

ML@LSE - Bootcamp 3: Support Vector classification methods

Ref. for today: Chapter 9, Hastie et al.

Introduction

A- Linear Classification

1. Maximal Margin Classifier

a. Hyperplanes

b. Maximal Margin Hyperplane

2. Soft Margin Classifier

a. The Optimization Problem of SMC

b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

a. What if we enlarged the Feature Space?

b. The New Optimization Problem

2. Kernel Tricks

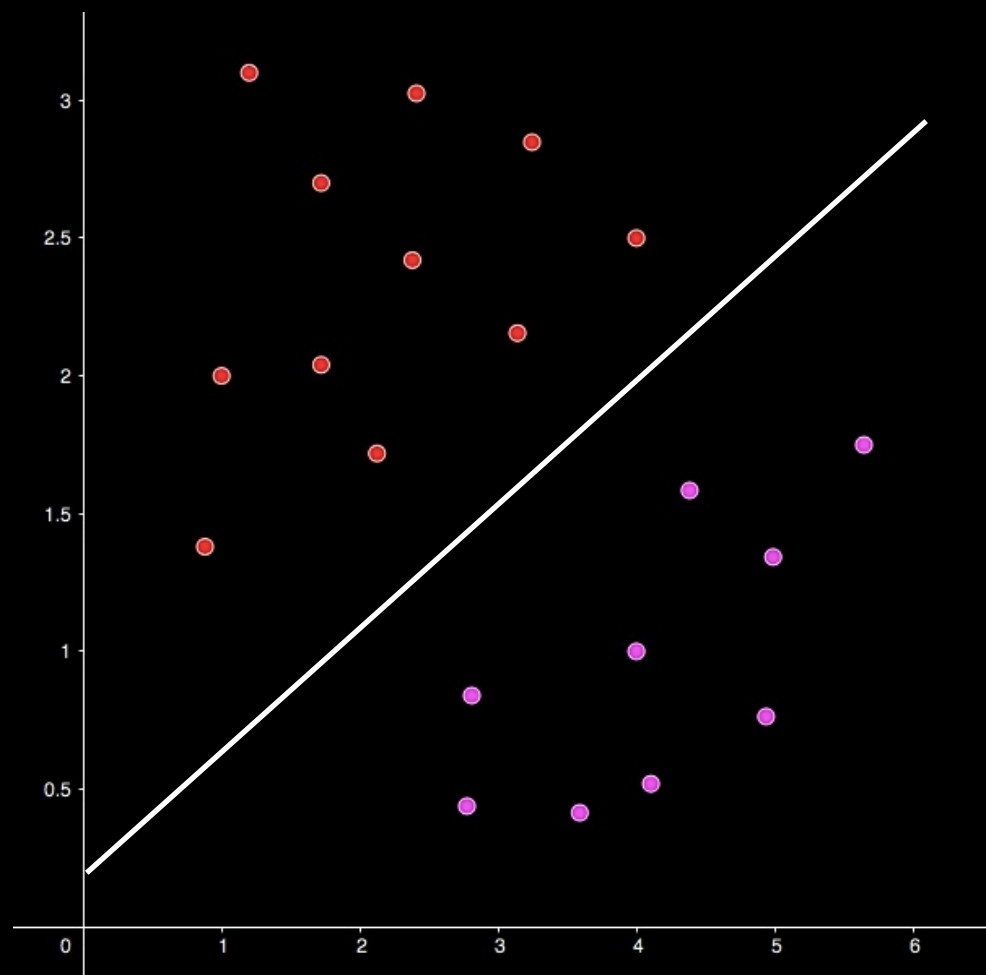
a. Inner Product and the Kernel Idea

b. Support Vector Machines

Last time we studied a supervised learning method that could be applied to regression and classification.

Today we're going to talk about another supervised learning method used for classification: Support Vector Classification.

What is this about?



What is a hyperplane?

A- Linear Classification

1. Maximal Margin Classifier

a. Hyperplanes

b. Maximal Margin Hyperplane

2. Soft Margin Classifier

a. The Optimization Problem of SMC

b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

a. What if we enlarged the Feature Space?

