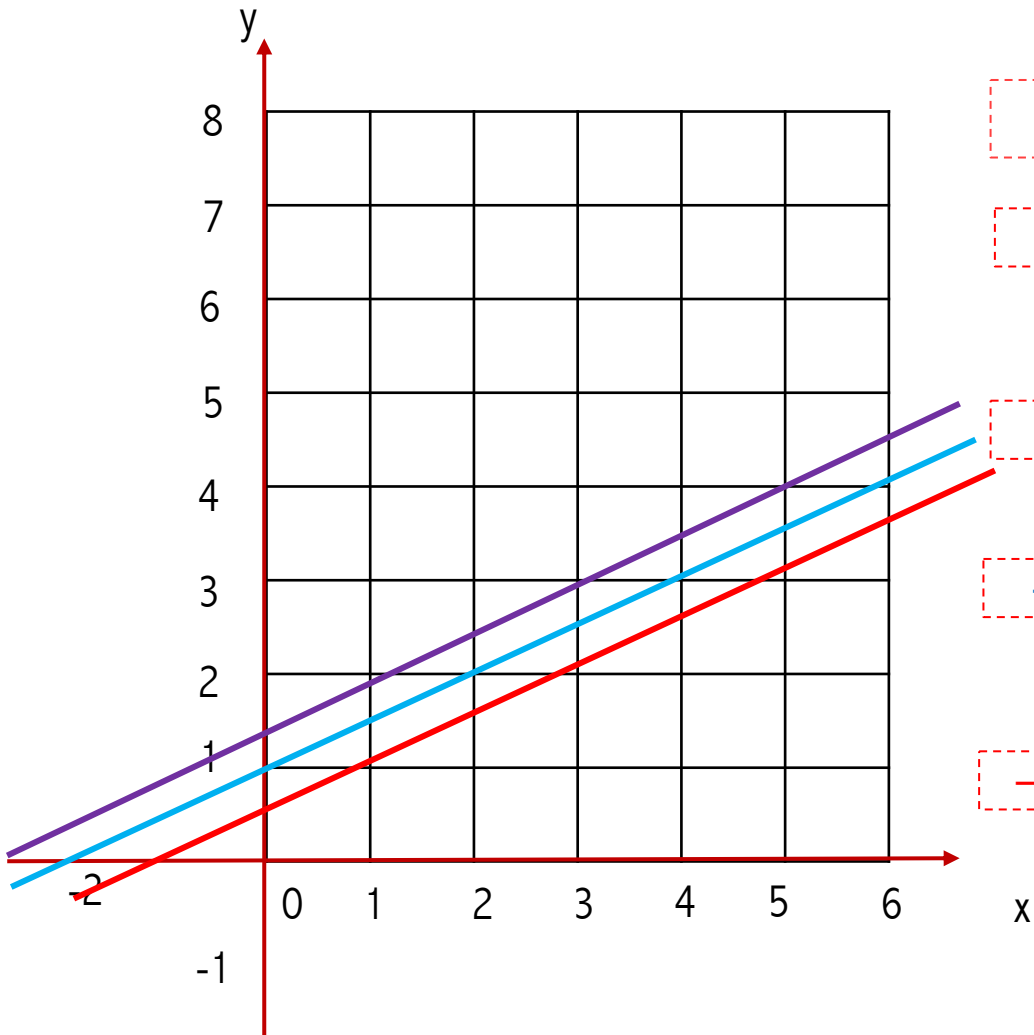
$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2 * 2 + 2 * 4 - 4 = 0$$

