# Day 13: Linear Algebra Review and SVMs

# Linear Algebra

$$\text{Matrix } A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

<span style="color:red">2×3 matrix</span>

Matrices can represent many things...

$$\begin{bmatrix} \ddots & \ddots \\ \ddots & \ddots \end{bmatrix}$$

images

| | Col1 | Col2 | Col3 |
|---|---|---|---|
| dp1 | 34 | 3 | 9 |
| dp2 | 48 | -3 | 9 |
| ... | | | |

tabular data

...

$$A^T = \text{"transpose of matrix A"} = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$$

$A^{-1} =$ "inverse of matrix $A$"

$$AA^{-1} = A^{-1}A = I = \text{"identity matrix"}$$

Identity matrix: 0's everywhere except 1's in the diagonal

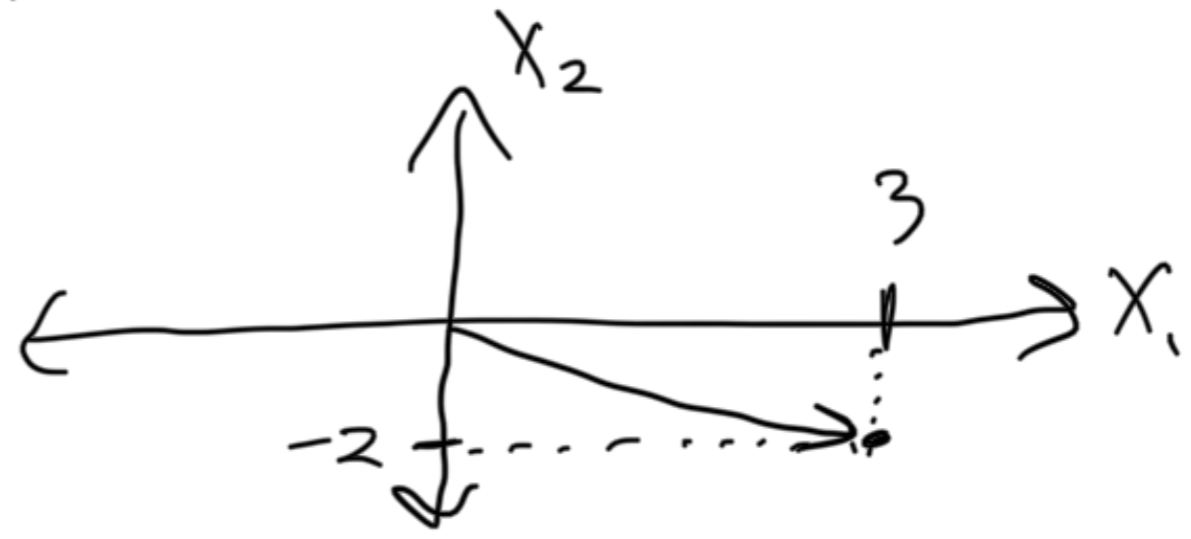$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$AI = A$

$IA = A$

Vectors are $n \times 1$ matrices

$$\begin{bmatrix} a \\ b \end{bmatrix}, \quad \begin{bmatrix} a \\ b \\ c \\ \cdots \\ 2 \end{bmatrix}, \quad \text{etc.}$$
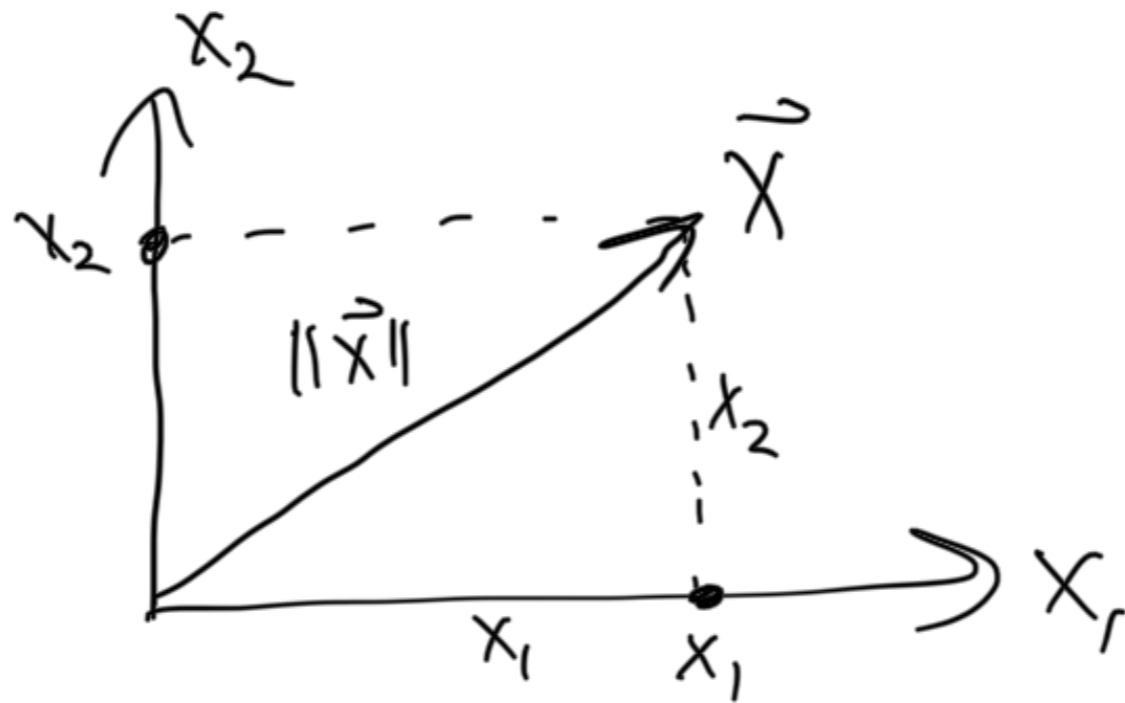
Represents a point in space:

Represent

$$\begin{bmatrix} 3 \\ -2 \end{bmatrix}$$



## Norm of a vector: magnitude of a vector

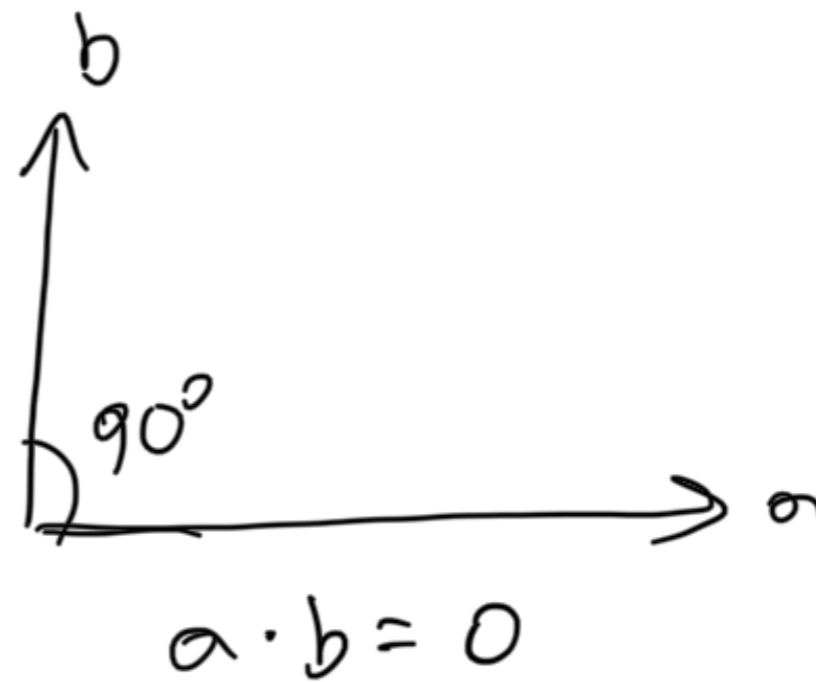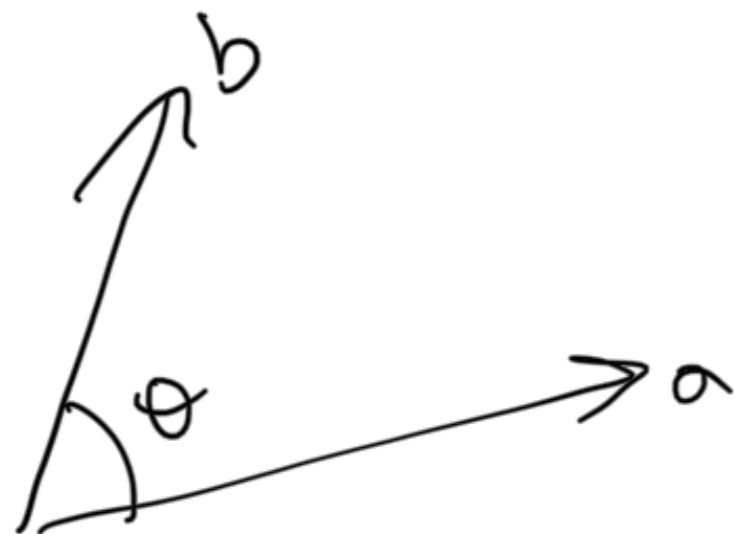$$\| X \|_2 = \sqrt{X_1^2 + X_2^2}$$



## Dot Product:

$$\begin{bmatrix} 4 & 3 & 2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 3 & -7 \end{bmatrix}$$

$$= 4 \cdot 1 + 3 \cdot 3 + 2 \cdot -7 = -1$$

$$a \cdot b = \|a\| \|b\| \cos \theta$$



$$a \cdot b = 0$$

## Matrix Multiplication

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix}$$

$$= \begin{bmatrix} a_{11} \cdot b_{11} + a_{12} \cdot b_{21} & \oslash & \oslash \\ \oslash & \oslash & \oslash \\ \oslash & \oslash & a_{31} b_{13} + a_{32} b_{23} \end{bmatrix}$$

(dot products between rows of $A$
and columns of $B$)

## Linear Regression as Linear Algebra

$$y = w_1 x_1 + w_2 x_2 + \ldots w_n x_n + b$$

$$= \vec{w}^T \vec{x} + b$$

For a dataset with 3 points:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \approx \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = X\theta$$

$$X\theta = \begin{bmatrix} \theta_0 + \theta_1 x_1 \\ \theta_0 + \theta_1 x_2 \\ \theta_0 + \theta_1 x_3 \end{bmatrix}$$

Minimize $(y - X\theta)^T (y - X\theta)$

$$= \ldots = y^T y - 2\theta^T X^T y + \theta^T X^T X \theta$$

So, to find the minimum;

$$\frac{\partial MSE}{\partial \theta} = -2\theta^{\top} x^{\top} y + 2x^{\top}x = 0$$

Solve for $\theta$:

$$\theta = \left(x^{\top}x\right)^{-1} x^{\top} y$$
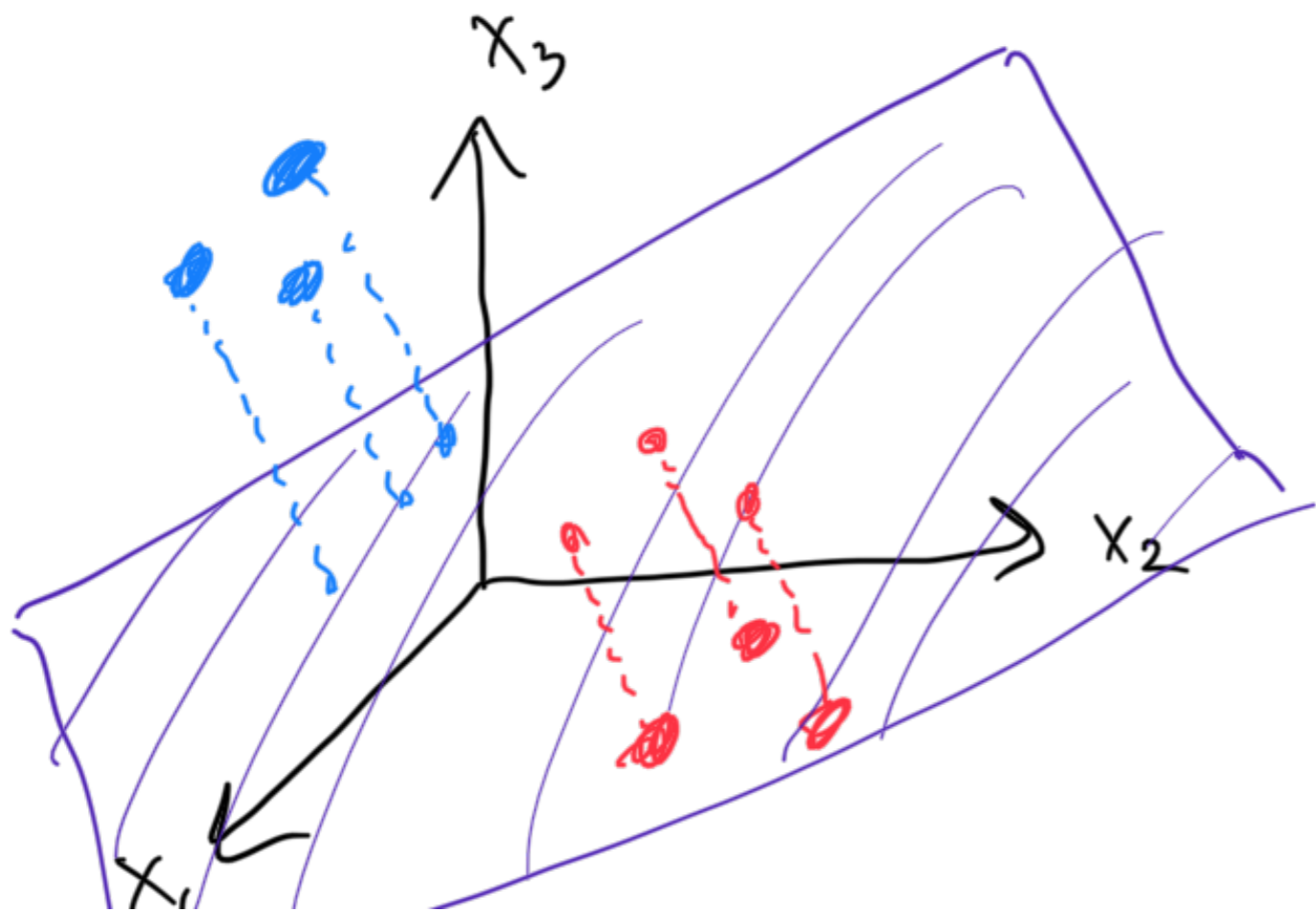
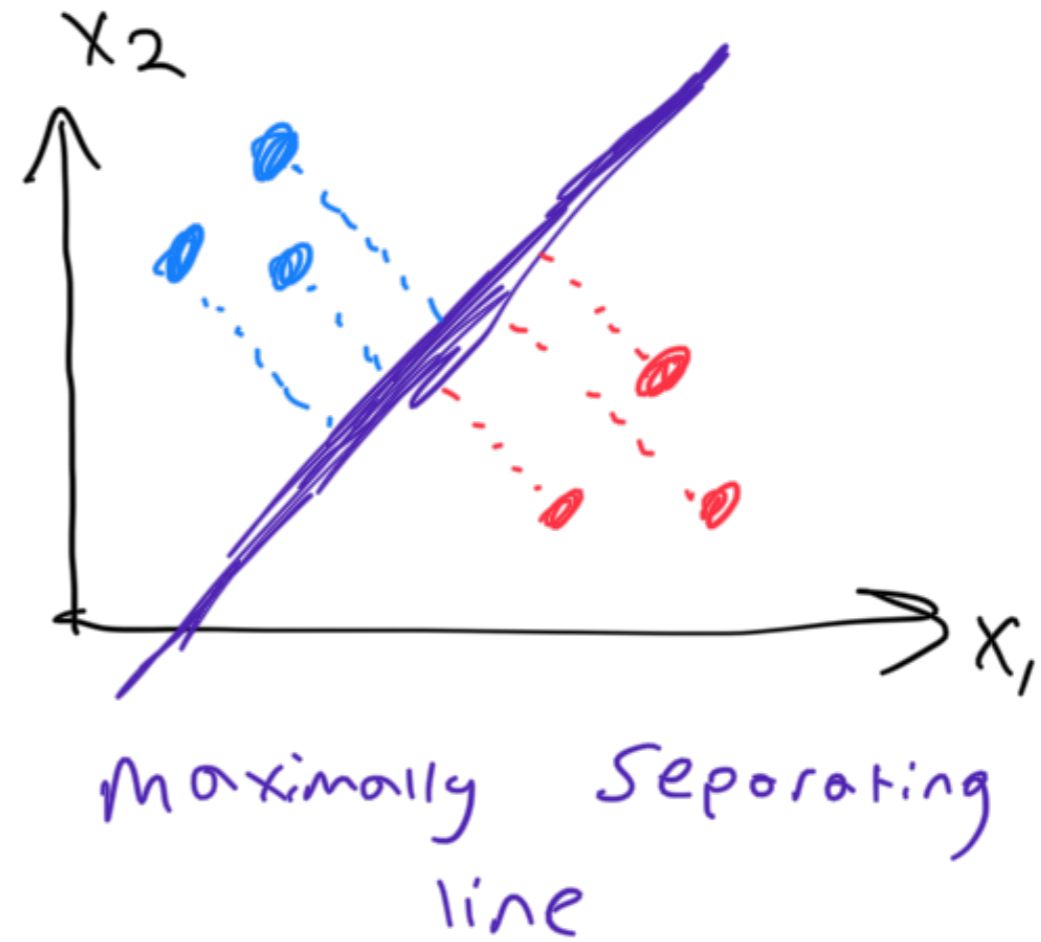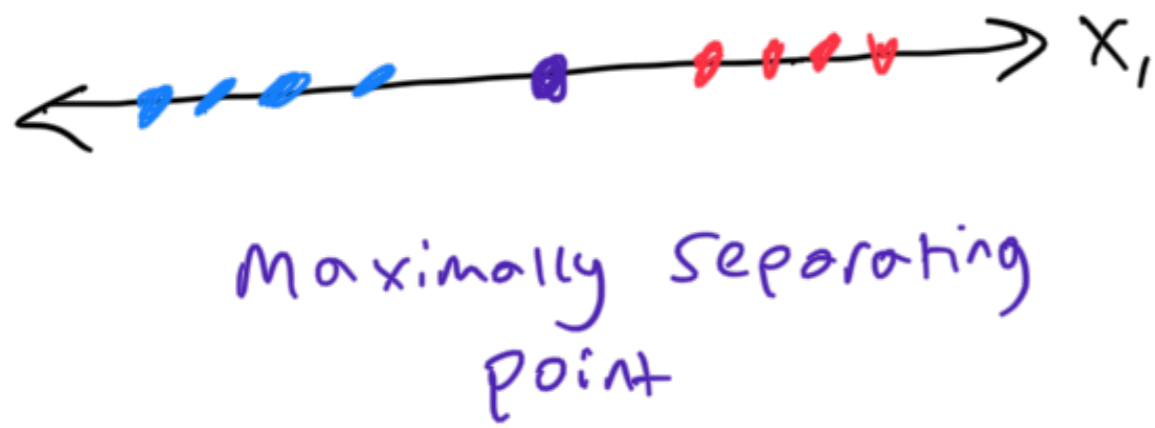<span