b. The New Optimization Problem

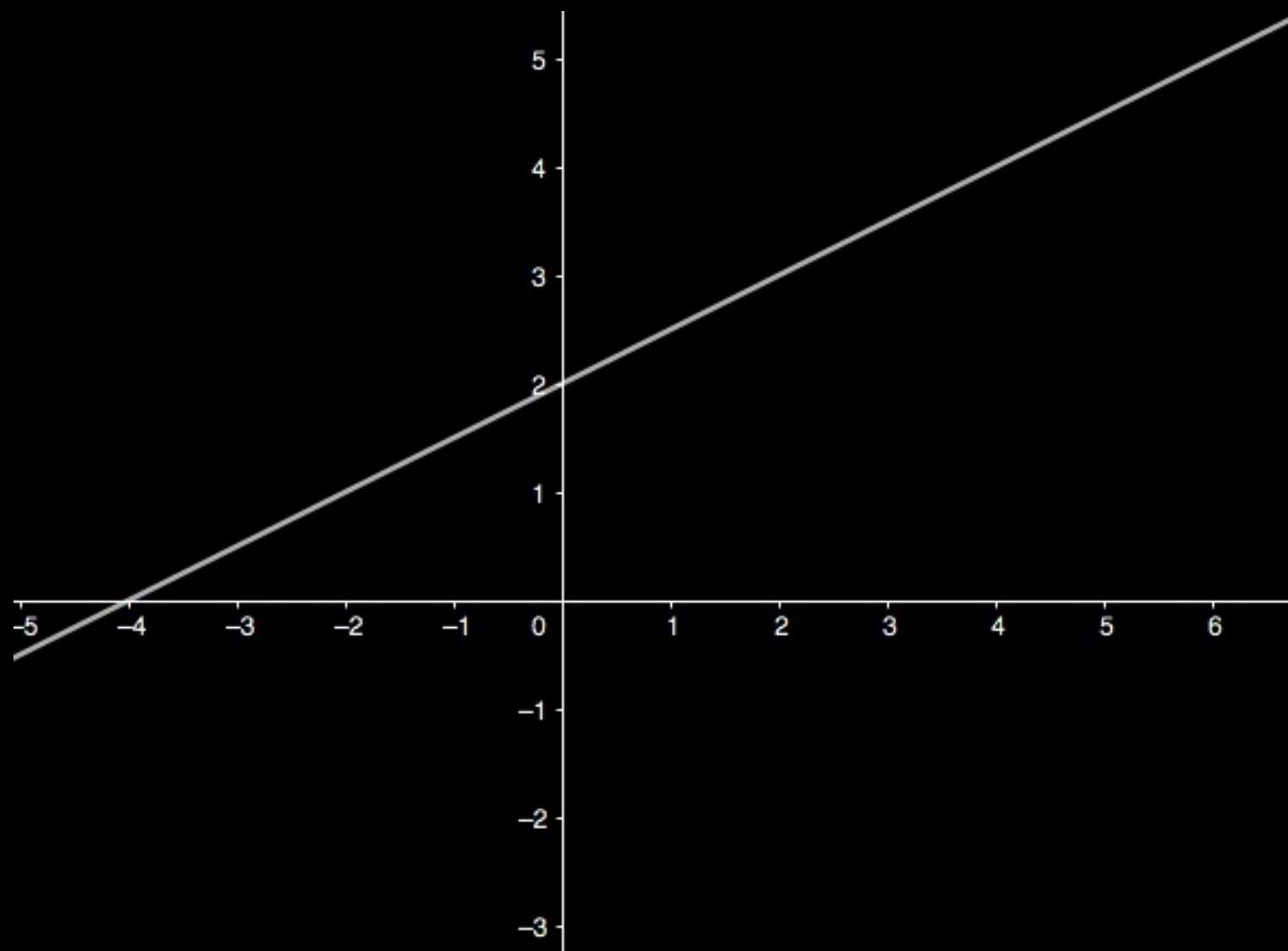
2. Kernel Tricks

a. Inner Product and the Kernel Idea

b. Support Vector Machines

Abstract definition: A flat affine subspace that is one dimension less than its ambient space.

What does that mean??



Hyperplanes and maximal margin

A- Linear Classification

1. Maximal Margin Classifier

a. Hyperplanes

b. Maximal Margin Hyperplane

2. Soft Margin Classifier

a. The Optimization Problem of SMC

b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

a. What if we enlarged the Feature Space?

b. The New Optimization Problem

2. Kernel Tricks

a. Inner Product and the Kernel Idea

b. Support Vector Machines

Mathematically, a hyperplane in p-dimension can be described by the equation:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0$$