C is on the line

General Form Line Equation



$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2x + 4y - 4 = 1$$

$$-2x + 4y - 4 = 0$$

$$-2x + 4y - 4 = -1$$

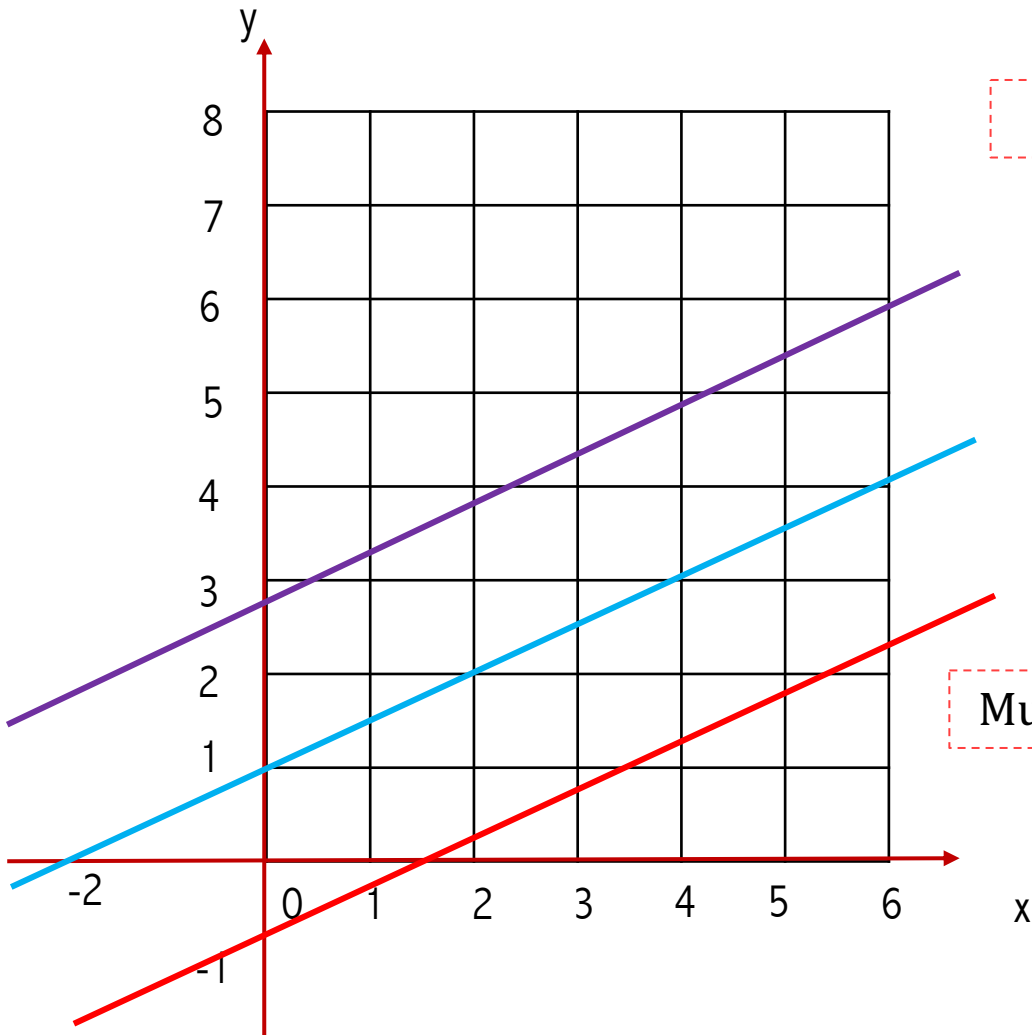
$$y = 0.5x + 1$$

$$y = 0.5x + 1.25$$

$$y = 0.5x + 1$$

$$y = 0.5x + 0.75$$

General Form Line Equation



$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2x + 4y - 4 = 1$$

$$-2x + 4y - 4 = -1$$



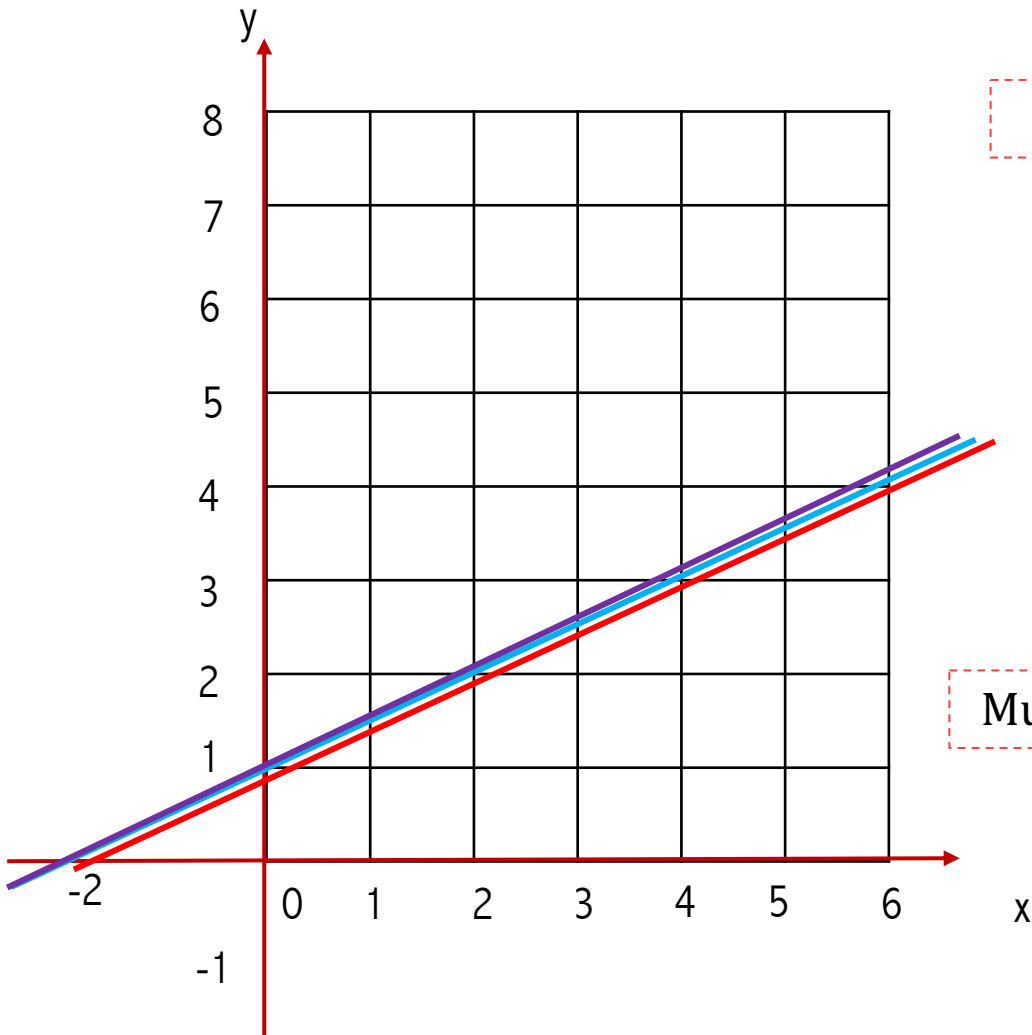
Multiple with a factor smaller than 1

$$y = 0.5x + 1$$

$$y = 0.5x + 1.25$$

$$y = 0.5x + 0.75$$

General Form Line Equation



$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$-2x + 4y - 4 = 1$$

$$-2x + 4y - 4 = -1$$

$$y = 0.5x + 1$$

$$y = 0.5x + 1.25$$

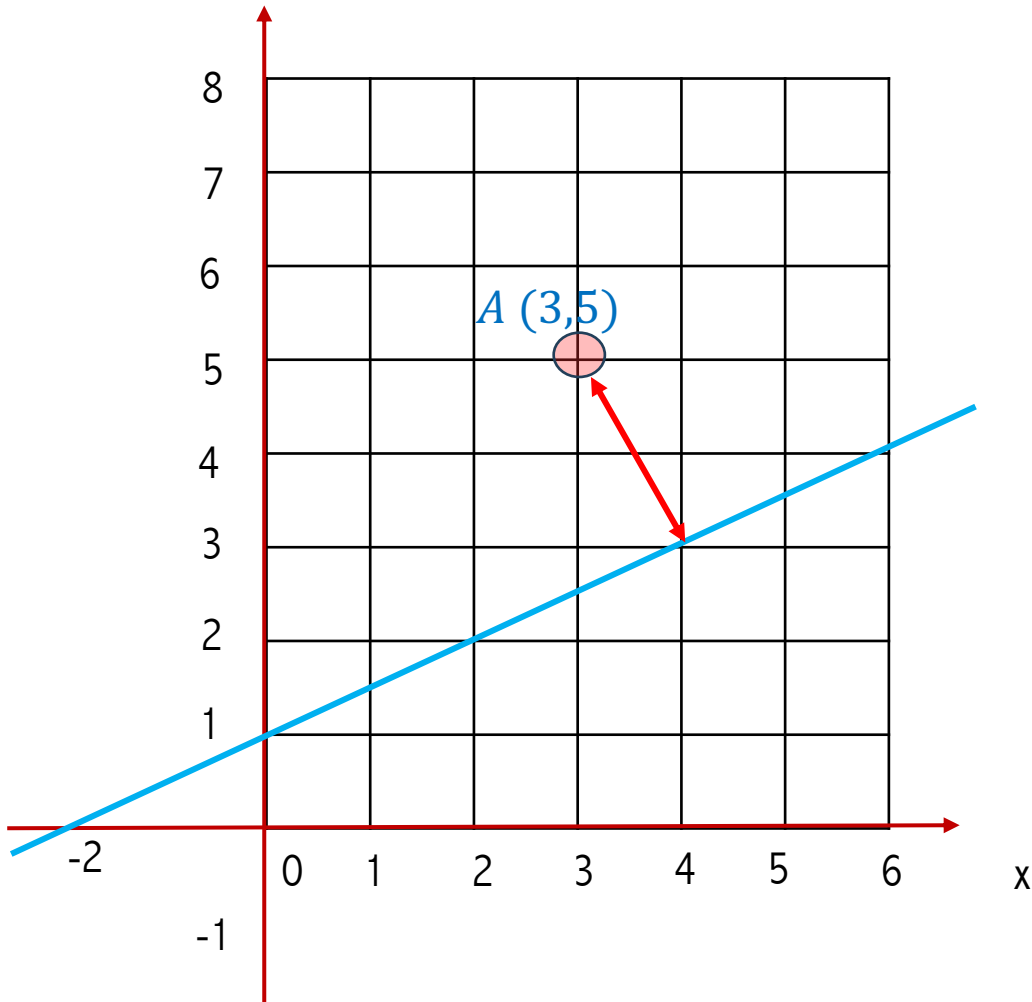
$$y = 0.5x + 0.75$$



Multiple with a factor greater than 1

This method is used on SVM to increase or decrease the Margin

Distance between a point and a line



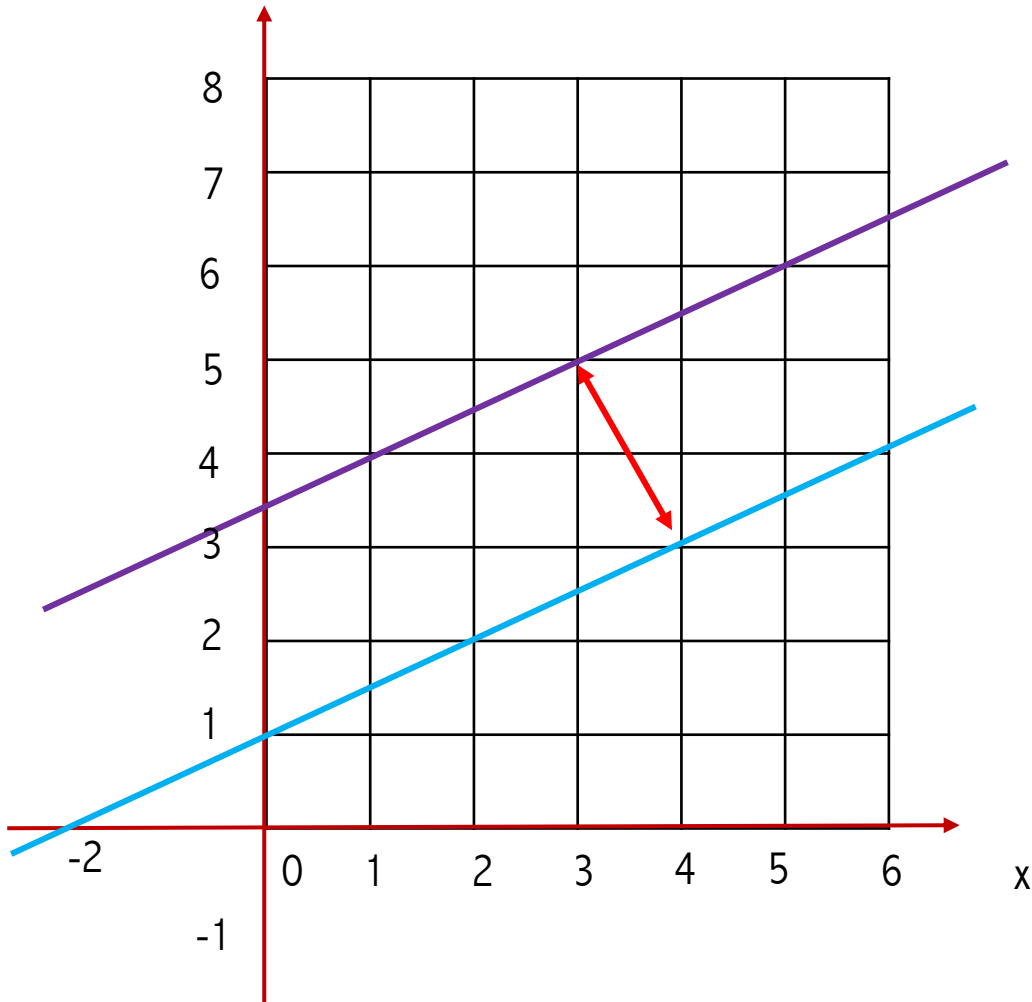
$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 0$$