style="color:red">"Normal Equation"</span>
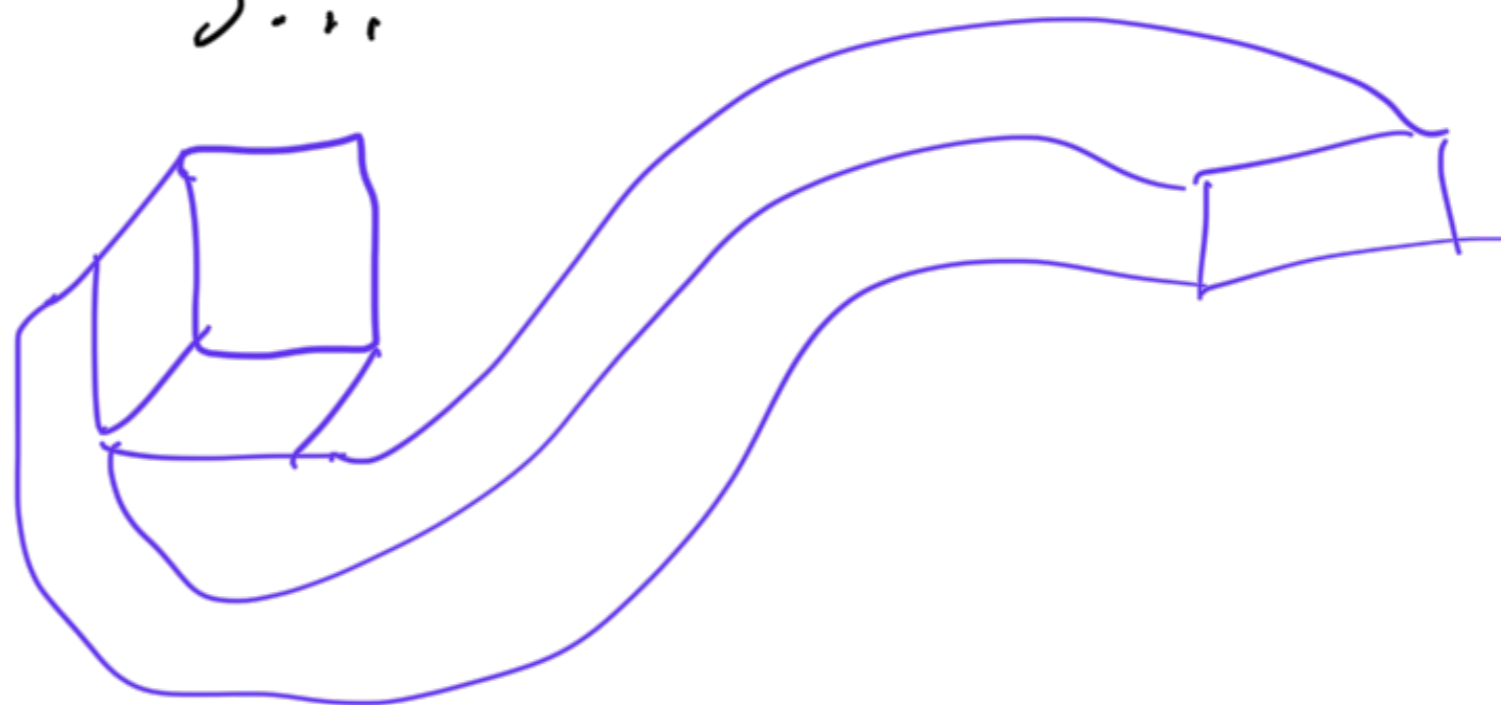
<span style="color:red">analytical solution to linear regression</span>