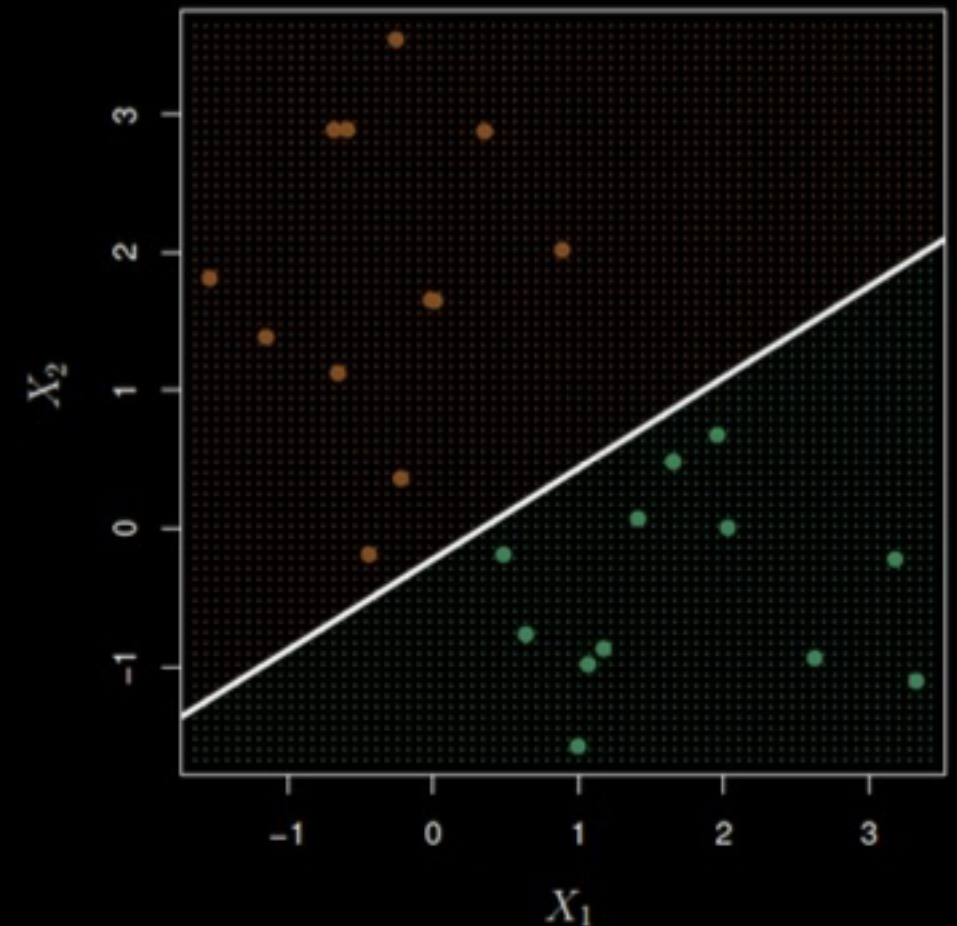
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p < 0$$

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p > 0$$

Let $Y_i=1$ if brown and $Y_i=-1$ if green.

$$\text{Let } f(\mathbf{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

Then given an observation \mathbf{x}_0 we can classify \mathbf{x}_0 according to the sign of $f(\mathbf{x}_0)$.



Maximum Margin Hyperplane

A- Linear Classification

1. Maximal Margin Classifier