$$d = \frac{|Ax + By + C|}{\sqrt{A^2 + B^2}}$$

$$d = \frac{|(-2)*3 + 4*5 + (-4)|}{\sqrt{(-2)^2 + (4)^2}} = 2.236$$

Distance between parallel lines



$$Ax + By + C = 0$$

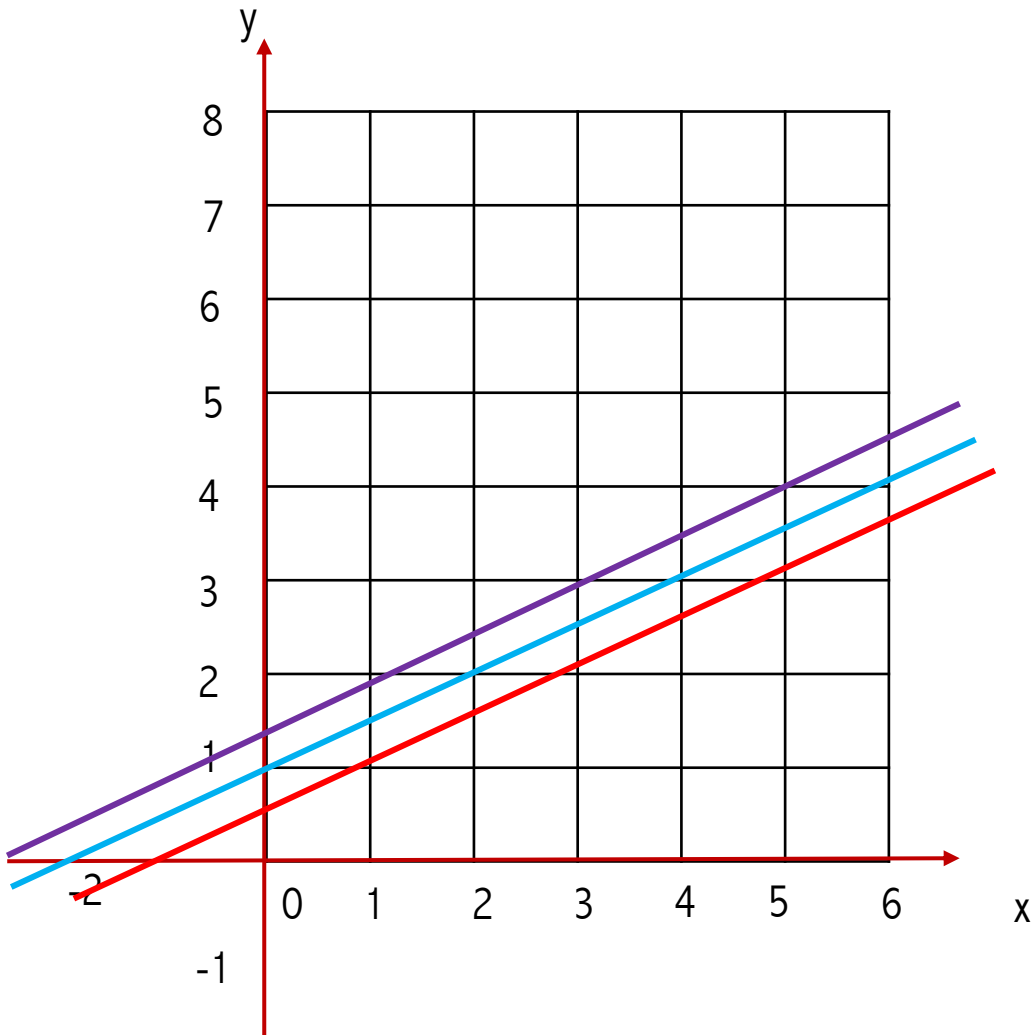
$$-2x + 4y - 14 = 0$$

$$-2x + 4y - 4 = 0$$

$$d = \frac{|C_1 - C_2|}{\sqrt{A^2 + B^2}}$$

$$d = \frac{|(-4) - (-14)|}{\sqrt{(-2)^2 + (4)^2}} = \frac{10}{\sqrt{20}} = 2.236$$

Distance between parallel lines



$$Ax + By + C = 0$$

$$-2x + 4y - 4 = 1$$

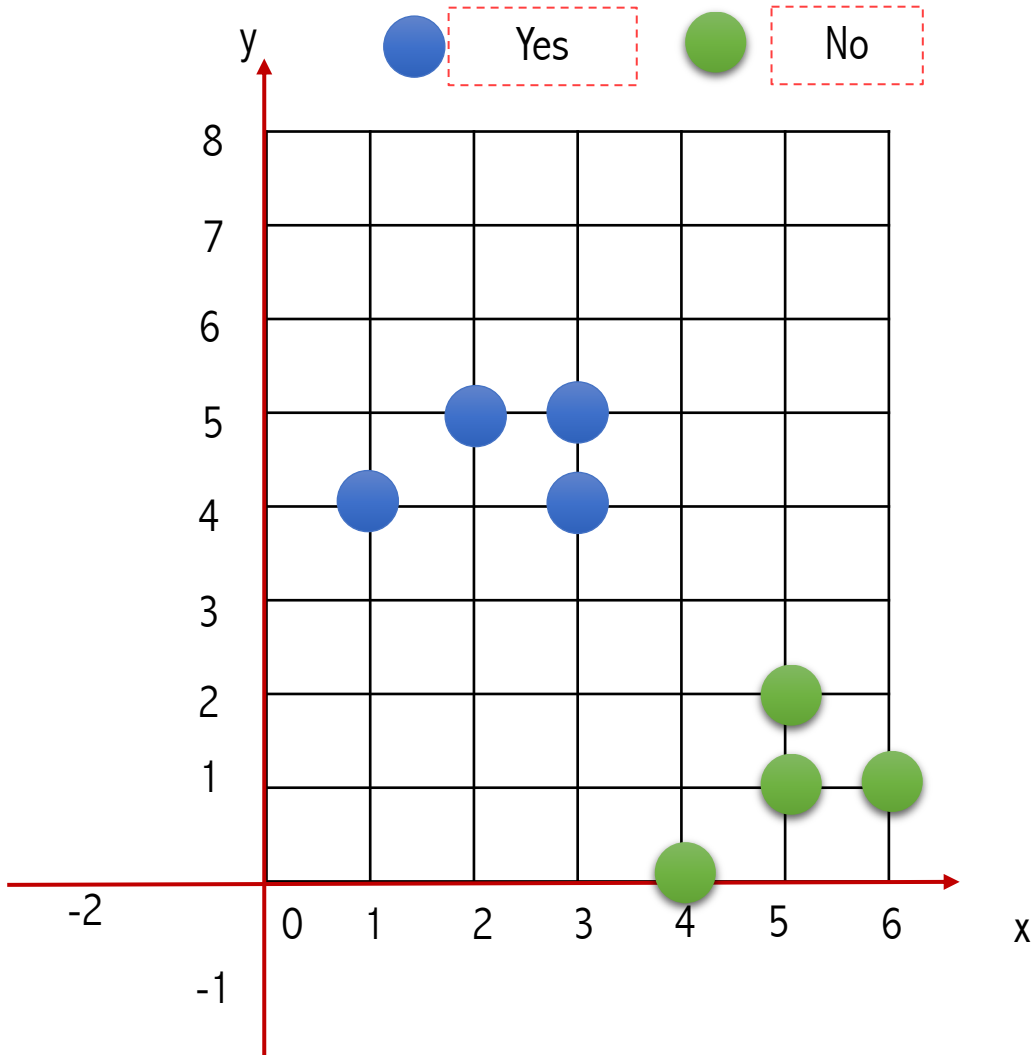
$$-2x + 4y - 4 = 0$$

$$-2x + 4y - 4 = -1$$

$$d = \frac{|(-3) - (-5)|}{\sqrt{(-2)^2 + (4)^2}} = \frac{2}{\sqrt{20}} = 0.447$$

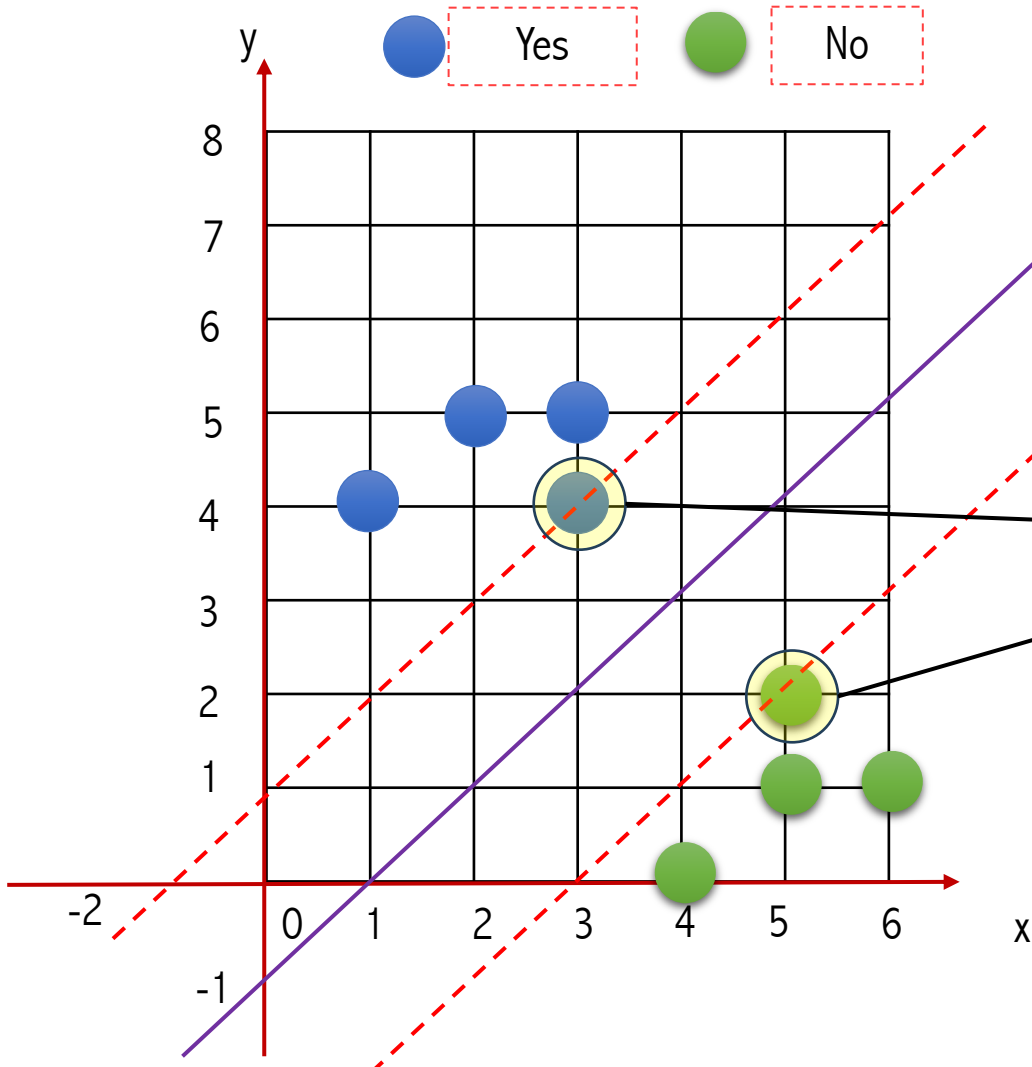
$$d = \frac{|2|}{\sqrt{A^2 + B^2}} = \frac{2}{\|W\|_2}$$

Example



Blood Pressure	Cholesterol Level	Disease
1	2	Yes
2	5	Yes
3	5	Yes
3	4	Yes
6	1	No
4	0	No
5	2	No
5	1	No