# Support Vector Machines

# Basic idea: find the maximally Separating hyperplane



Maximally Separating point
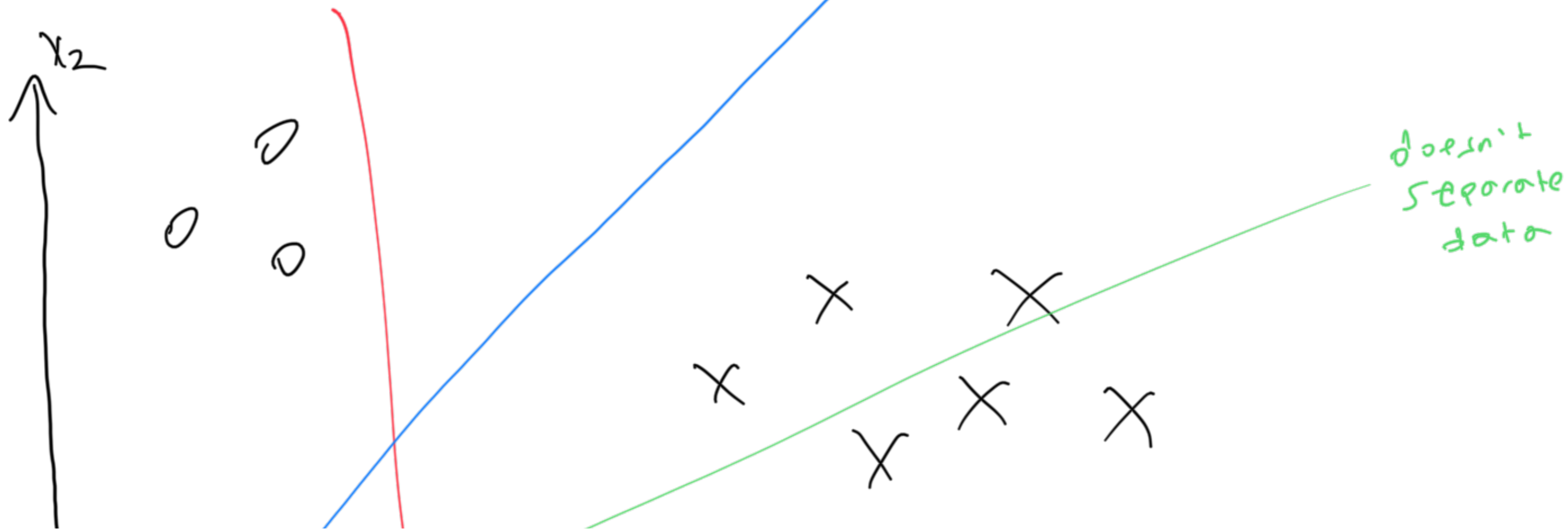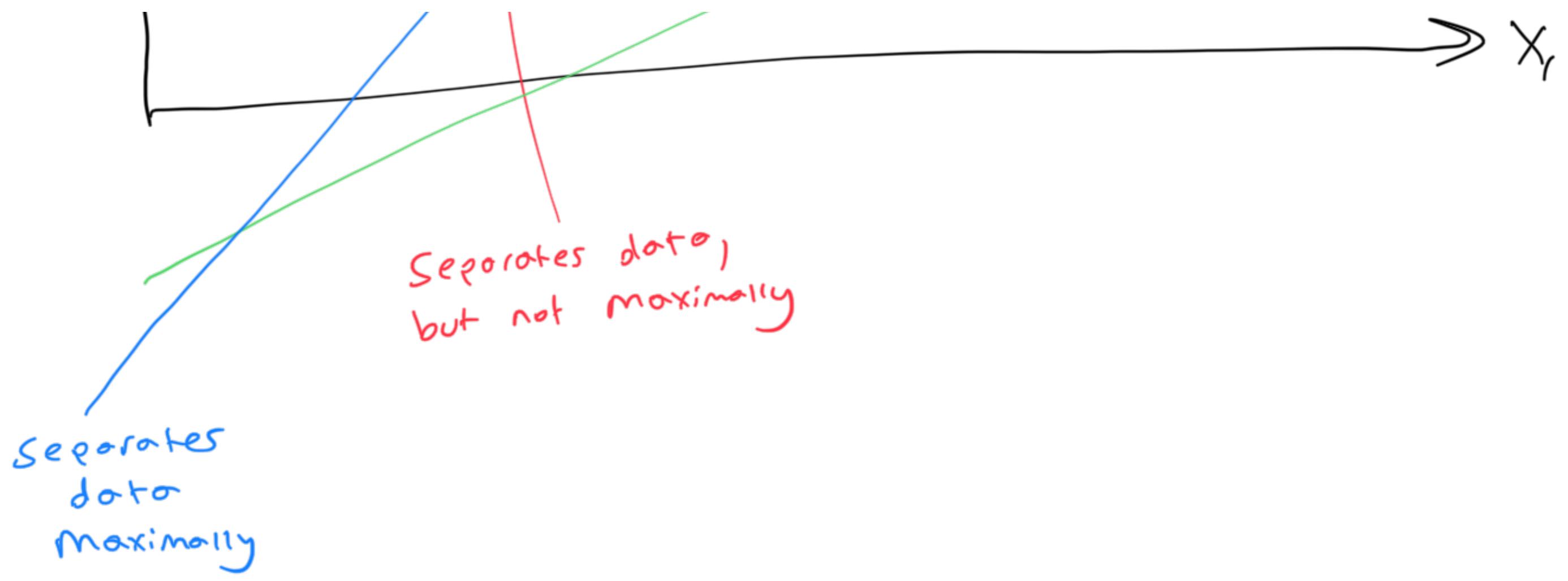


Maximally Separating line

When input dimensionality (# of input features)
is greater than 3...