a. Hyperplanes

b. Maximal Margin Hyperplane

2. Soft Margin Classifier

a. The Optimization Problem of SMC

b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

a. What if we enlarged the Feature Space?

b. The New Optimization Problem

2. Kernel Tricks

a. Inner Product and the Kernel Idea

b. Support Vector Machines

How can we find this separating hyperplane?

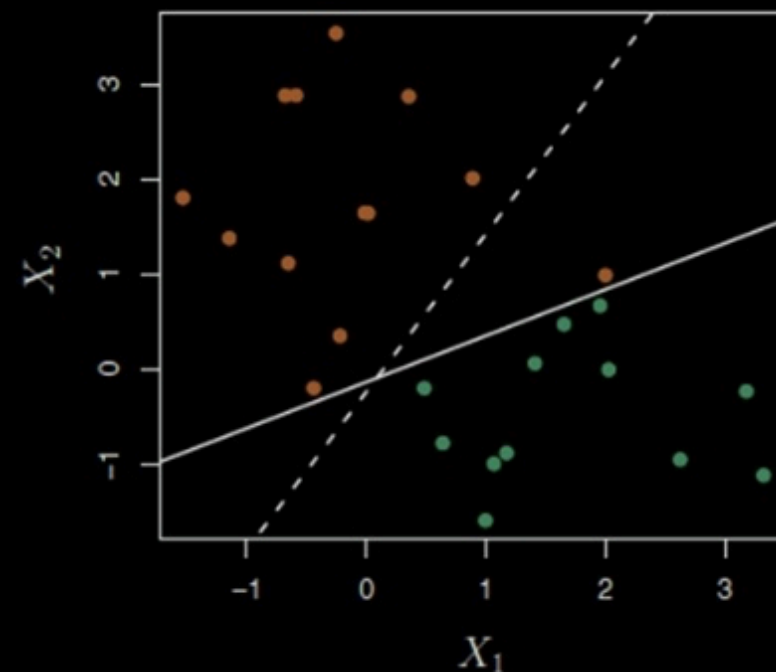
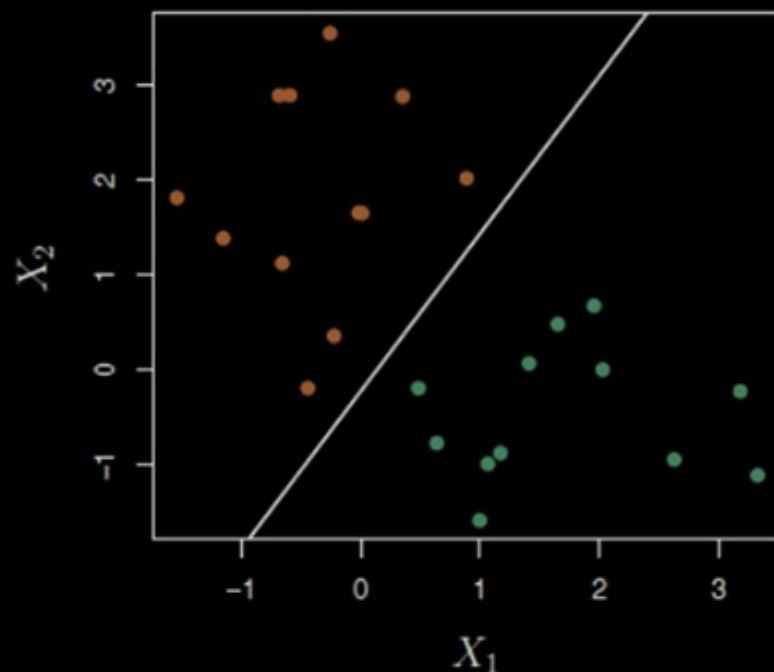
Let's try to find one manually