Example



$$-4x + 4y + (-4) = 0$$

$$-4x + 4y + 4 = 0$$

$$-4x + 4y + 12 = 0$$

$$-0.5x + 0.5y + 0.5 = 1$$

$$-0.5x + 0.5y + 0.5 = 0$$

$$-0.5x + 0.5y + 0.5 = -1$$

Support vector

Find a value of K so that the left-hand side is equal to one

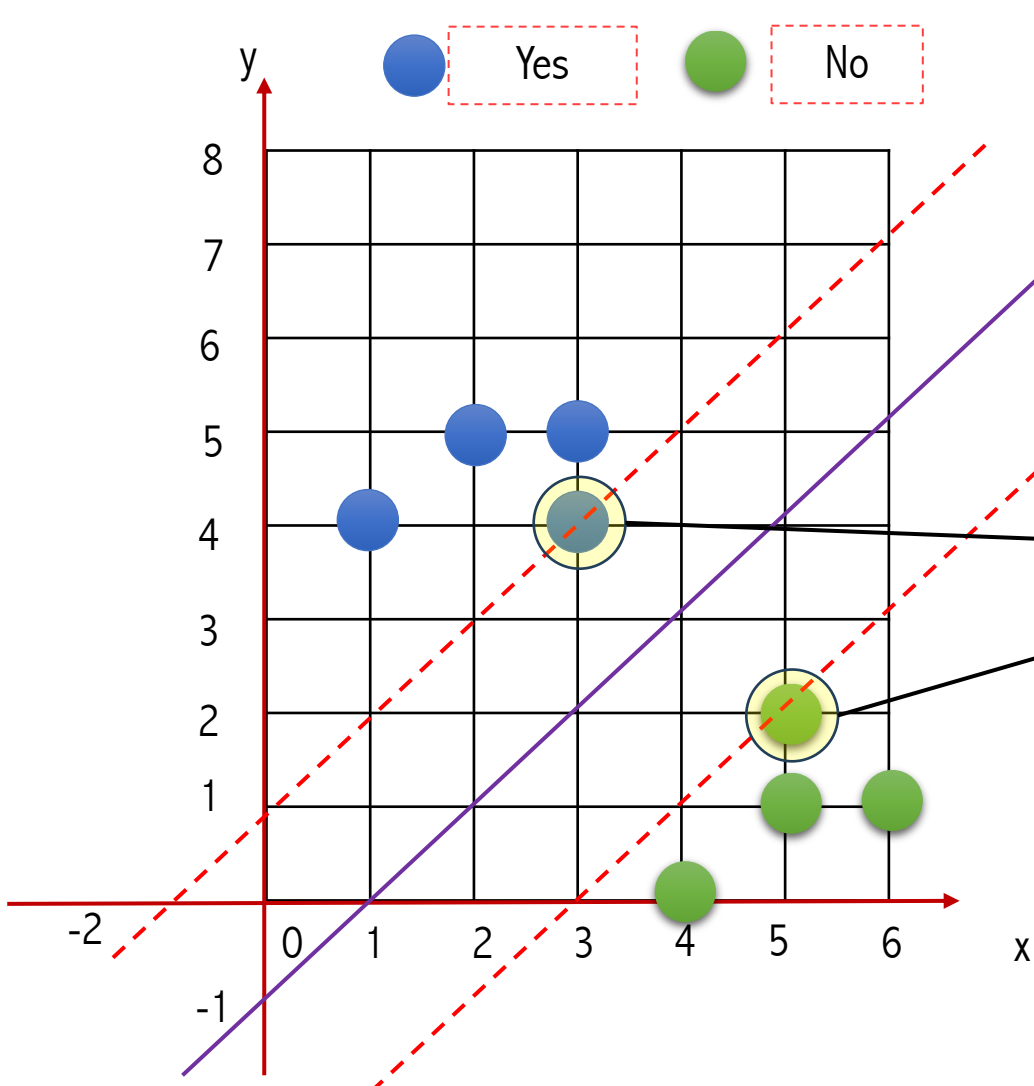
$$k(-4x + 4y + 4) = 1$$

$$k(-4 \cdot 3 + 4 \cdot 4 + 4) = 1$$

$$k \cdot 8 = 1$$

$$k = 1/8$$

Example



$$-4x + 4y + (-4) = 0$$

$$k = 1/8$$

$$-0.5x + 0.5y + 0.5 = 1$$

$$-4x + 4y + 4 = 0$$

$$-0.5x + 0.5y + 0.5 = 0$$

$$-4x + 4y + 12 = 0$$

$$-0.5x + 0.5y + 0.5 = -1$$