Maximally Separating hyperplane

$x_2$

O O O O

X X X X X X X X

doesn't separate data

$\longrightarrow X_1$

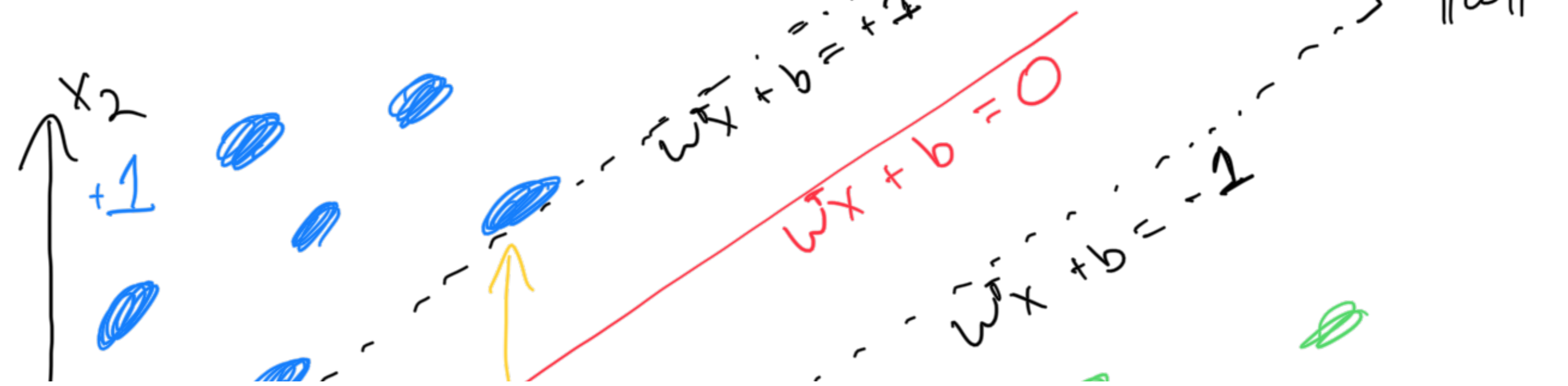<span style="color:red">Separates data, but not maximally</span>

<span style="color:blue">Separates data maximally</span>

# Hard - Margin SVC (Support Vector Classifier)

Find margin which doesn't allow... for violations of the margin

$\vec{w}\vec{x} + b = +1$

$\vec{w}x + b = 0$

$\vec{w}x + b = -1$

margin $= \dfrac{2}{\|w\|}$

$X_2$

+1

Support Vectors

$-1$

$\rightarrow X_1$

(y is coming off the page !!!)

$$w^T x_1 + b = +1$$
$$-(w^T x_2 + b = -1)$$
$$\overline{\frac{w^T(x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}}$$

$$\text{margin} = \frac{2}{\|w\|}$$

Math goal:

maximize margin $\left(\dfrac{2}{\|w\|}\right)$ subject to:

- margin $\geq 0$

Classify everything correctly

- $x_i w + b \geq +1$ if $y_i = +1$
- $x_i w + b \leq -1$ if $y_i = -1$

More concise: $y_i \left(w^T x_i + b\right) \geq 1$

Re-state this as a minimization problem:

minimize $\|w\|$ subject to:

<span style="color:red">Goal of hard margin</span>

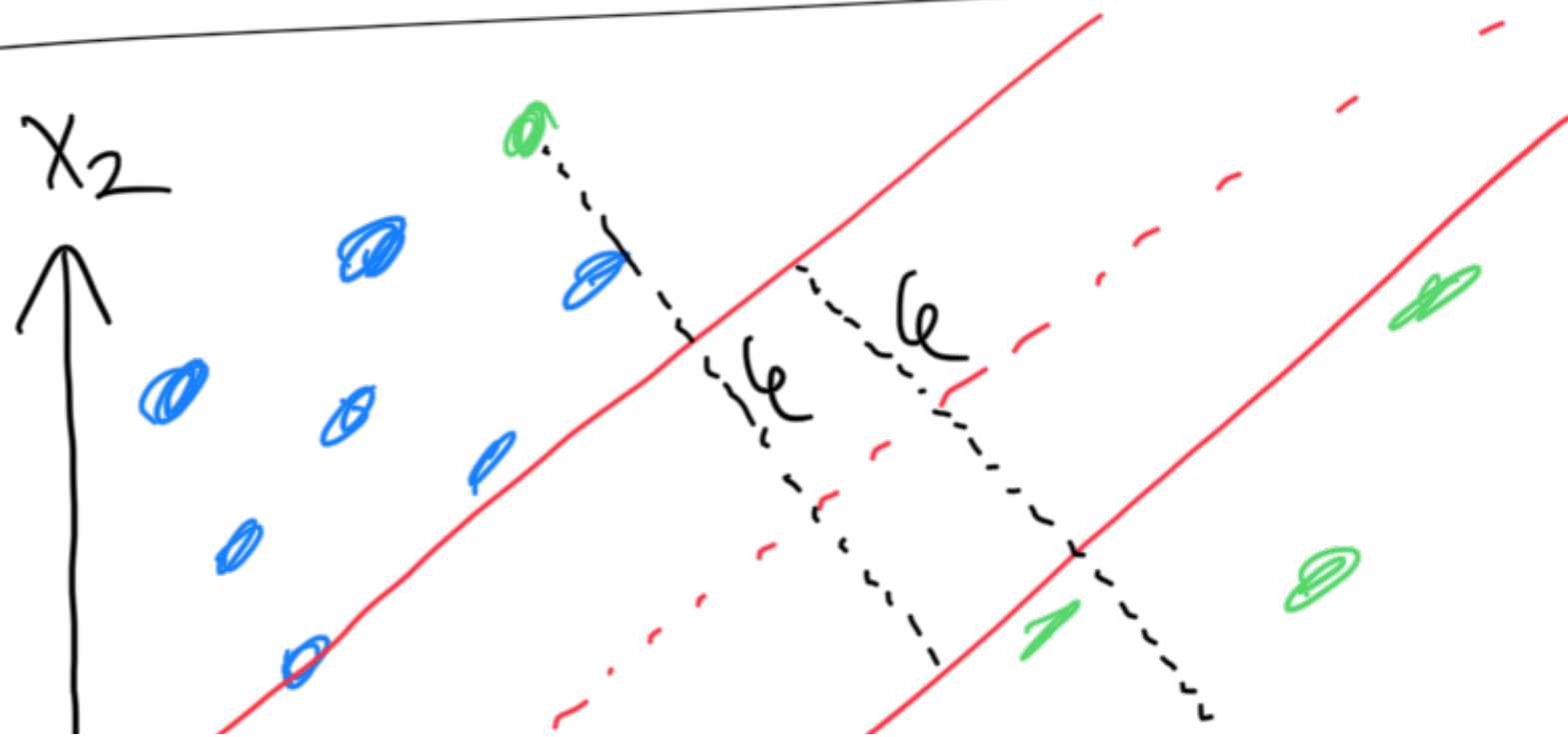<span style="color:red">SVM</span>

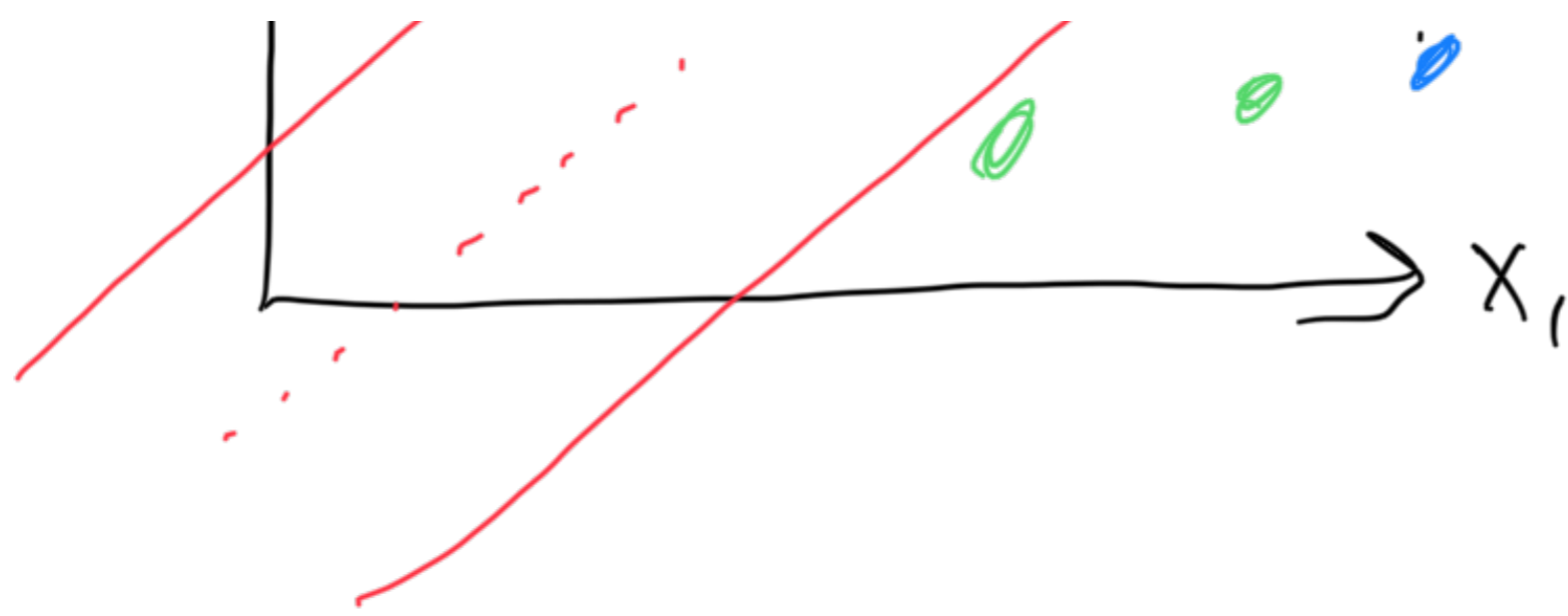$$y_i(w \, x_i + b) \,\underline{=}\, 1$$

## Constrained Optimization Problems[a]:
## Lagrange Multipliers

What if the data are not perfecty separable?

# Soft - Margin SVC

Introduce a "slack" variable $\xi$

New objective function is:

hyperparameter

$$\text{minimize } \|w\| + C \sum_{i=1}^{c} \xi_i$$

subject to : $y_i \left( w^T x_i + b \right) \geq 1 - \xi_i$

$$\cdot \, \xi_i \geq 0$$

# Hinge Loss Function

$$\text{Hinge Loss} = \max\left(0, \, 1 - y \cdot \hat{y}\right)$$

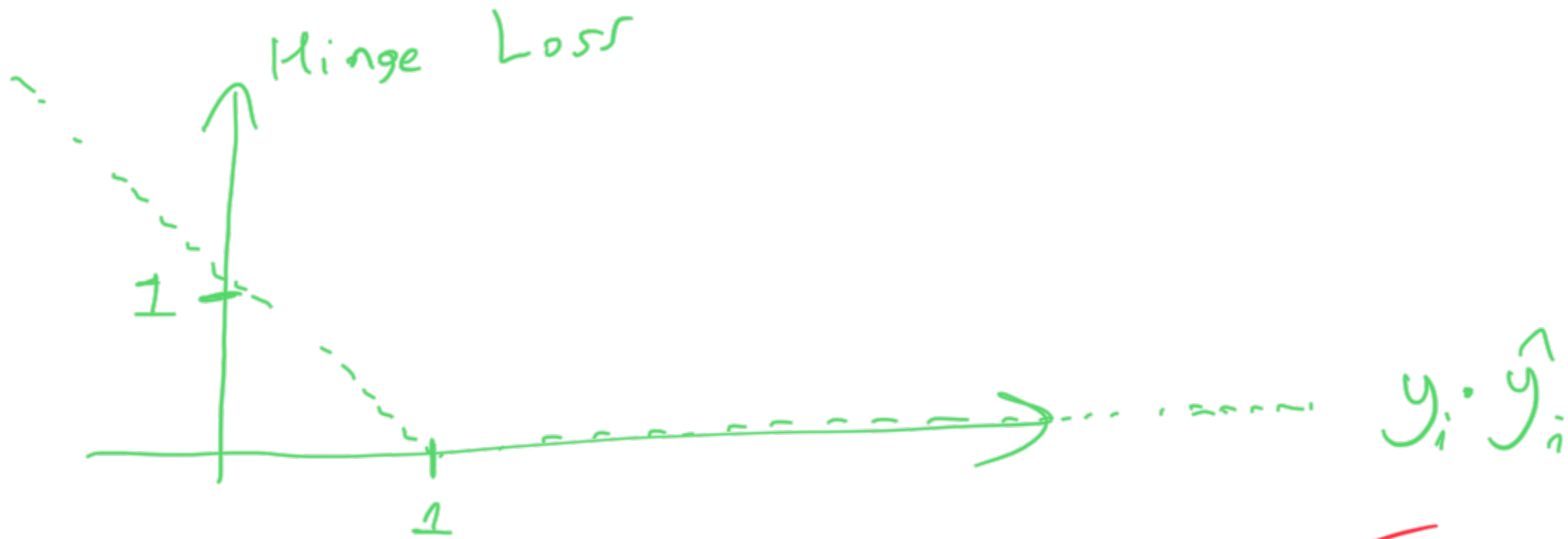(recall that $\hat{y}$ is the predicted value and $y$ is the actual value)