Maximum Margin Hyperplane

How can we find this separating hyperplane?

Let's try to find one manually

$$\max_{\beta_0, \beta_1, \dots, \beta_p} M \quad s.t. \quad \sum_{j=1}^p \beta_j^2 = 1, \quad y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M$$



Soft Margin Classifier

A more robust approach: Soft Margin Classifier (or support vector classifier):

We allow for *some* observations to be within the margin, or even to be on the other side of the hyperplane.

A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

a. The Optimization Problem of SMC

- b. The Bias-Variance tradeoff

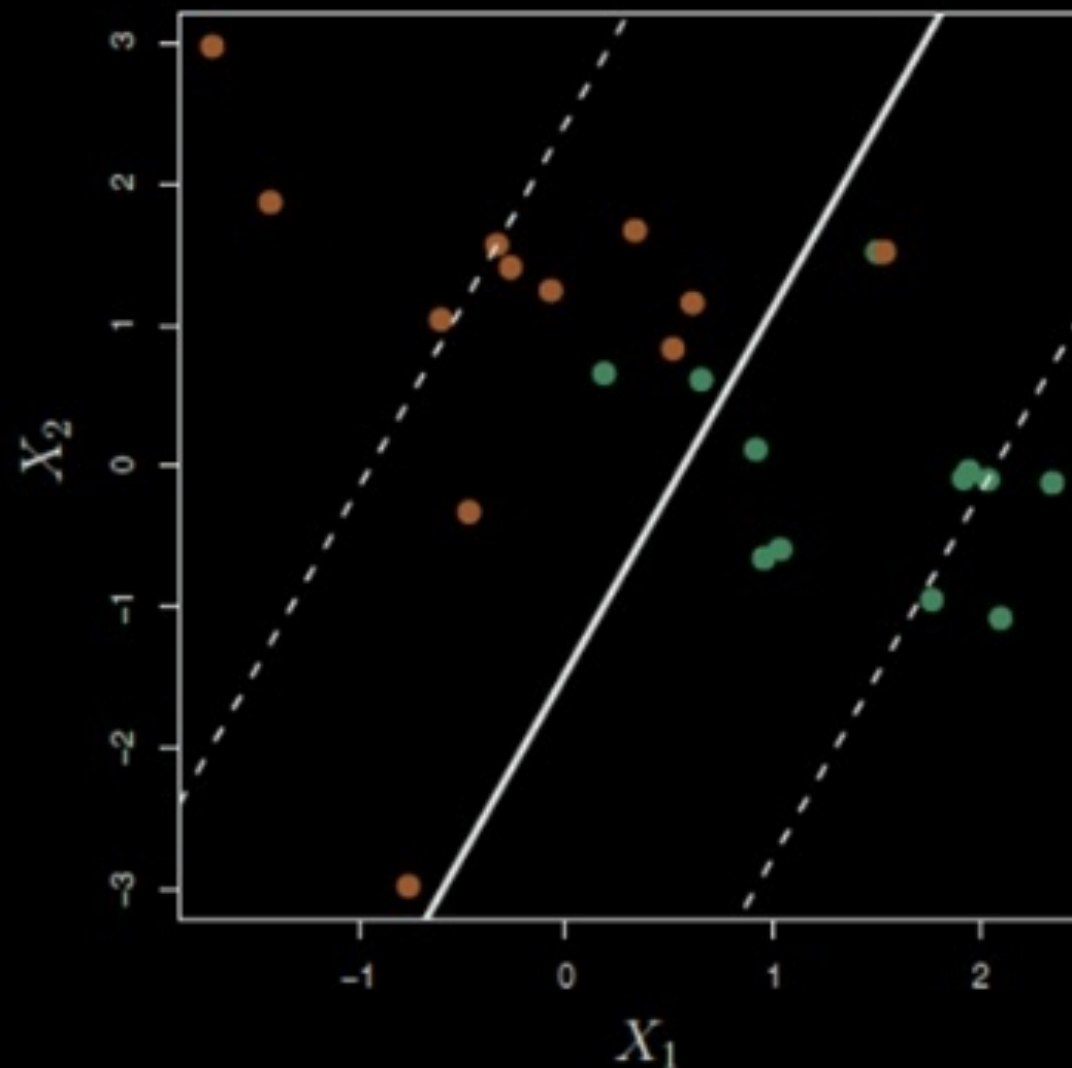
B- Support Vector Machines

1. Non-linear Decision Boundaries

- a. What if we enlarged the Feature Space?
- b. The New Optimization Problem

2. Kernel Tricks

- a. Inner Product and the Kernel Idea
- b. Support Vector Machines



Soft Margin Classifier: Optimization problem

A- Linear Classification

1. Maximal Margin Classifier

- Hyperplanes
- Maximal Margin Hyperplane

2. Soft Margin Classifier

a. The Optimization Problem of SMC

- The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

- What if we enlarged the Feature Space?
- The New Optimization Problem

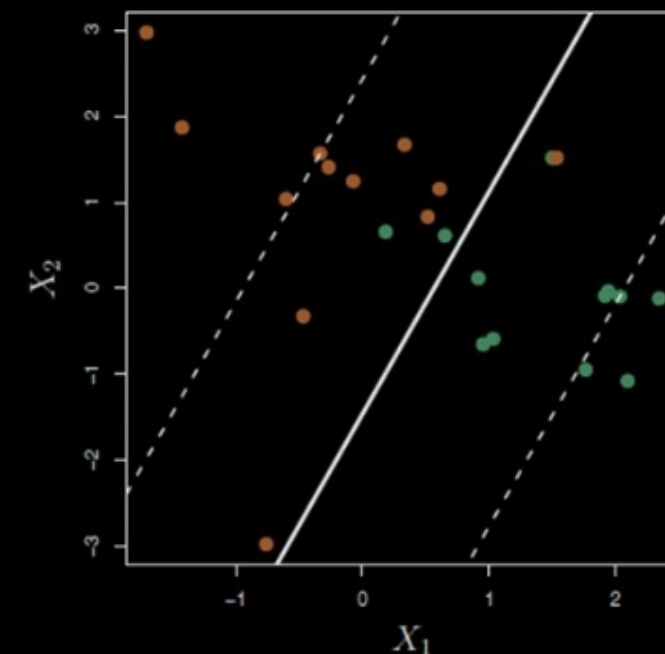
2. Kernel Tricks

- Inner Product and the Kernel Idea
- Support Vector Machines

The optimization problem:

$$\max_{\beta_0, \beta_1, \dots, \beta_p, s_1, \dots, s_n} M \quad s.t. \quad \begin{cases} \sum_{j=1}^p \beta_j^2 = 1 \\ y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M(1 - s_i) \\ s_i \geq 0, \sum_{i=1}^n s_i \leq C \end{cases}$$

- if $s_i = 0$, the i -th observation is on the correct side of the margin.
- if $0 < s_i < 1$, the i -th observation is on the wrong side of the margin but on the right side of the hyperplane.
- if $s_i < 0$, the i -th observation is on the wrong side of the hyperplane.