Find a value of K so that the left-hand side is equal to 1

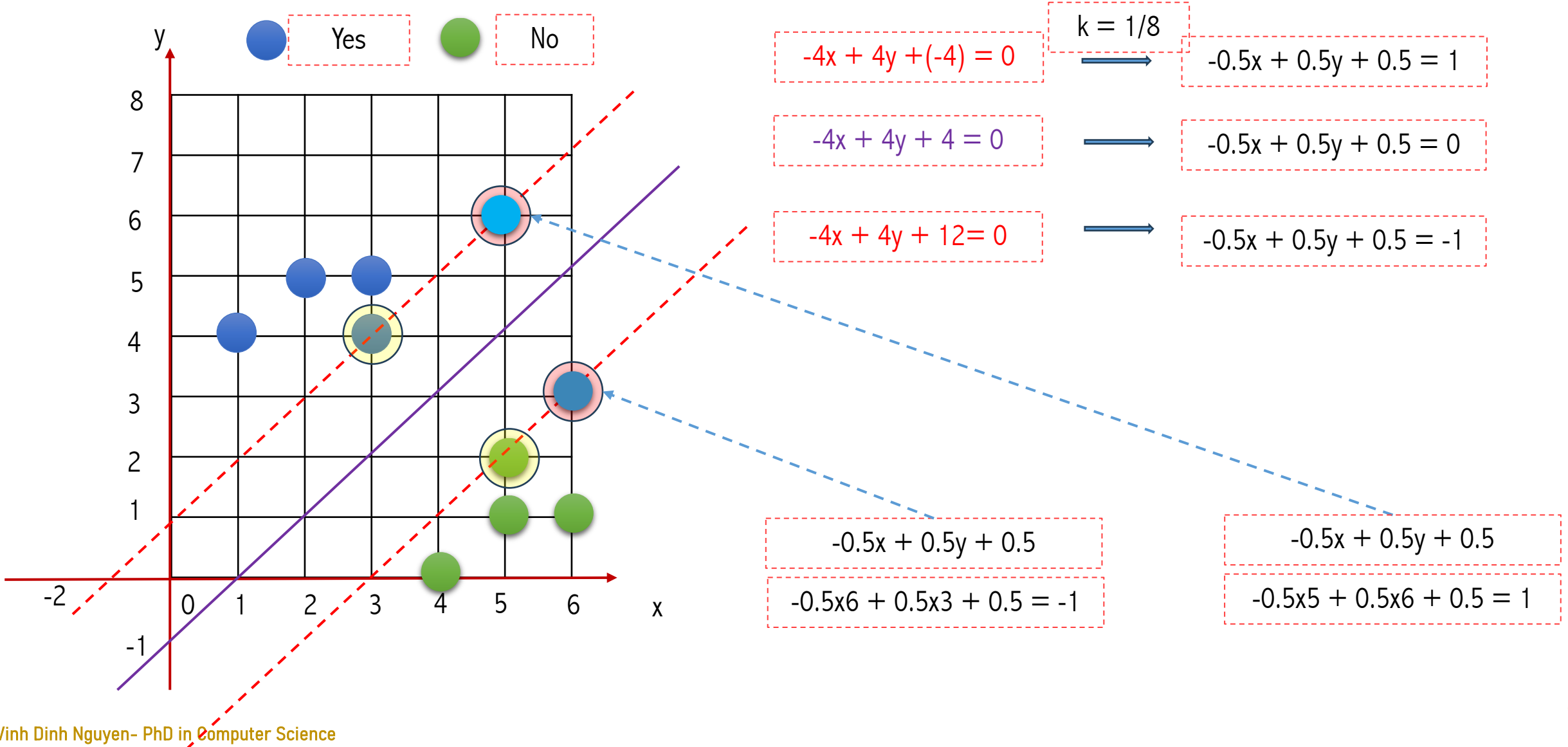
$$k(-4x + 4y + 4) = 1$$

$$k(-4 \cdot 3 + 4 \cdot 4 + 4) = 1$$

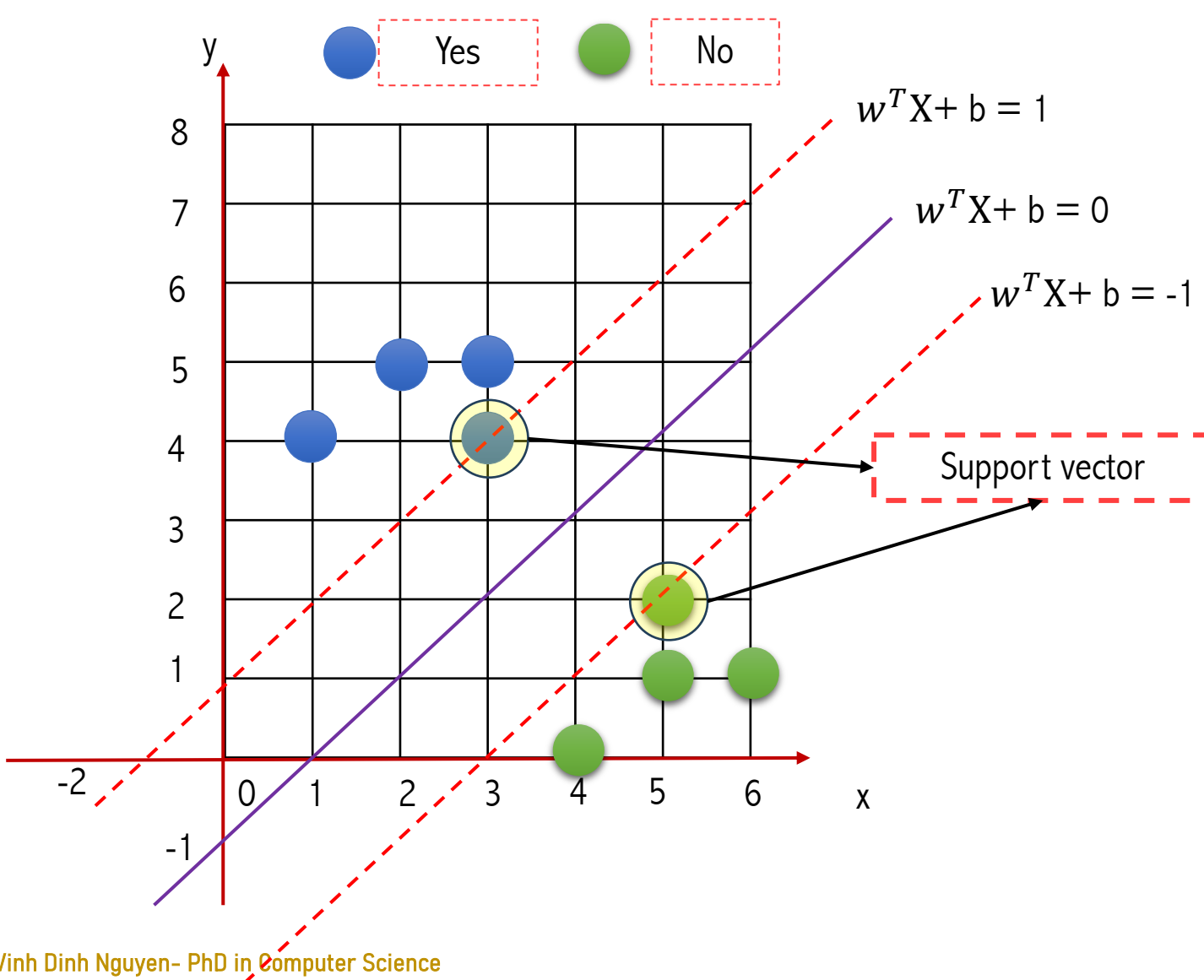
$$k \cdot 8 = 1$$

$$k = 1/8$$

Example



Example



$$b = 0.5$$

$$-0.5x + 0.5y + b = 1$$

$$-0.5x + 0.5y + b = 0$$

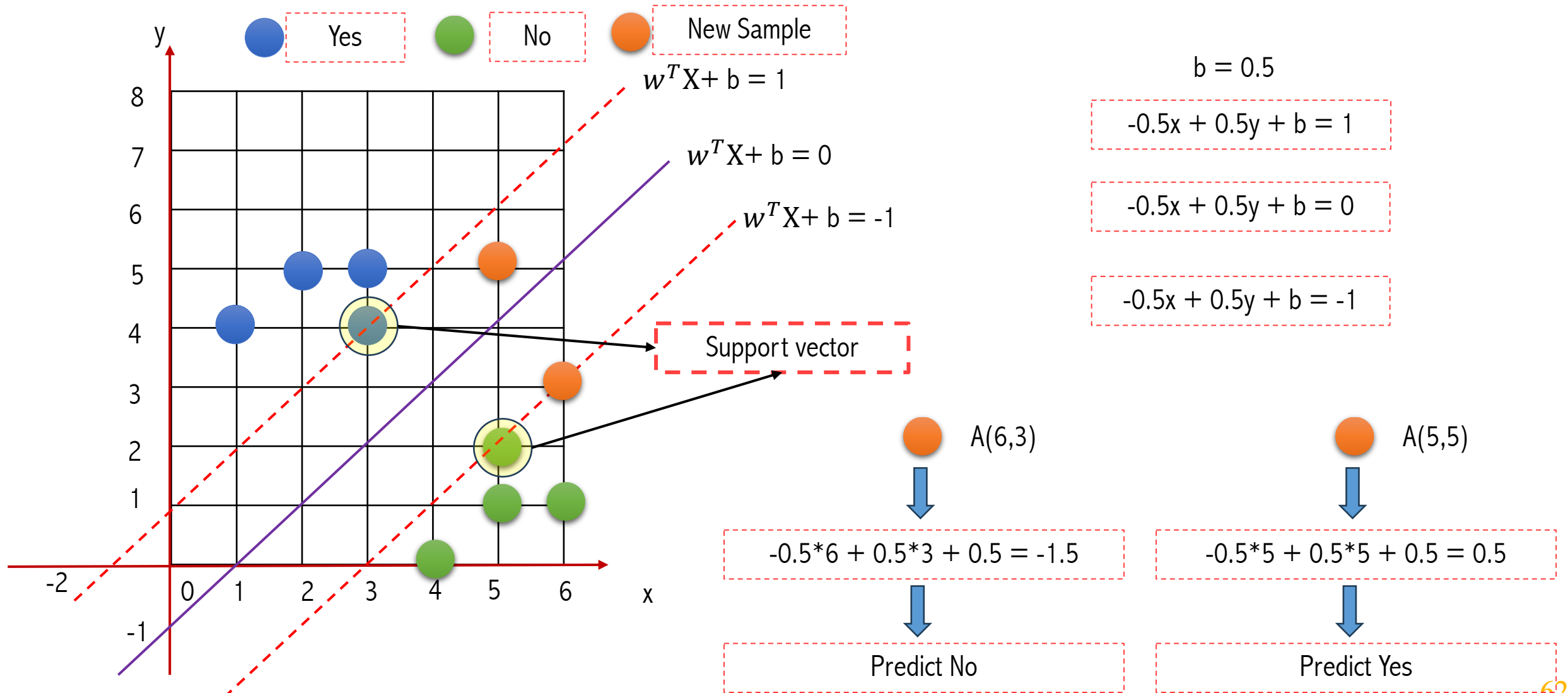
$$-0.5x + 0.5y + b = -1$$

$$w^T X + b = 0$$

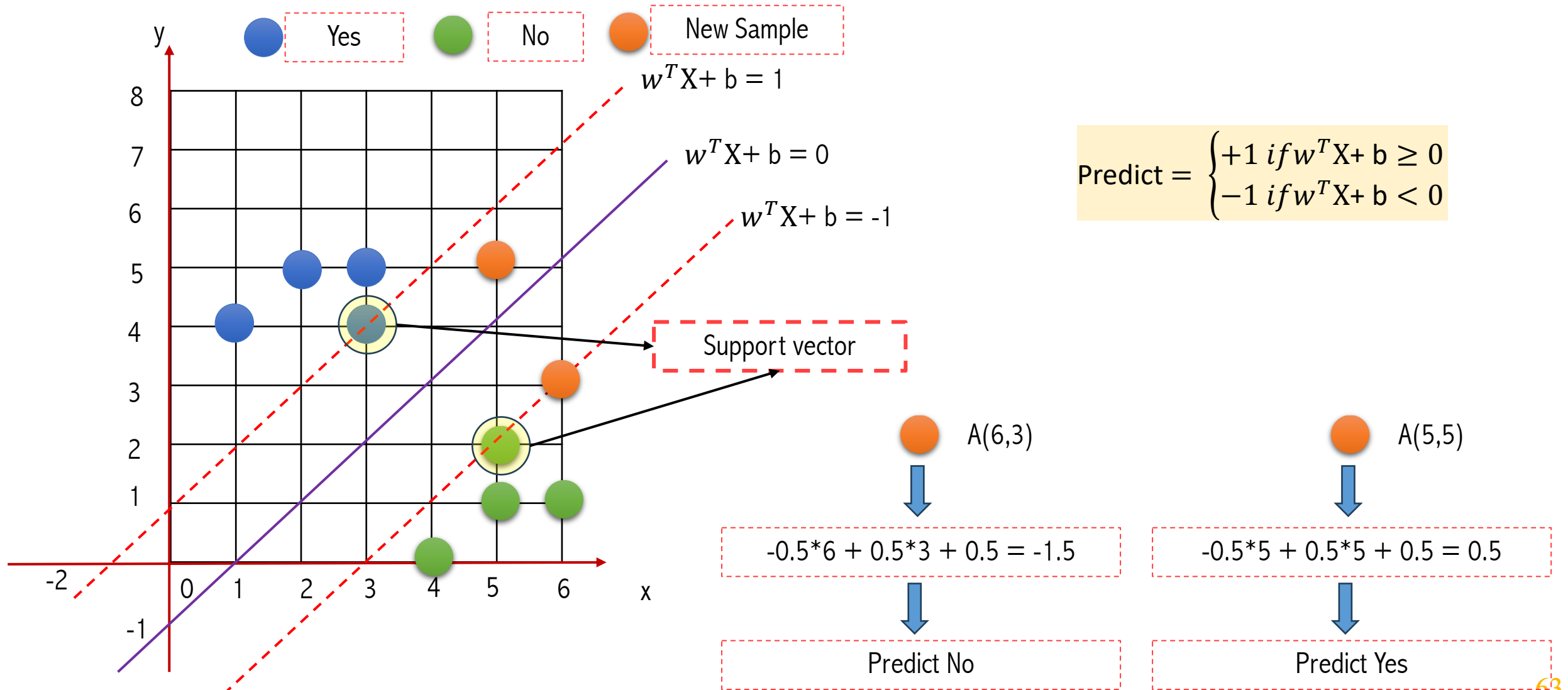
$$X = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$W = \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix}$$

