

Data Science and Artificial Intelligence

Machine Learning

Support Vector Machine

Lecture No. 6



By- SIDDHARTH SABHARWAL SIR

Recap of Previous Lecture



Topic

SVM's \Rightarrow advantages

Topic

Kernels

Topic

Soft margin SVM's.

Topic

Topic

Topics to be Covered



Topic

Svm advantage / disadvantage

Topic

Svm for regression

Topic

Questions

Topic

Topic



Nothing will work
unless you do.

Maya Angelou



Soft Margin SVMs

⇒ • For noisy data we need to use soft margin SVM's

- $\min_{w, b, \xi_i} \frac{1}{2} \|w\|^2 + C \sum \xi_i$ What was effect of C

St. $\xi_i \geq 0$

$y_i(w x_i + b) \geq 1 - \xi_i$

→ $\xi_i = 0$ for SV's and other points

→ $0 < \xi_i < 1$

→ $\xi_i > 1$



Soft Margin SVMs

$\xi_i \approx$ distance of point from marginal plane

KKT 1) $\frac{\partial L}{\partial \omega} = 0 \Rightarrow \omega = \sum \lambda_i y_i x_i$

2) $\frac{\partial L}{\partial b} = 0 \Rightarrow \sum \lambda_i y_i = 0$

3) $\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i + \lambda_i = C$

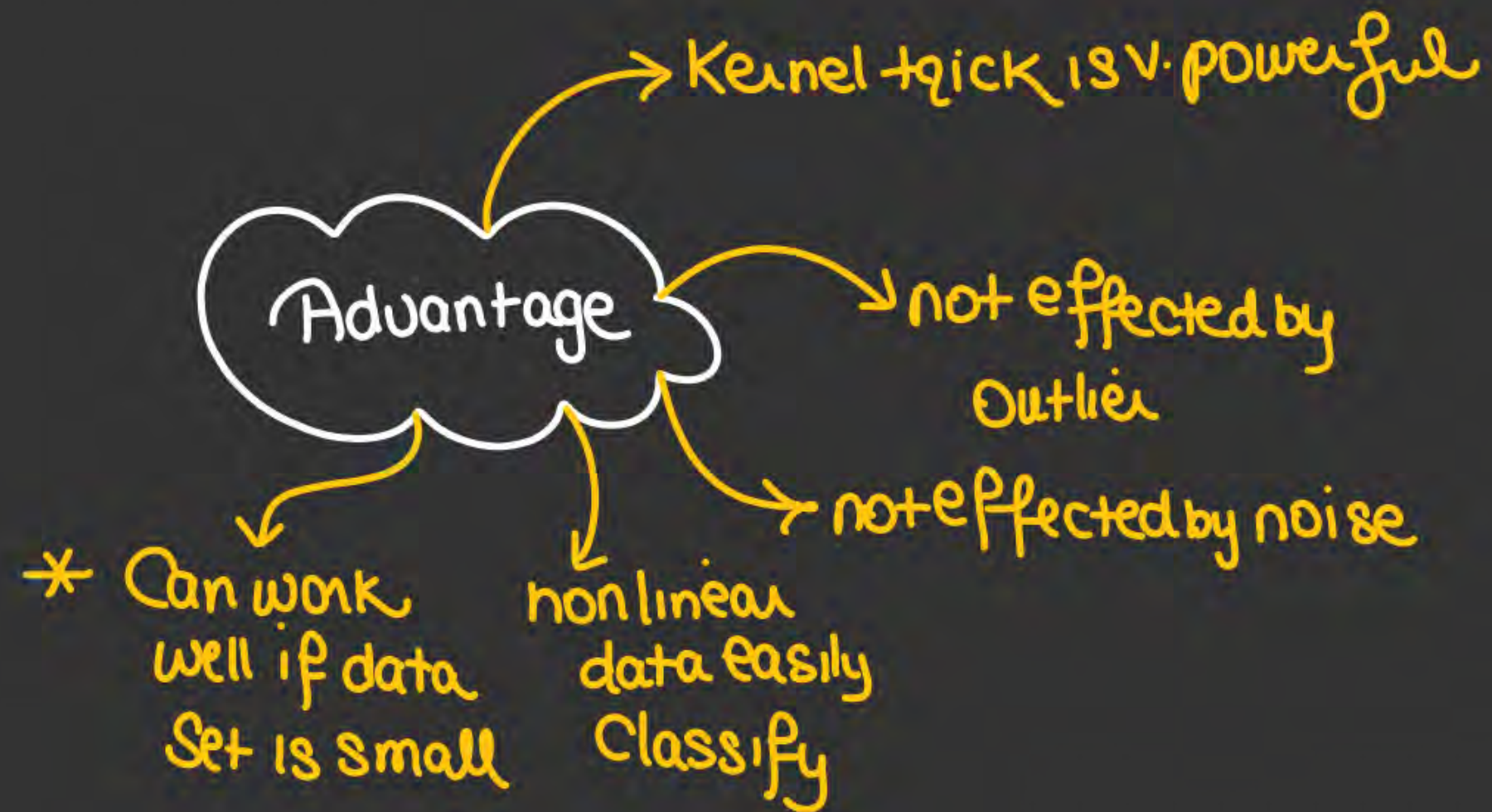
4) $\mu_i \xi_i = 0$

5) $\lambda_i (1 - y_i(\omega x_i + b) - \xi_i) = 0$



Soft Margin SVMs

- Point away from marginal plane $\Rightarrow \begin{matrix} \xi_i = 0 \\ \lambda_i = 0 \\ \mu_i = C \end{matrix}$
- SV's $\lambda_i \neq 0, \xi_i = 0, \mu_i \neq 0, \lambda_i + \mu_i = C$
- Other noisy points $\lambda_i = C, \xi_i \neq 0, \mu_i = 0$



Advantages of SVM

- ✓ Handling high-dimensional data: SVMs are effective in handling high-dimensional data, which is common in many applications such as image and text classification.
- ✓ Handling small datasets: SVMs can perform well with small datasets, as they only require a small number of support vectors to define the boundary.
- ✓ Modeling non-linear decision boundaries: SVMs can model non-linear decision boundaries by using the kernel trick, which maps the data into a higher-dimensional space where the data becomes linearly separable.
- ✓ Robustness to noise: SVMs are robust to noise in the data, as the decision boundary is determined by the support vectors, which are the closest data points to the boundary.