Soft Margin Classifier: Optimization problem

A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

- a. The Optimization Problem of SMC
- b. **The Bias-Variance tradeoff**

B- Support Vector Machines

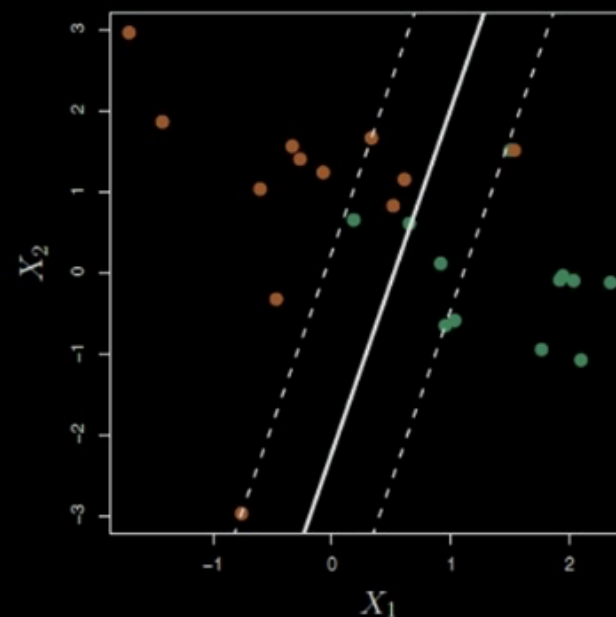
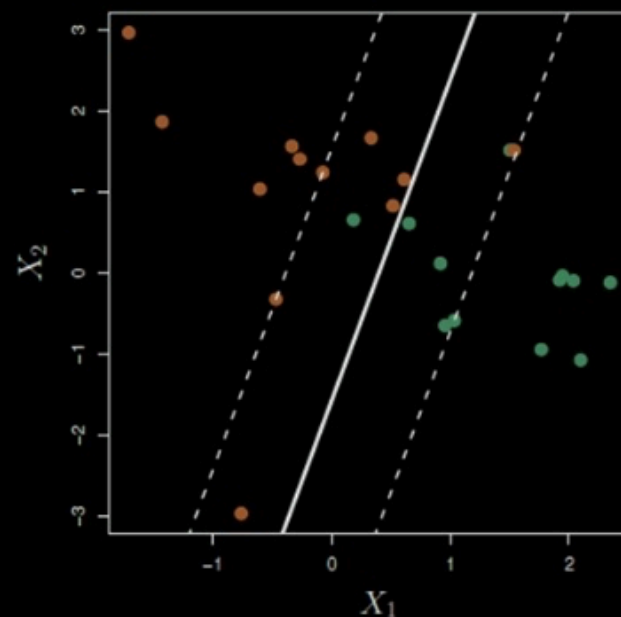
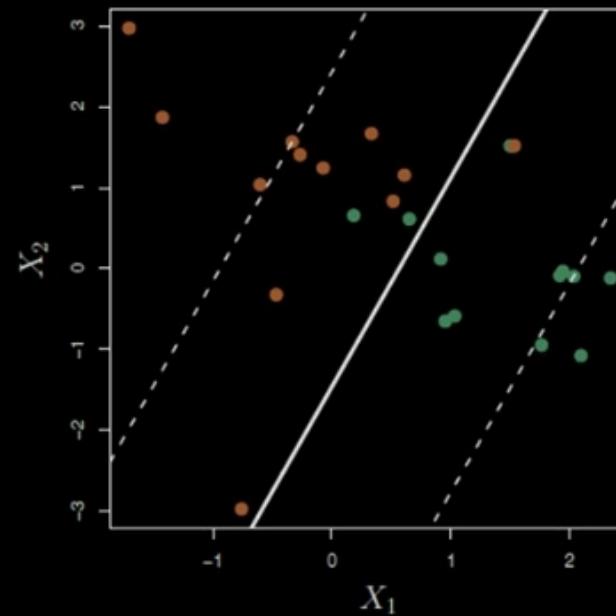
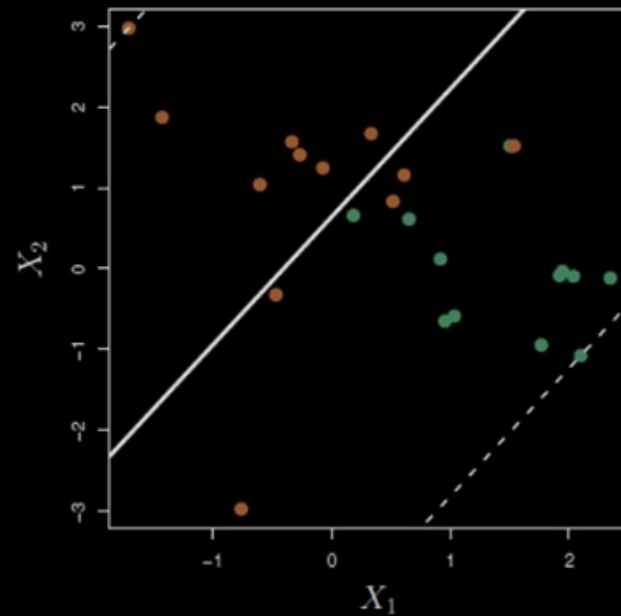
1. Non-linear Decision Boundaries

- a. What if we enlarged the Feature Space?
- b. The New Optimization Problem

2. Kernel Tricks

- a. Inner Product and the Kernel Idea
- b. Support Vector Machines

$$\max_{\beta_0, \beta_1, \dots, \beta_p, s_1, \dots, s_n} M \quad s.t. \quad \begin{cases} \sum_{j=1}^p \beta_j^2 = 1 \\ y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M(1 - s_i) \\ s_i \geq 0, \sum_{i=1}^n s_i \leq C \end{cases}$$



C is a regularization parameter

Non-linear Decision Boundaries

What can we do when the data is not linearly separable?