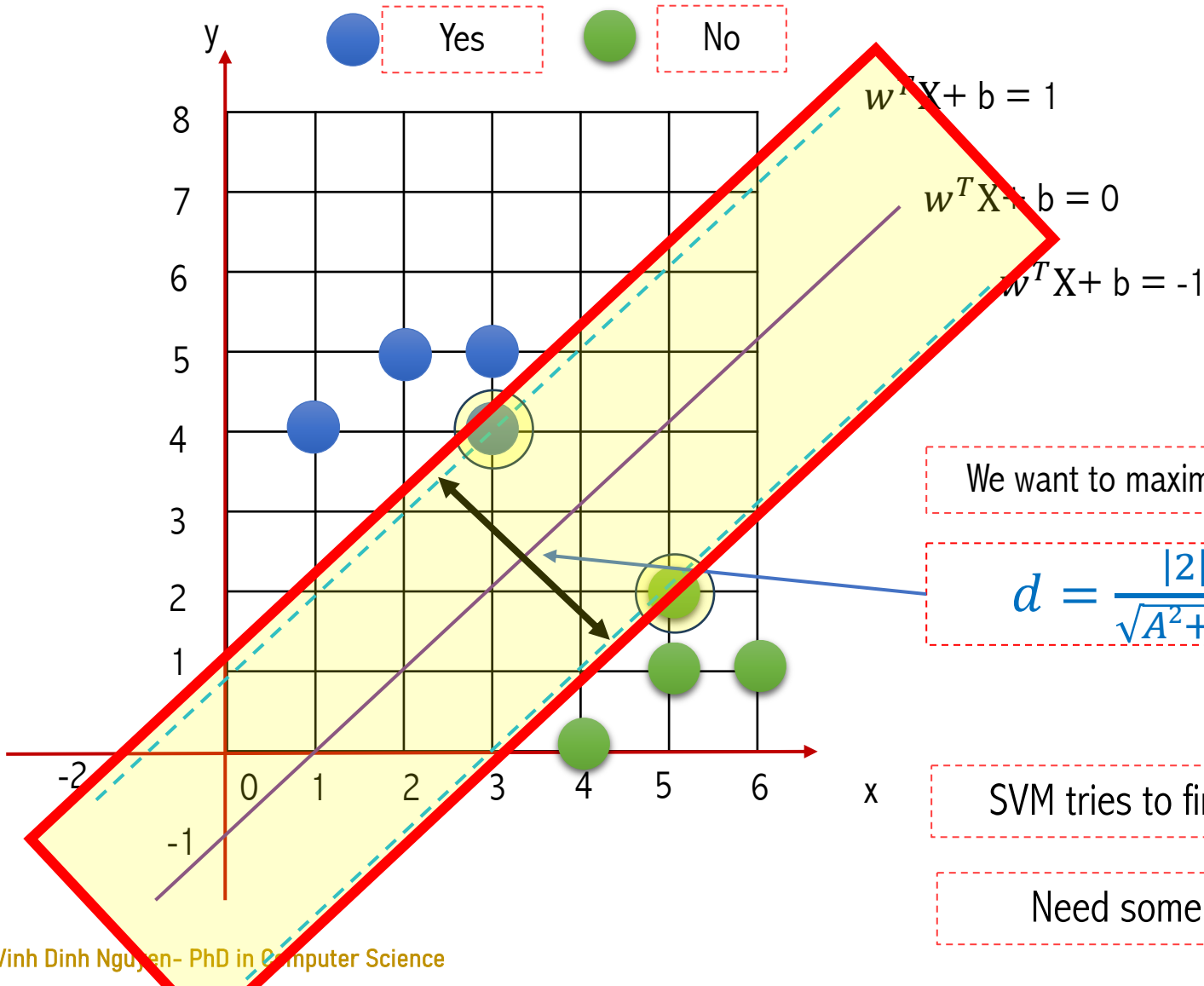
Example



Example



Example



$$Y = \begin{cases} +1 & \text{if } w^T X + b \geq 0 \\ -1 & \text{if } w^T X + b < 0 \end{cases}$$

We want to maximize the distance d

We want to minimize $\|W\|_2$

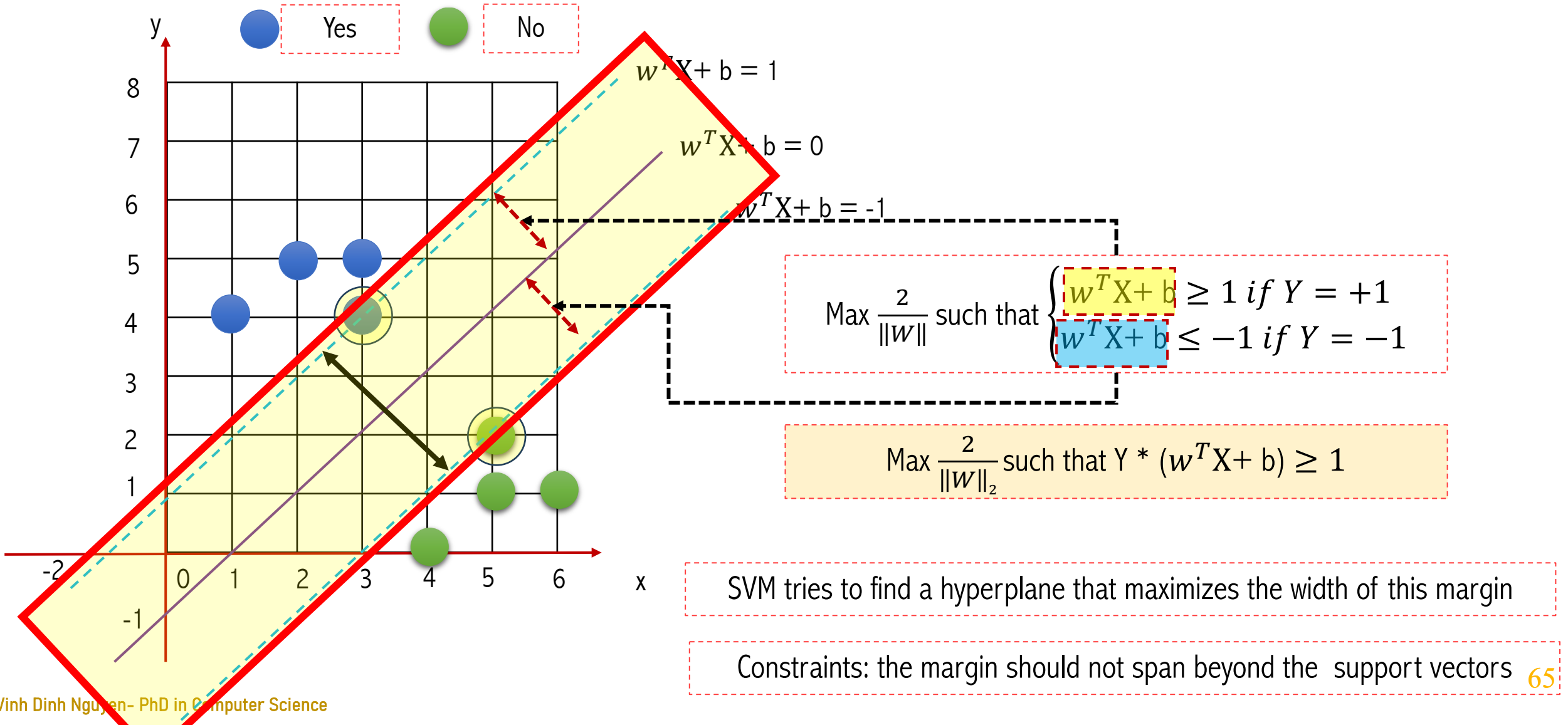
$$d = \frac{|2|}{\sqrt{A^2 + B^2}} = \frac{2}{\|W\|_2}$$

$$d = \frac{|2|}{\sqrt{A^2 + B^2}} = \frac{2}{\|W\|_2}$$

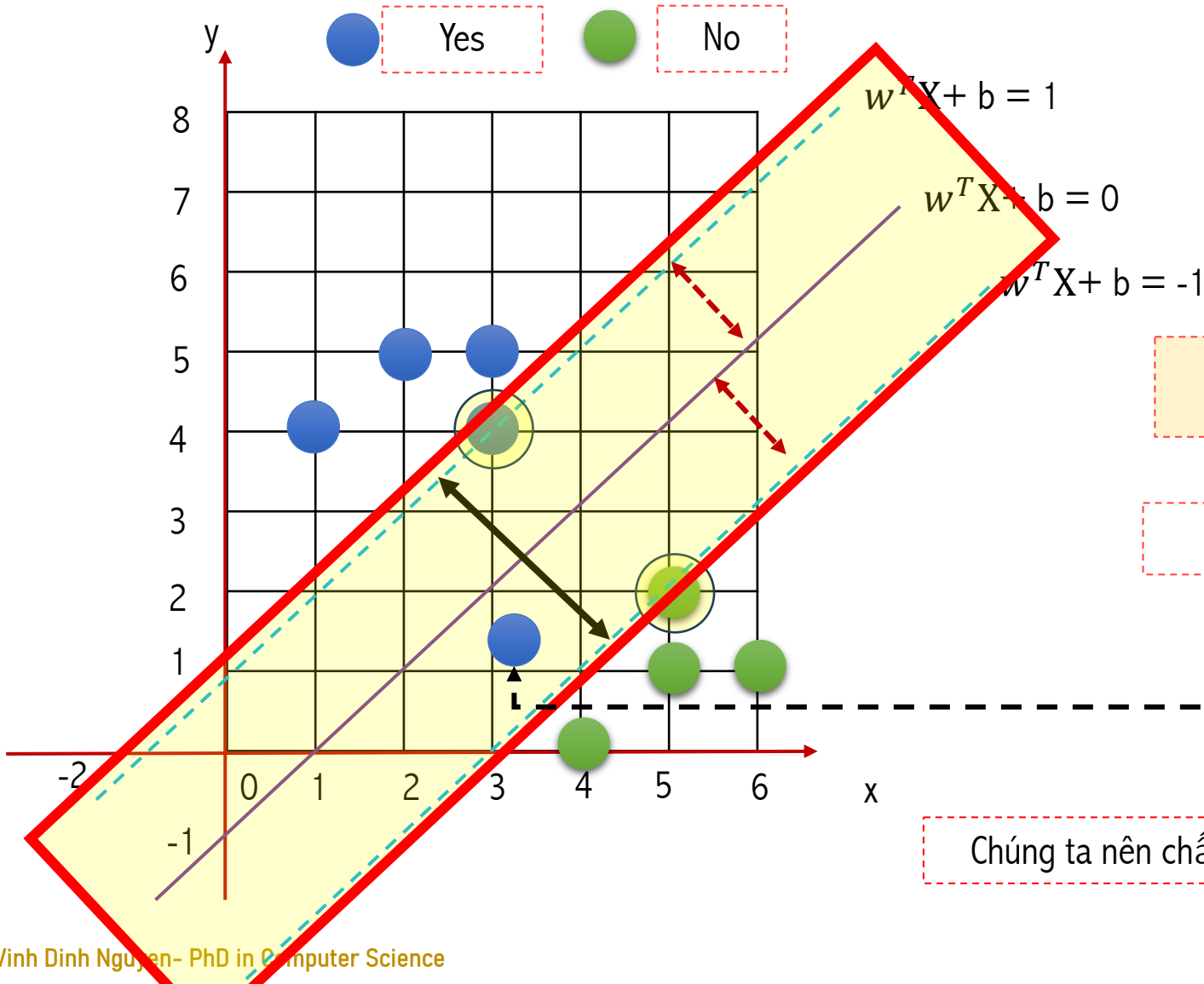
SVM tries to find a hyperplane that maximizes the width of this margin