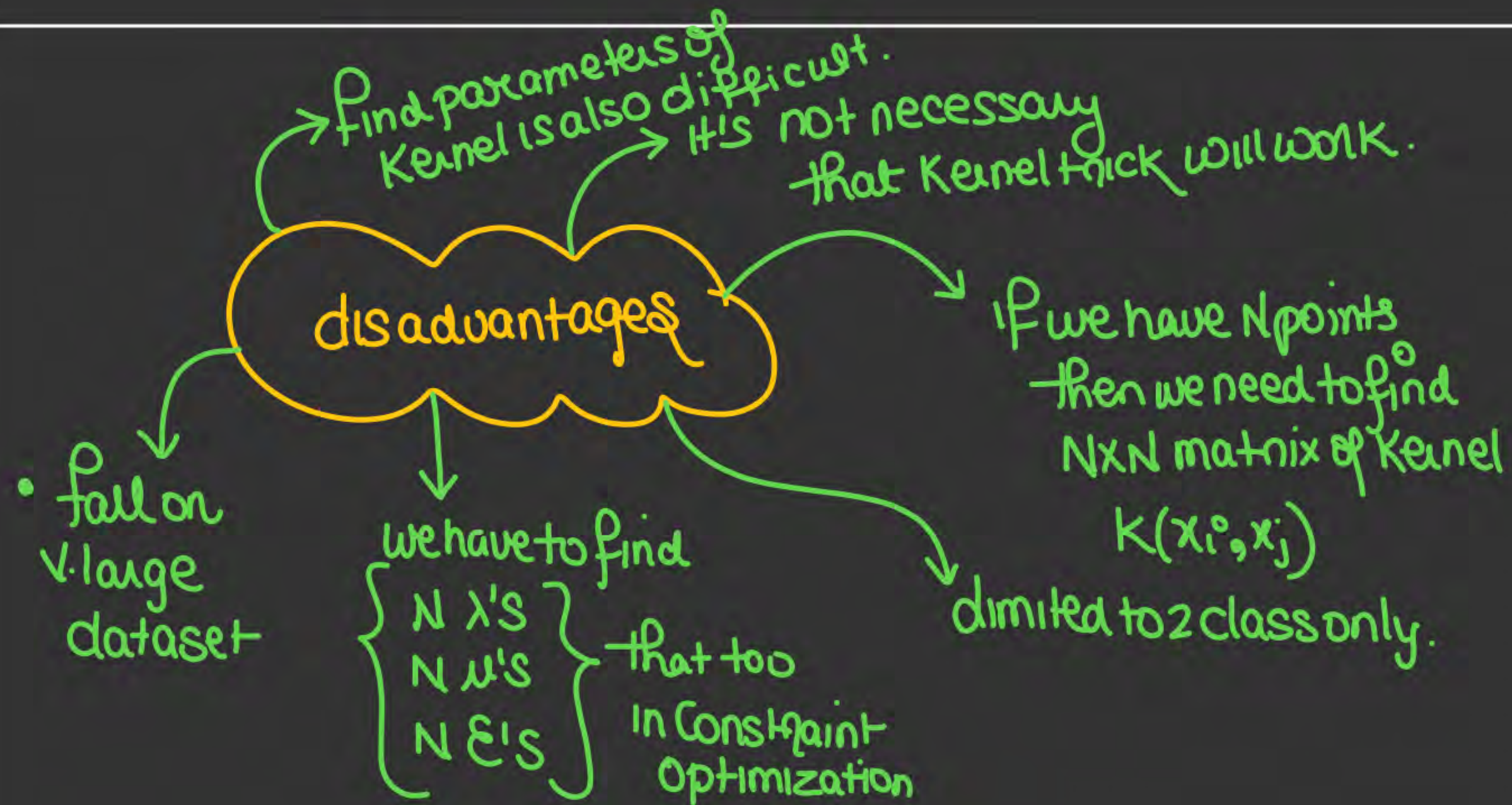
Advantages of SVM

Reduce dimension \Rightarrow use subset of data

SVM \Rightarrow use SV's

in lasso also \Rightarrow Some $\beta_i = 0$ so it uses few dimension.

- ✓ Sparse solution: SVMs have sparse solutions, which means that they only use a subset of the training data to make predictions. This makes the algorithm more efficient and less prone to overfitting.
- ✓ Regularization: SVMs can be regularized, which means that the algorithm can be modified to avoid overfitting.



Disadvantages of SVM

→ Parameters of Kernel γ , 'C'

- ✓ **Computationally expensive:** SVMs can be computationally expensive for large datasets, as the algorithm requires solving a quadratic optimization problem.
- ✓ **Choice of kernel:** The choice of kernel can greatly affect the performance of an SVM, and it can be difficult to determine the best kernel for a given dataset.
- ✓ **Sensitivity to the choice of parameters:** SVMs can be sensitive to the choice of parameters, such as the regularization parameter, and it can be difficult to determine the optimal parameter values for a given dataset.
- ✓ **Memory-intensive:** SVMs can be memory-intensive, as the algorithm requires storing the kernel matrix, which can be large for large datasets.
- ✓ **Limited to two-class problems:** SVMs are primarily used for two-class problems, although multi-class problems can be solved by using one-versus-one or one-versus-all strategies.



Disadvantages of SVM

- Not suitable for large datasets with many features: SVMs can be very slow and can consume a lot of memory when the dataset has many features.
- Not suitable for datasets with missing values: SVMs requires complete datasets, with no missing values, it can not handle missing values.

Hinge loss \Rightarrow

$$\Rightarrow H.L + w(x_i + b) = 1$$

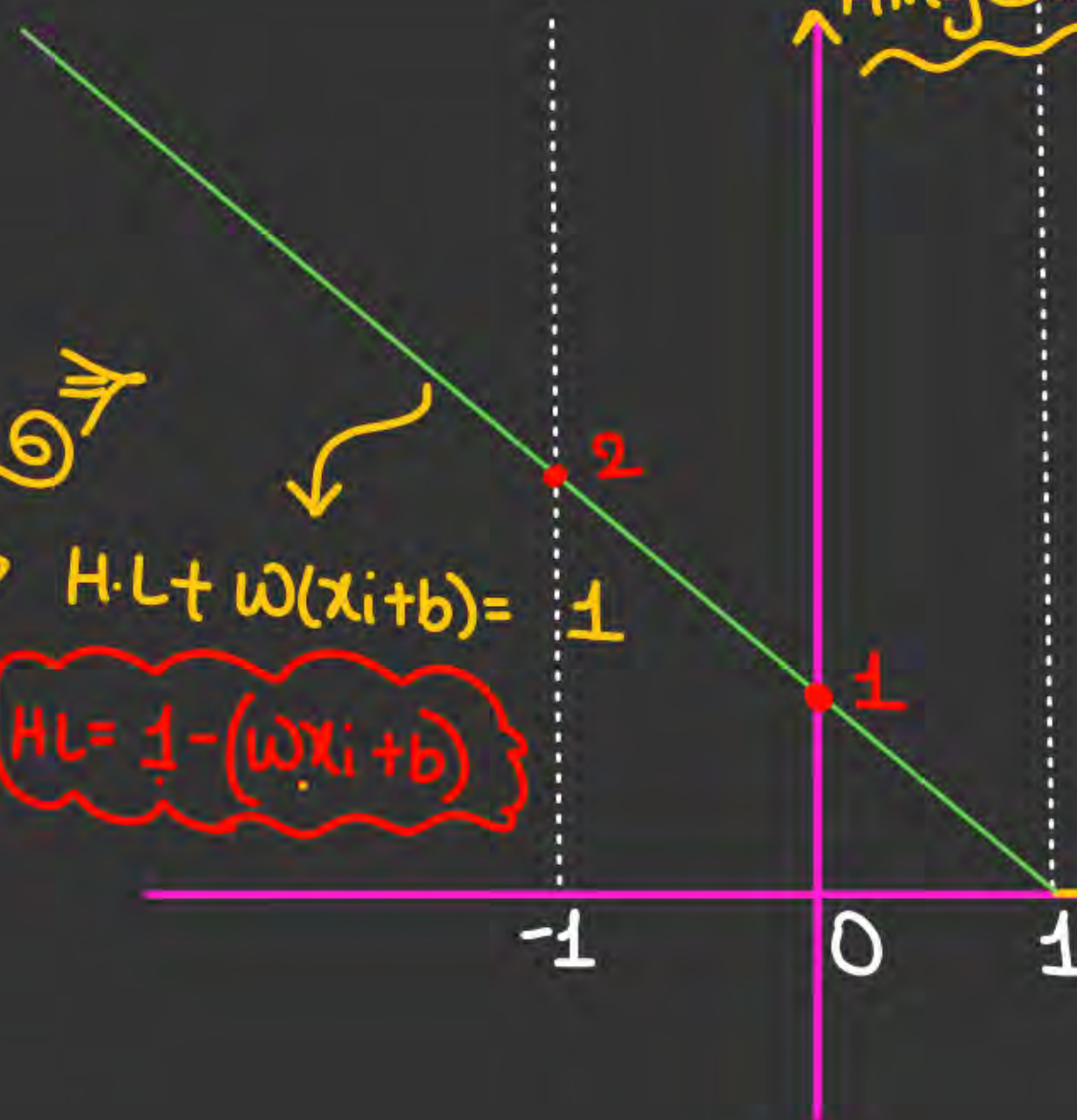
$$\Rightarrow \boxed{HL = 1 - (wx_i + b)}$$