Properties:

- if $\hat{y}$ is correct and $|\hat{y}| \geq 1$:

  Hinge Loss $= 0$

- if $\hat{u}$ is correct and $|\hat{y}| < 1$:

$$0 < \text{Hinge Loss} < 1$$

- is $\hat{y}$ is incorrect:

$$\text{Hinge Loss} \geq 1$$
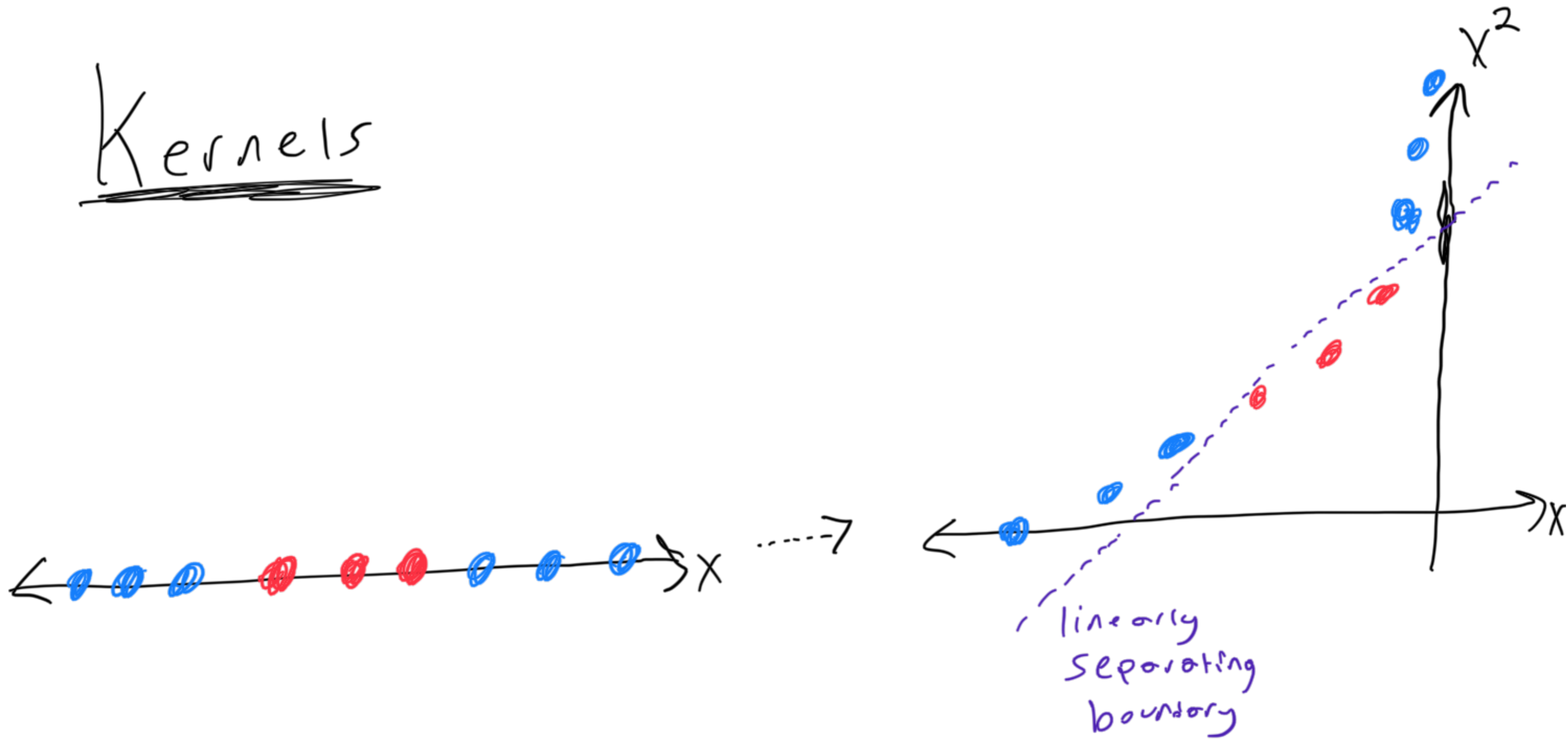
Hinge Loss

$1$

$1$

$y_i \cdot \hat{y}_i$

incorrect side of the boundary

correct side of the boundary but inside margin

correct side of the boundary and outside the margin

What is the data are not even close to being separable?

# Kernels



linearly separating boundary

Map data to high dimensional space

Key idea: non-linearly separable data can be

separable in high dimensions

**Fundamental Issue:** moving a reasonably sized feature space (e.g., 10,000 features) into a huge feature space which may be required to linearly separate the data (e.g, $10^{10}$ features) is computationally infeasible

The Solution:

## The Kernel Trick

[ Outside the scope of this class, an equivalent formulation of SVM is:

1

$$\max_{\alpha} \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \left( X_i^T X_j \right)$$

Scalars

matrices