Need some constraints because the margin can be infinitely large

Example



Example

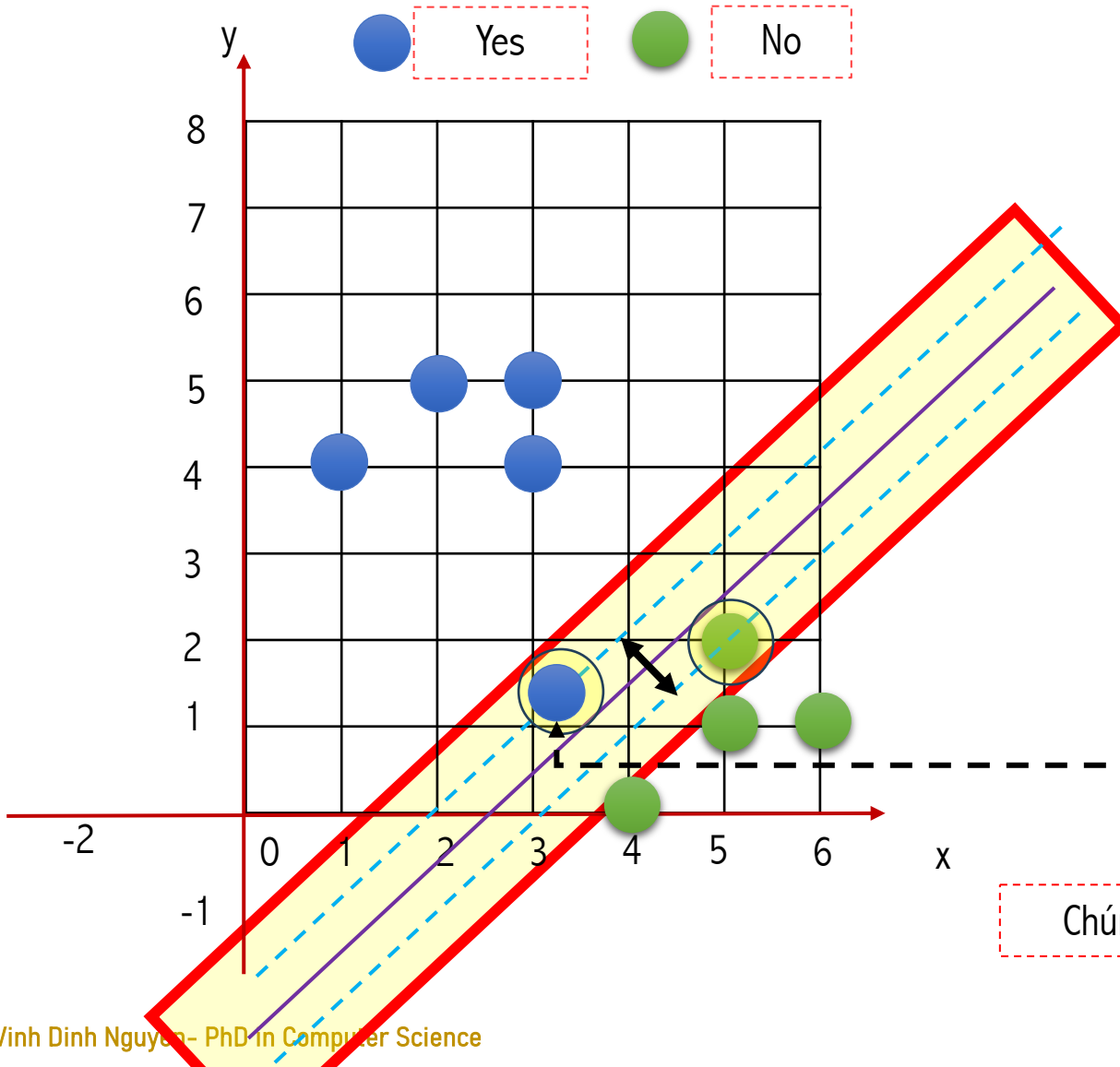


$$\text{Max } \frac{2}{\|w\|_2} \text{ such that } Y * (w^T X + b) \geq 1$$

What should we do if a green data point is noise?