Hinge loss

Point x_i of
Class 1

• $w x_i + b > 1 \Rightarrow$ no error

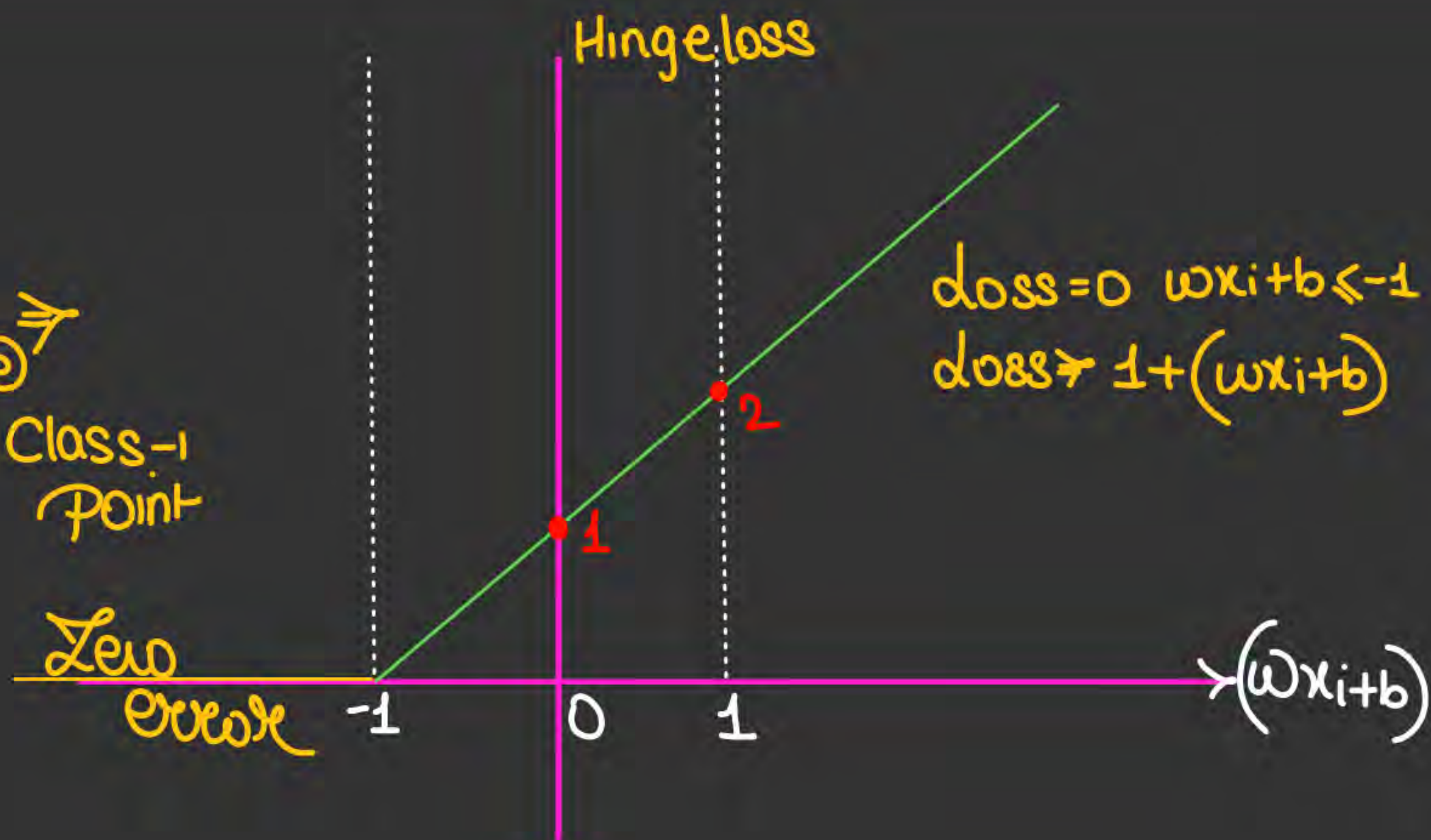
Zero error $\rightarrow (wx_i + b)$



Hinge loss \Rightarrow

Class-1
Point

Zero
error

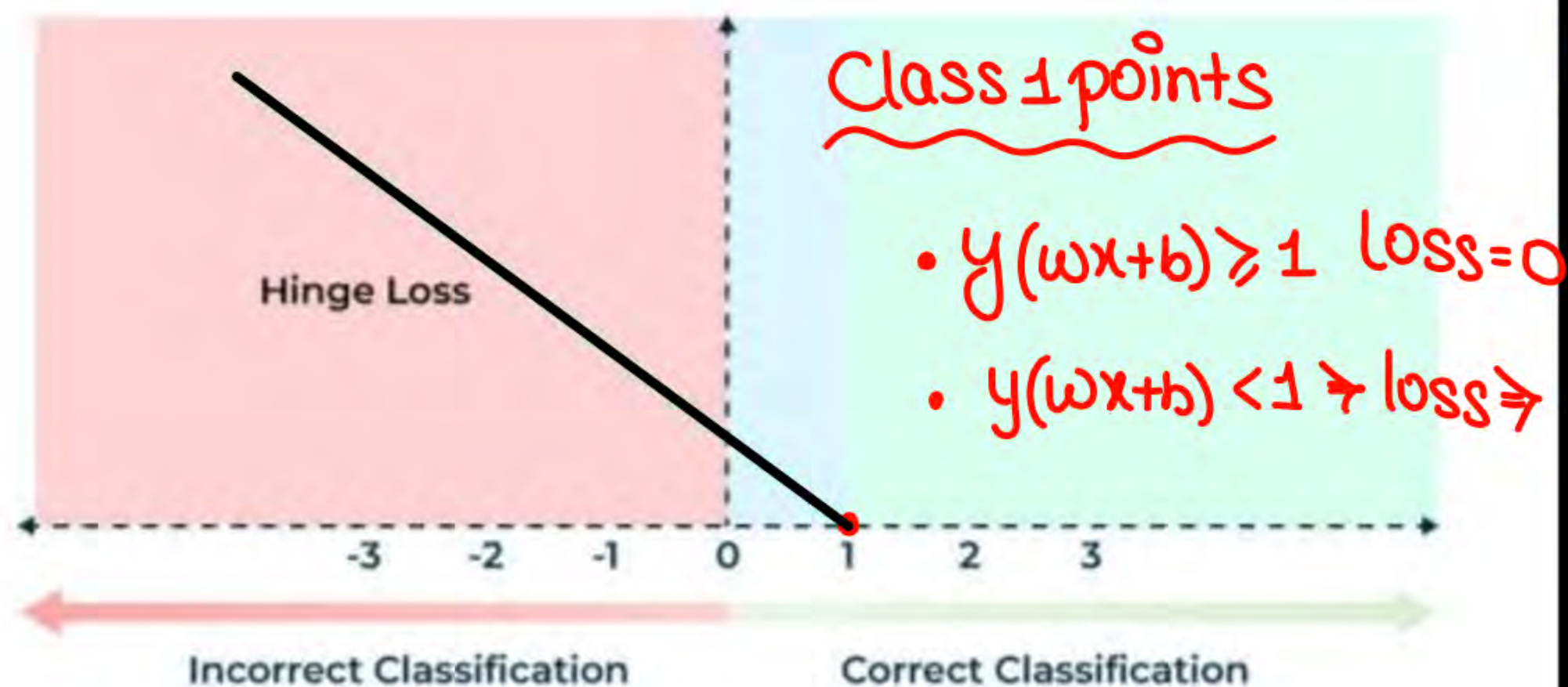




Hinge Loss in SVMs

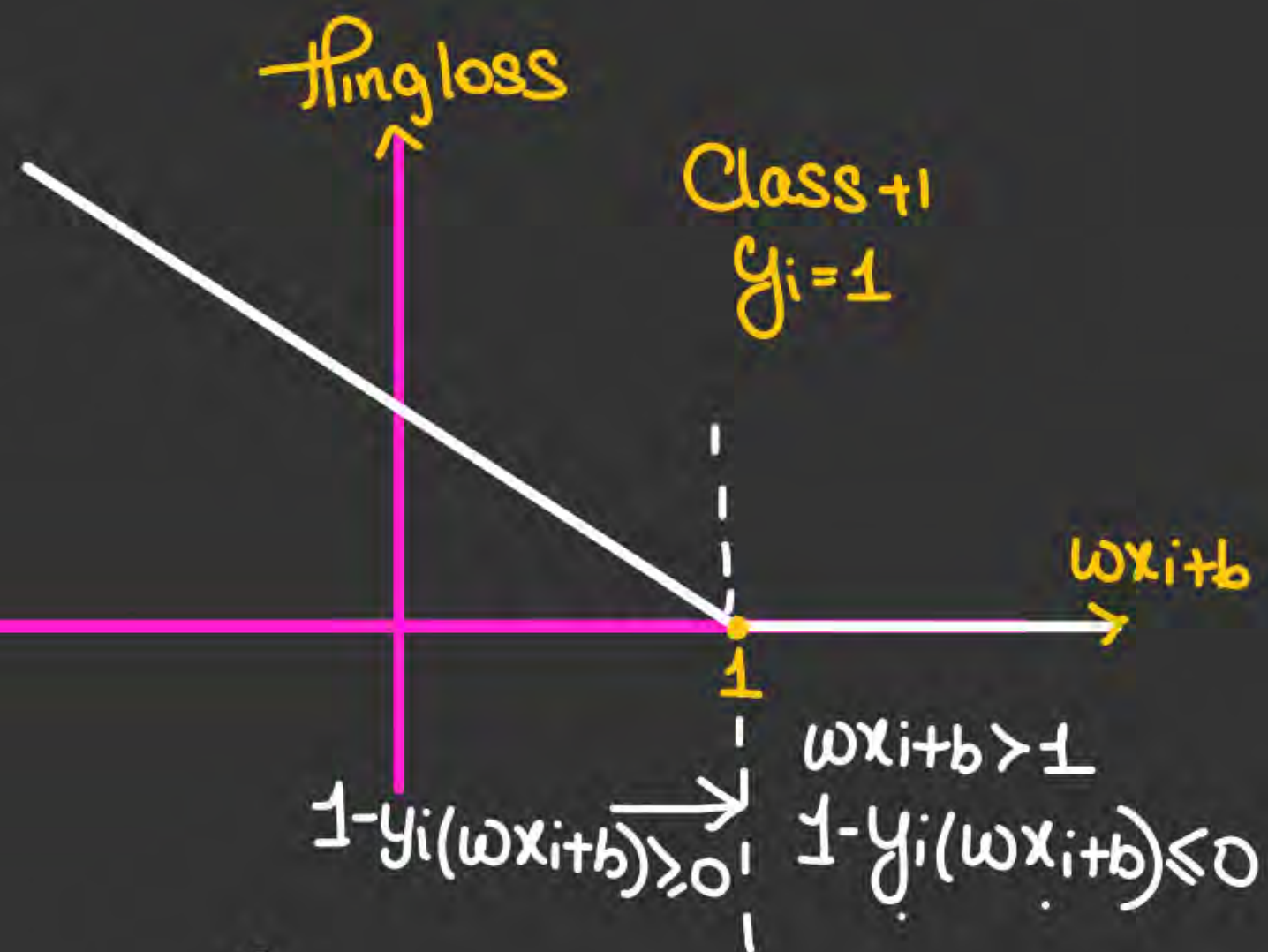
Mathematically, Hinge loss for a data point can be represented as :

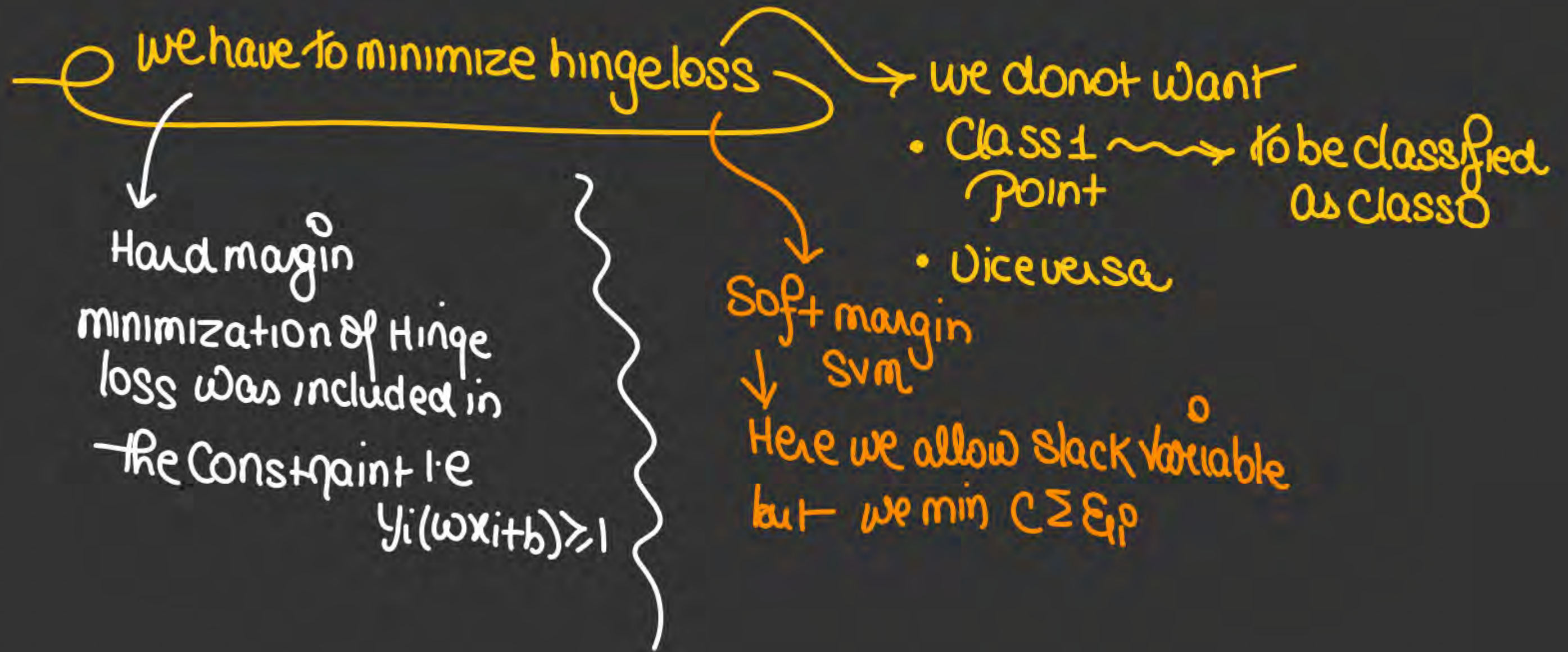
$$L(y, f(x)) = \max(0, 1 - y * f(x))$$



Combine our
general eq,
 $d = 1 - y_i(\omega x_i + b)$
Class 1 $y_i = 1$
 $d = 1 - (\omega x_i + b)$
Class -1 $y_i = -1$
 $d = 1 + (\omega x_i + b)$

So the eq of Hinge loss
 $\Rightarrow \max(0, 1 - y_i(\omega x_i + b))$







Hinge Loss in SVMs

- ✓ If we look at the mathematical formulation the *hinge loss is effectively present in the constraints* of a hard margin. This ensures that the decision boundary (the hyperplane) is positioned in such a way that it maximizes the margin without allowing any data points to be within or on the wrong side of the margin.
- ✓ Here the hinge loss component, is part of the objective function itself through slack variable.

↪ in soft margin SVMs



Practice

The soft margin SVM is more preferred than the hard-margin svm when:

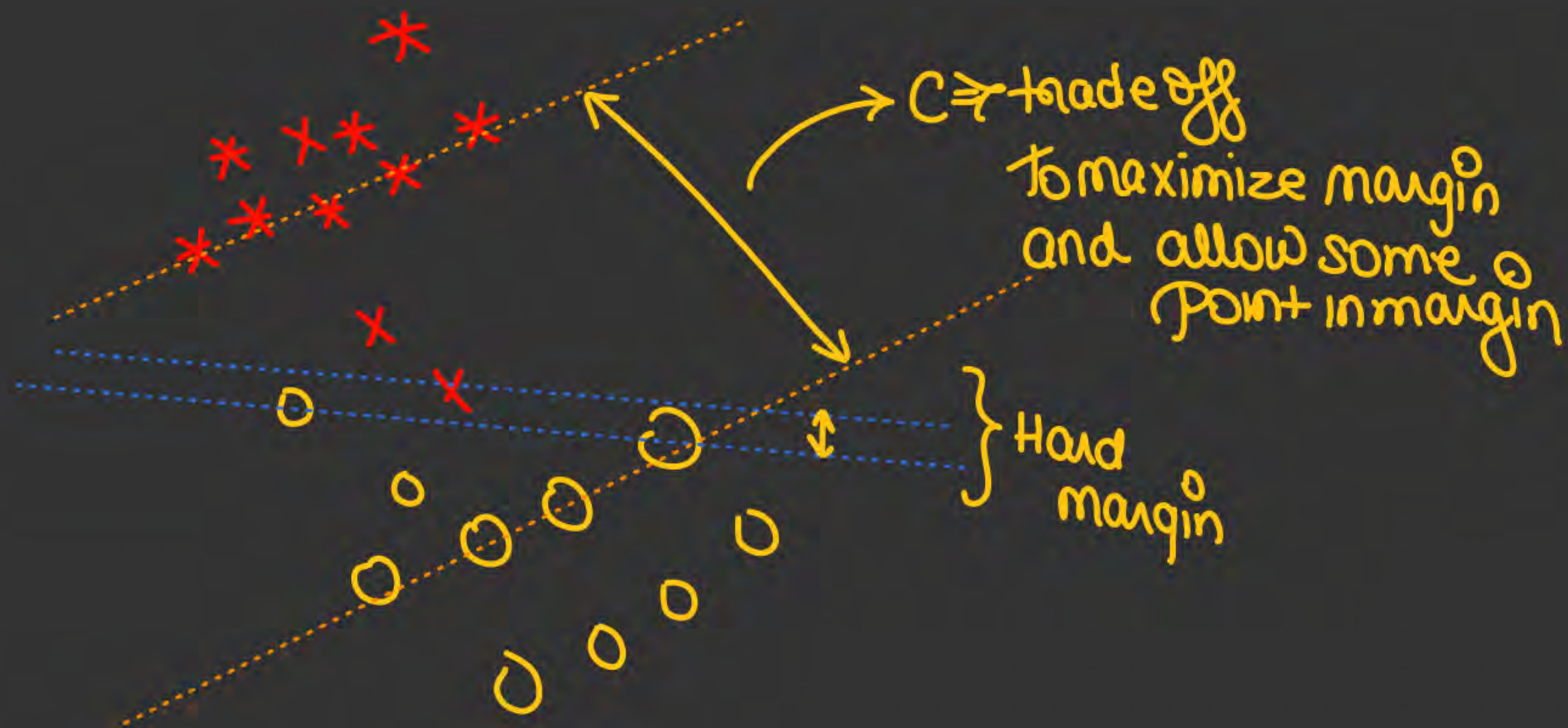
1. The data is linearly separable
- ✓ 2. The data is noisy and contains overlapping point



Practice

In the linearly non-separable case, what effect does the C parameter have on the SVM mode.

- a. it determines how many data points lie within the margin
- b. it is a count of the number of data points which do not lie on their respective side of the hyperplane
- ☒ c. it allows us to trade-off the number of misclassified points in the training data and the size of the margin
- d. it counts the support vectors



Practice

SVM is a supervised Machine Learning can be used for Options :

☐ Regression

☐ Classification

☒ both a or b

☐ None of These

Practice

Closest Point to the hyperplane are support vectors

→ Classifier

☐ True

Hard margin

☐ Unpredictable

☐ False

soft+margin

☐ None of these

Practice

In SVM, the dimension of the hyperplane depends upon which one? ⑥

Hyperplane $wx+b=0$.

- ☒ the number of features
- ☐ the number of samples
- ☐ the number of target variables
- ☐ All of the above



$$y_i(\omega^T x_i + b) \text{ always } \geq 0$$

Practice

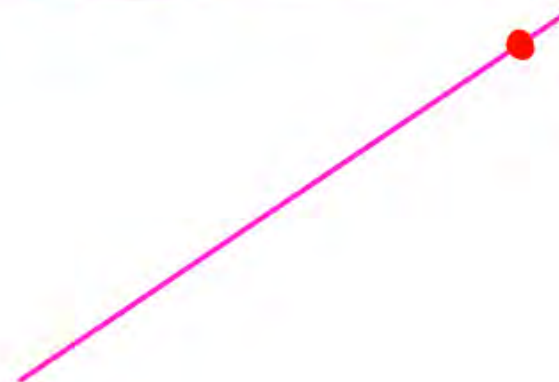
Choose the correct option regarding classification using SVM for two classes

Statement i : While designing an SVM for two classes, the equation $y_i(\omega^T x_i + b) \geq 1$ is used to choose the separating plane using the training vectors.

~~Statement ii~~ : During inference, for an unknown vector x_j , if $y_j(\omega^T x_j + b) \geq 0$, then the vector can be assigned class 1.

Statement iii : During inference, for an unknown vector x_j , if $(\omega^T x_j + b) > 0$, then the vector can be assigned class 1.