Chúng ta nên chấp nhận như là miss-classification hay là thay đổi hyperplane

Example

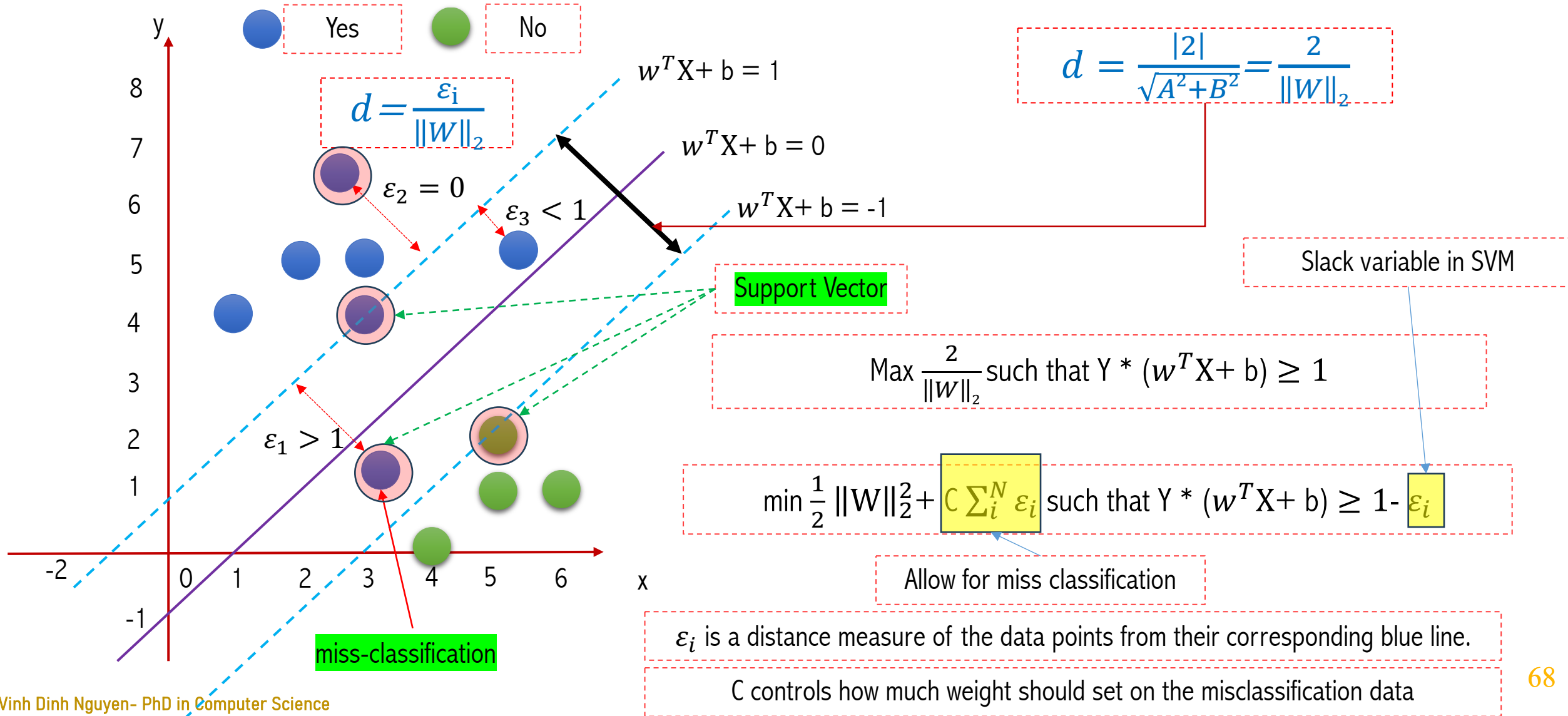


$$\text{Max } \frac{2}{\|w\|_2} \text{ such that } Y * (w^T X + b) \geq 1$$

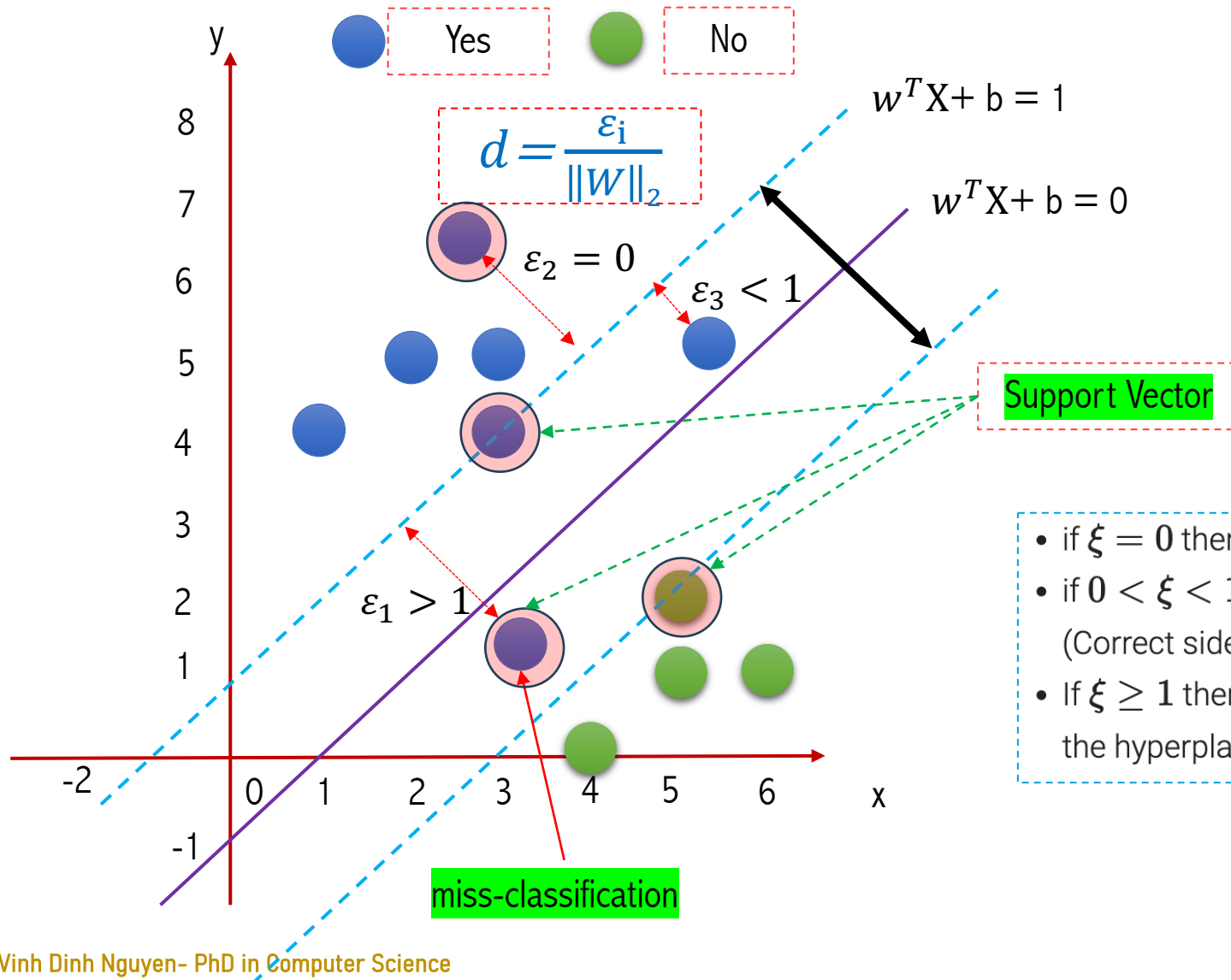
What should we do if a green data point is noise?

Chúng ta nên chấp này như là miss-classification hay là thay đổi hyperplane

Example



Example



- if $\xi = 0$ then the corresponding point ξ is on the margin or further away.
- if $0 < \xi < 1$ then the point ξ is within the margin and classified correctly (Correct side of the hyperplane).
- If $\xi \geq 1$ then the point is misclassified and present at the wrong side of the hyperplane.

Dicussion

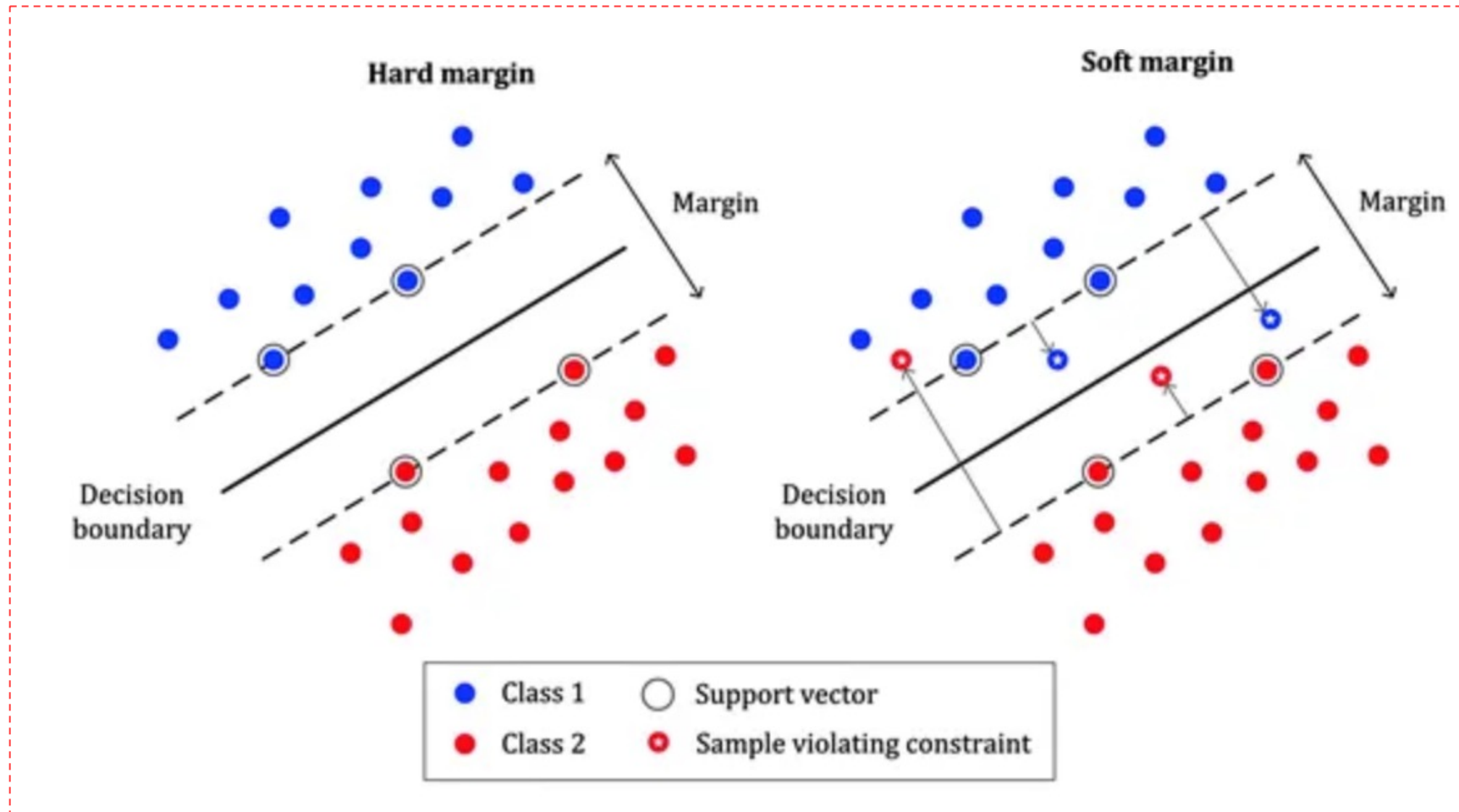
$$\min \frac{1}{2} \|W\|^2 + C \sum_i^N \varepsilon_i \text{ such that } Y * (w^T X + b) \geq 1 - \varepsilon_i \quad \varepsilon_i \geq 0$$



How about the case C is small?
Soft Margin Classifier

How about the case C is large?
Hard Margin Classifier

Hard Margin vs Soft Margin



Further Study

Primal Problem

$$\min \frac{1}{2} \|W\|_2^2 + C \sum_i^N \varepsilon_i \text{ such that } Y * (w^T X + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0$$

$$\text{Objective Function : } \min_{\beta, b, \xi_i} \left\{ \frac{\|\beta^2\|}{2} + C \sum_{i=1}^n \xi_i^2 \right\}$$

$$\text{s.t Linear Constraint : } y_i(\beta^T x_i + b) \geq 1 - \xi_i$$

The **Lagrangian** can be defined as below. Notice we only need one **Lagrange Multiplier** due to the dropped constraint.

$$L = \frac{\|\beta^2\|}{2} + C \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i (y_i(\beta^T x_i + b) - 1 + \xi_i)$$

The HyperParameter C is also called as *Regularization Constant*.

If $k = 1$, then the loss is named as Hinge Loss and if $k = 2$ then its called Quadratic Loss

Another form:

Kernel trick might apply here

$$\phi(x_i)^T \phi(x_j)$$

The Dual Objective can be written as,

$$\max_{\alpha} L_{dual} = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \left(x_i^T x_j - \frac{1}{2C} \delta_{ij} \right)$$

$$\text{s.t constraint : } \alpha_i \geq 0, \forall i \in D, \text{ and } \sum_{i=1}^n \alpha_i y_i = 0$$



However, this way we won't be able use the objective function to solve for **non-linear cases**. Hence, we will find an equivalent problem named **Dual Problem** and solve that using **Lagrange Multipliers**.