- a. Only Statement i is true
- ☒ b. Both Statements i and iii are true
- c. Both Statements i and ii are true
- d. Both Statements ii and iii are true



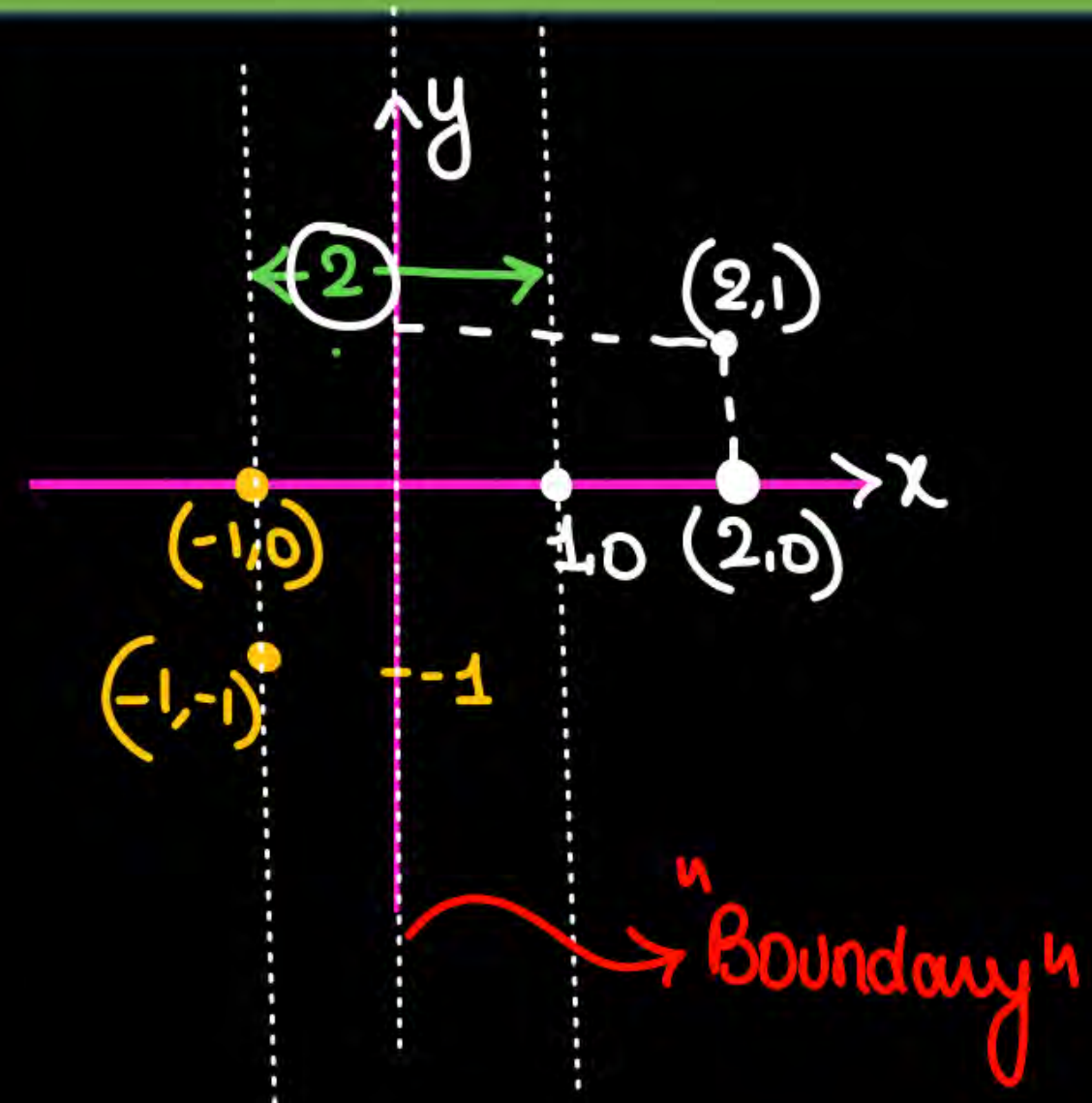
Practice

QUESTION 7:

Suppose we have the below set of points with their respective classes as shown in the table. Answer the following question based on the table.

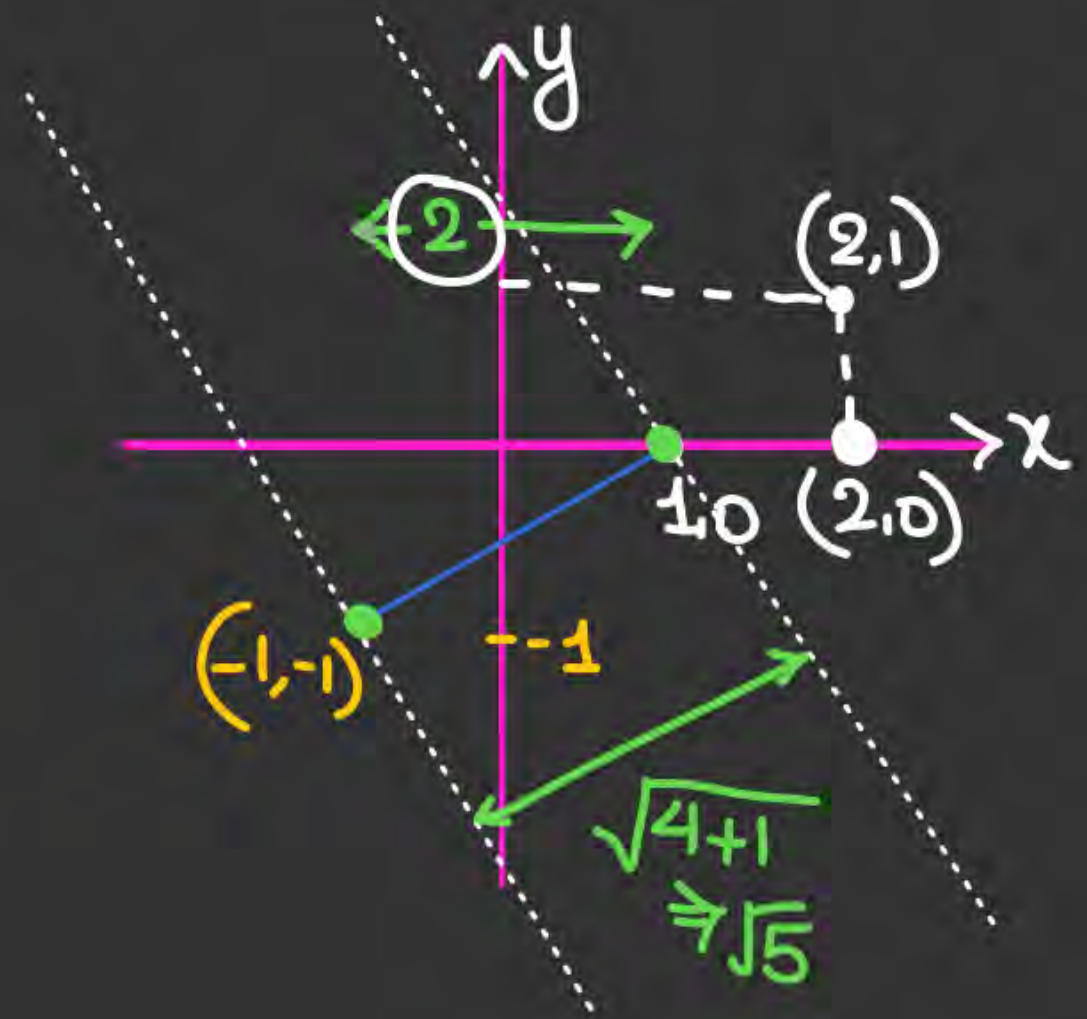
X	Y	Class Label
1	0	+1
-1	0	-1
2	1	+1
-1	-1	-1
2	0	+1

• To find margin simply find the euclidean distance b/w sv's



What will happen to maximum margin if we remove the point $(-1,0)$ from the training set?

- a. Maximum margin will decrease
- ☒ b. Maximum margin will increase
- c. Maximum margin will remain same
- d. Can not decide





Practice

Suppose we have the below set of points with their respective classes as shown in the table. Answer the following question based on the table.

X	Y	Class Label
1	0	+1
-1	0	-1
2	1	+1
-1	-1	-1
2	0	+1

What can be a possible decision boundary of the SVM for the given points?

- a. $y = 0$
- ☒ b. $x = 0$
- c. $x = y$
- d. $x + y = 1$



Practice

Suppose we have the below set of points with their respective classes as shown in the table. Answer the following question based on the table.

X	Y	Class Label
1	0	+1
-1	0	-1
2	1	+1
-1	-1	-1
2	0	+1

done
 $x=0$
 $x>0 \rightarrow +1$
 $x<0 \rightarrow -1$

Find the decision boundary of the SVM trained on these points and choose which of the following statements are true based on the decision boundary.

- ☒ The point $(-1, -2)$ is classified as -1
- ☒ The point $(1, -2)$ is classified as -1
- ☒ The point $(-1, -2)$ is classified as +1
- ☒ The point $(1, -2)$ is classified as +1



Practice

Which one of the following is a valid representation of hinge loss (of margin = 1) for a two-class problem?

y = class label (+1 or -1).

p = predicted (not normalized to denote any probability) value for a class.?

- a. $L(y, p) = \max(0, 1 - yp)$
- b. $L(y, p) = \min(0, 1 - yp)$
- c. $L(y, p) = \max(0, 1 + yp)$
- d. None of the above

we will see 6

(MCQ)



#Q. Consider the problem of finding an optimal hyperplane for non-separable patterns, we introduce a new set of variables, $\{\xi_i\}_{i=1}^N$ into the definition of the 2 points separating hyperplane as ~~$d_i(w^T x_i + b) > 1 - \xi_i$~~ . Choose the correct statements from the options given below. $y_i(w^T x_i + b) > 1 - \xi_i$

- ☒ **A** ~~The slack variable ξ_i can take both positive and negative values.~~ (b, c, d)
- ☒ **B** For $0 < \xi_i \leq 1$ the data point falls inside the region of separation, but on the correct side of the decision surface.
- ☒ **C** For $\xi_i > 1$ the data point falls on the wrong side of the separating hyperplane.
- ☒ **D** For support vectors ξ_i will be always zero.

#Q. For the nonseparable case, we minimize the cost function defined as

$$L = \frac{1}{2} \omega^T \omega + C \sum_{i=1}^N \xi_i$$

(True/False) The optimal value of C is obtained by minimizing the cost function with respect to C.

False

→ (we find C by Cross validation)

A True

B False

#Q. In continuation with question 2, consider the following statements:

- ☒ (a) The parameter C can be chosen using cross validation approach.
 - ☒ (b) When C is assigned a small value, the training samples are considered to be noisy, and less emphasis should therefore be placed on it.
 - ☒ (c) The optimization problem for linearly separable patterns can be considered as a special case of optimization problem for nonseparable patterns, by setting $\xi_i = 0$ for points all i .
 - ☒ (d) When C is assigned a large value, the implication is that the designer of the SVM has high confidence in the quality of the training samples.
- Which of the above statements are correct?

☐ **A** Only a and c

☐ **B** Only b and d

☐ **C** Only a, b and c

☒ **d** ✓

☐ **D** a, b, c and d

#Q. If we are using a kernel function k to evaluate the inner products in a feature space with feature map ϕ , the associated Gram matrix G has entries $G_{ij} = k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$. Then the kernel matrix G is

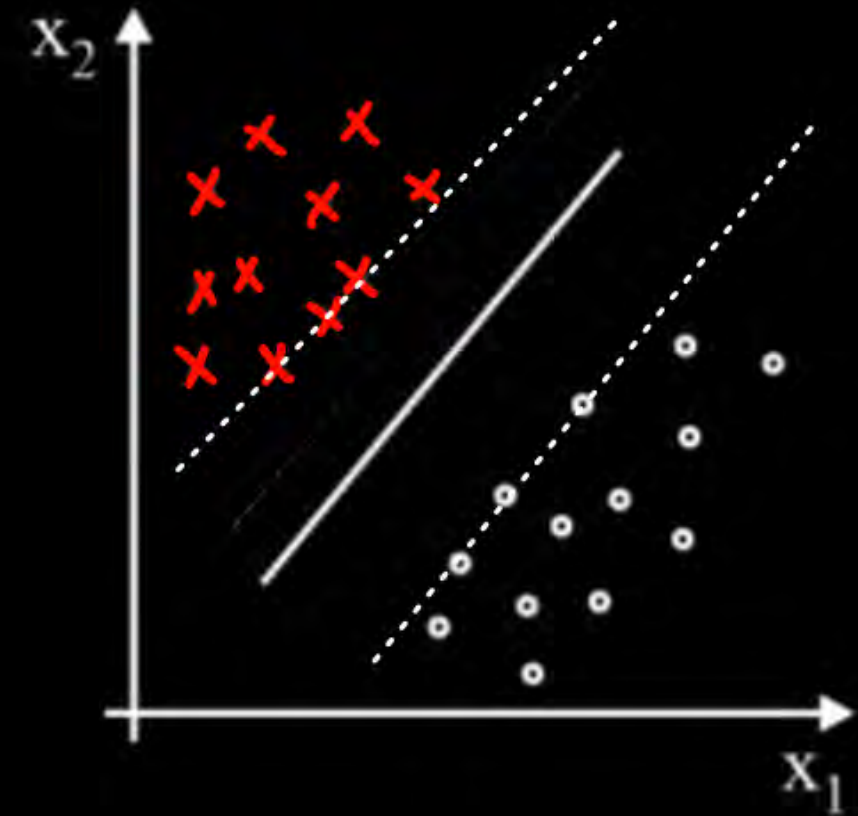
- A** Positive definite.
- B** Negative definite.
- C** Positive semi-definite.
- D** Negative semi-definite.

#Q. In the linearly non-separable case, what effect does the C parameter have on the SVM model?

(d) done

- A** it determines the count of support vectors
- B** it is a count of the number of data points which do not lie on their respective side of the hyperplane
- C** it determines how many data points lie within the margin
- D** it allows us to trade-off the number of misclassified points in the training data and the size of the margin

#Q. What is the leave-one-out cross-validation error estimate for maximum margin separation in the following figure?



if any one point go into
validation
and other for training
→ Classifier will not
change

A 0 ✓

B 2

C 3

D 6

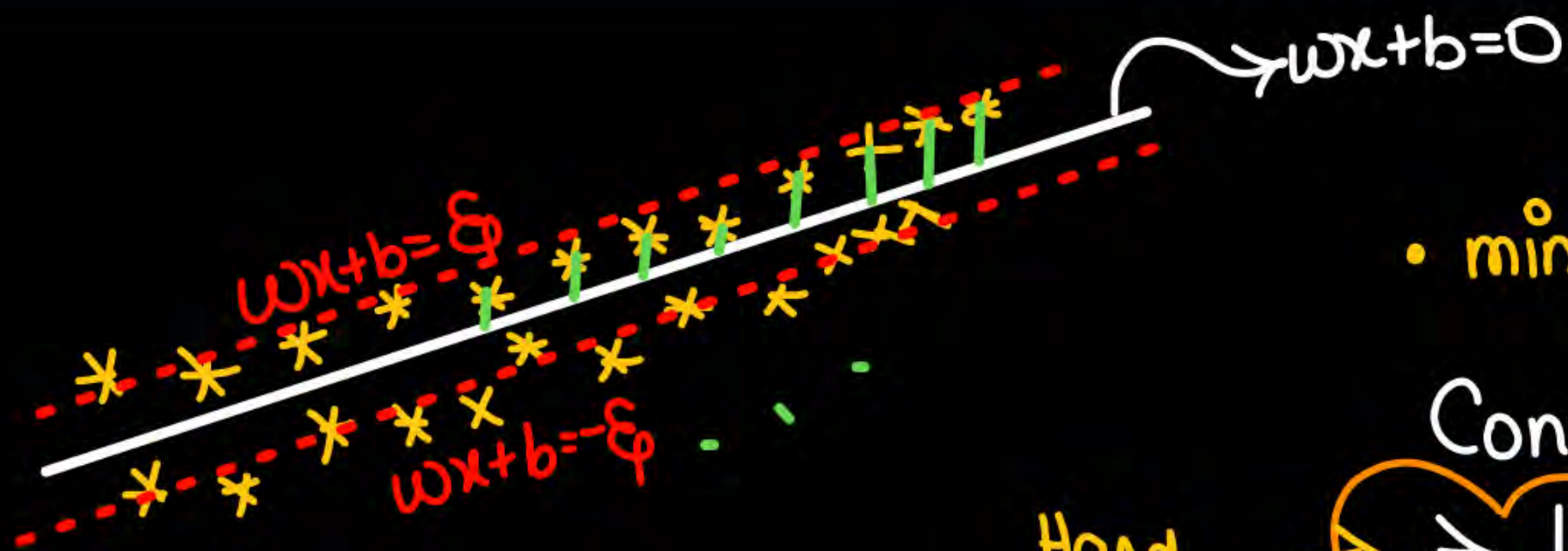
SVM for regression

- generally we use SVM for classification but we can use SVRegressor also.

Brief intro



SVM for regression



- $\min \frac{1}{2} \|w\|^2$

Constraint \Rightarrow

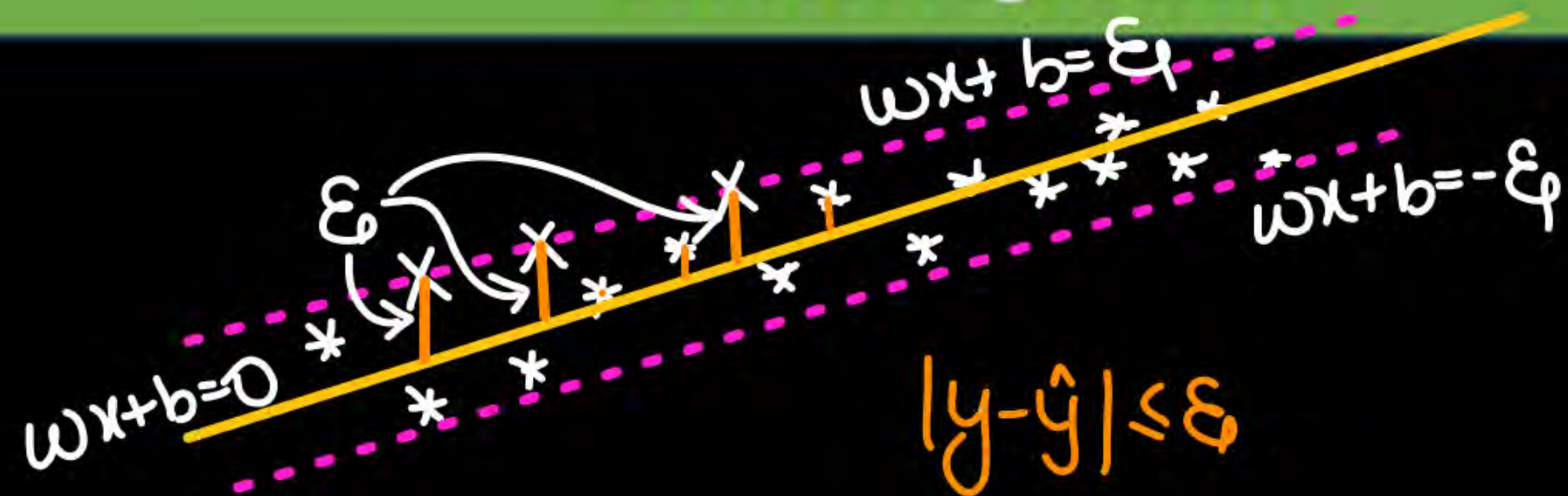
Hard margin \Rightarrow

$$\Rightarrow |y - \hat{y}| \leq \xi$$

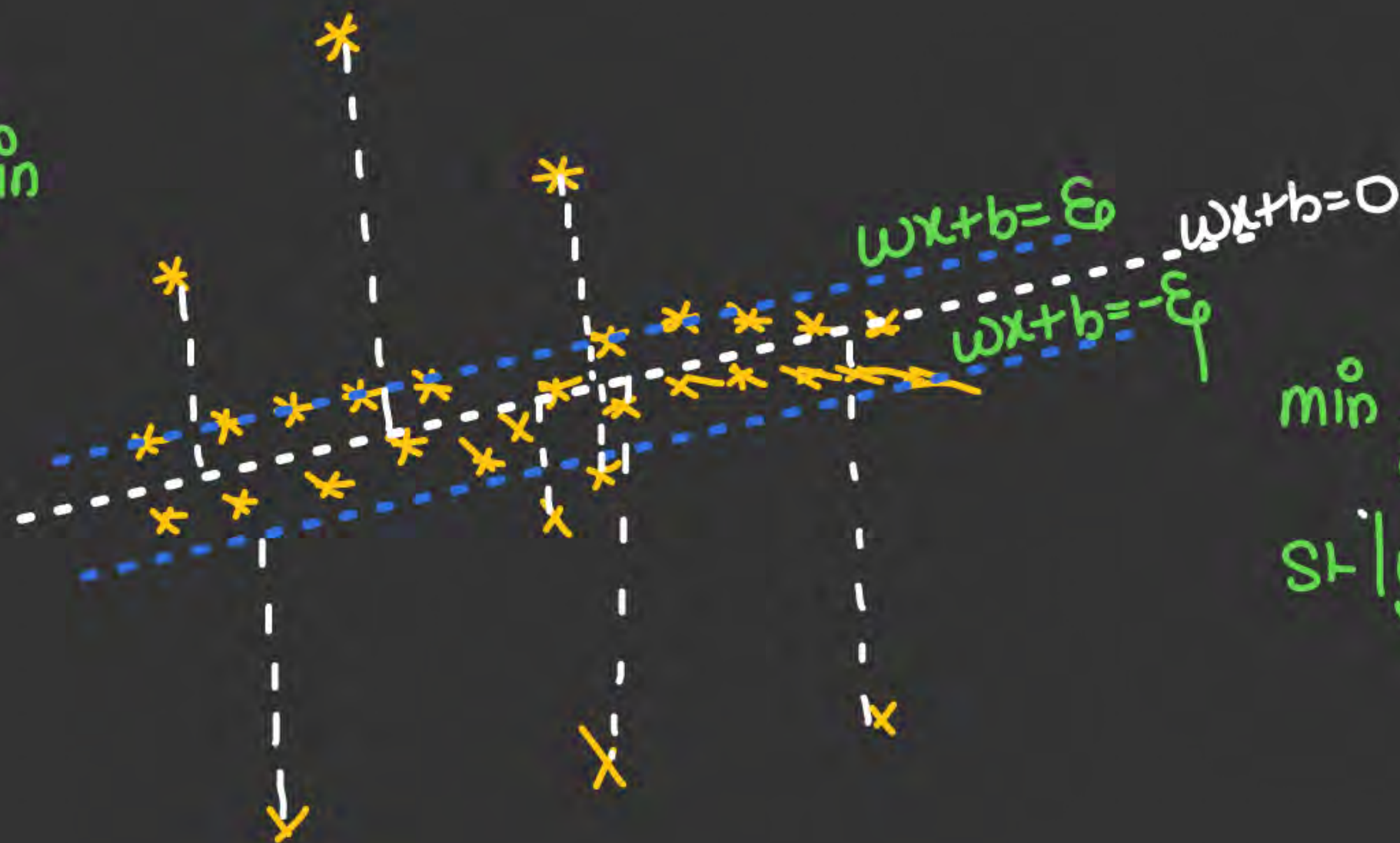
\rightarrow we want that all point shd be within the marginal plane



SVM for regression



Soft margin



$$\min \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^N \alpha_i$$
$$\text{s.t. } |y - \hat{y}| \leq \xi_0 + \alpha_i$$

THANK - YOU