A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

- a. The Optimization Problem of SMC
- b. The Bias-Variance tradeoff

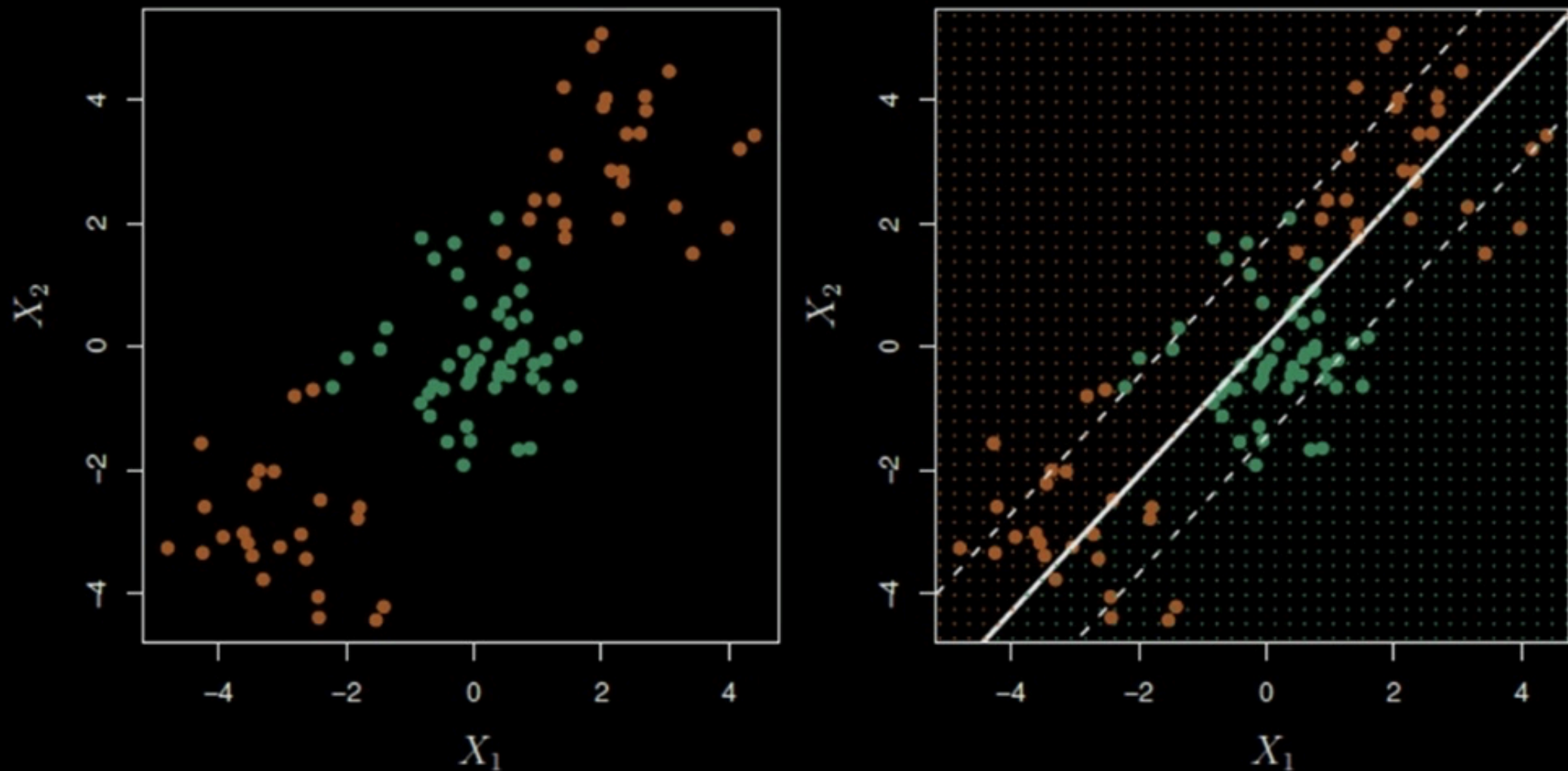
B- Support Vector Machines

1. Non-linear Decision Boundaries

- a. **What if we enlarged the Feature Space?**
- b. The New Optimization Problem

2. Kernel Tricks

- a. Inner Product and the Kernel Idea
- b. Support Vector Machines



Enlarge the feature space by adding powers of the features:

$$X_1, X_1^2, X_2, X_2^2, \dots, X_p, X_p^2$$

Non-linear Decision Boundaries

Enlarge the feature space by adding powers of the features:

$$X_1, X_1^2, X_2, X_2^2, \dots, X_p, X_p^2$$

A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

- a. The Optimization Problem of SMC
- b. The Bias-Variance tradeoff

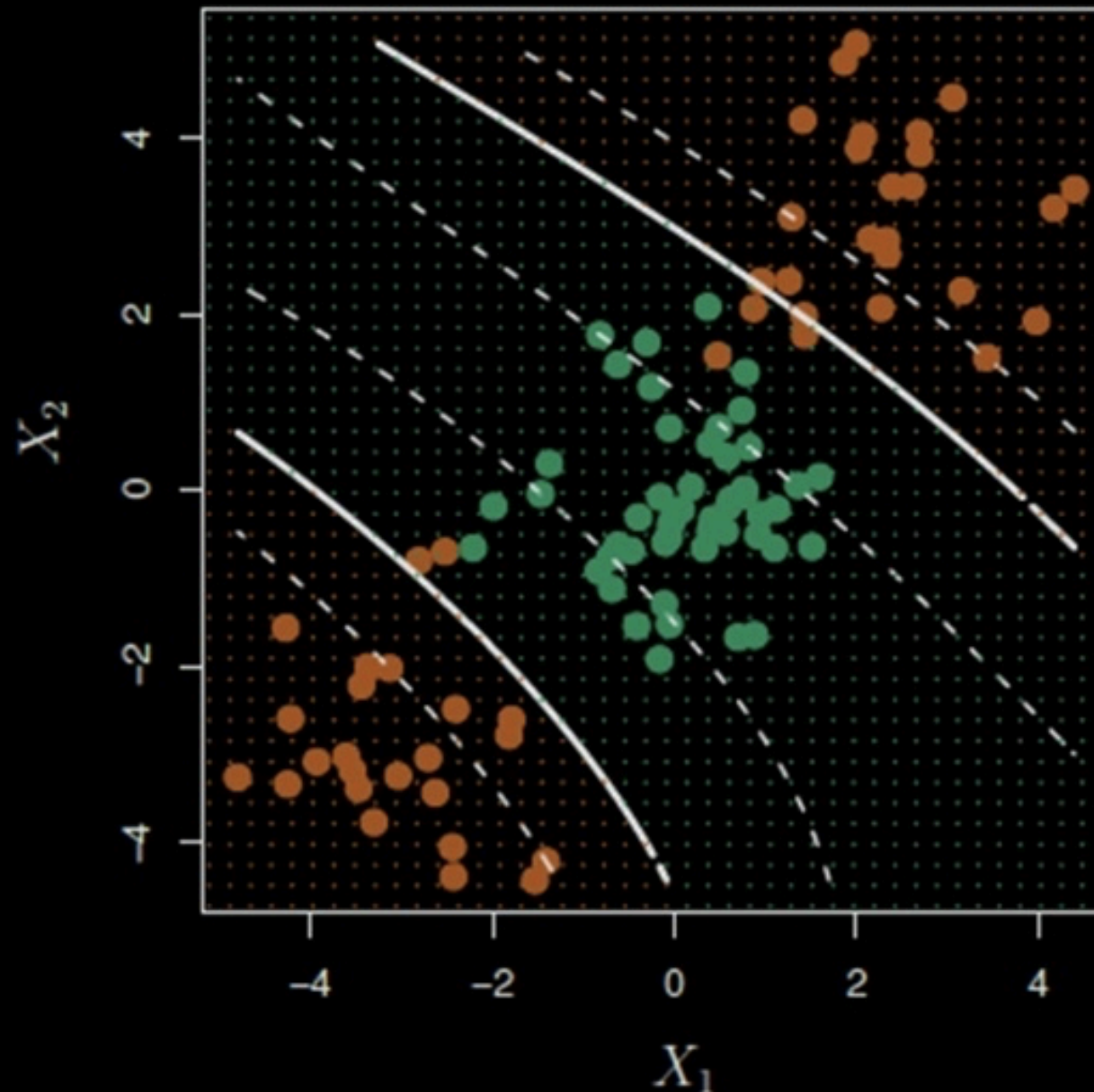
B- Support Vector Machines

1. Non-linear Decision Boundaries

- a. What if we enlarged the Feature Space?
- b. The New Optimization Problem

2. Kernel Tricks

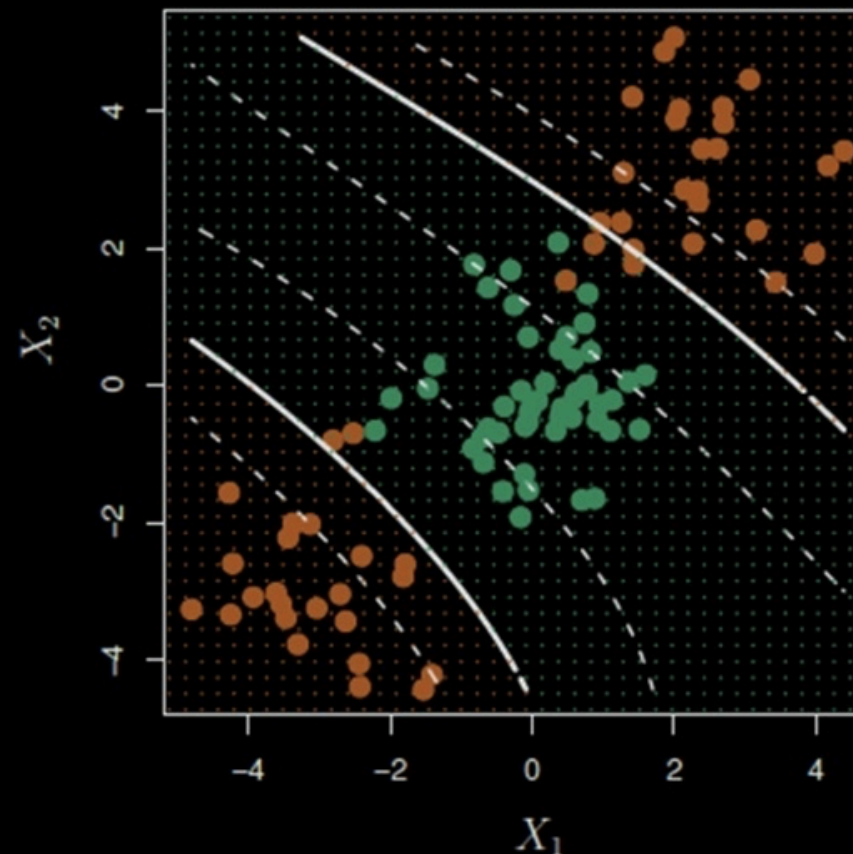
- a. Inner Product and the Kernel Idea
- b. Support Vector Machines



Optimization problem with polynomials

Here is the new optimization problem:

$$\max_{\beta_0, \beta_1, \dots, \beta_p, s_1, \dots, s_n} M \quad s.t. \quad \begin{cases} \sum_{j=1}^p \beta_j^2 = 1 \\ y_i(\beta_0 + \beta_{11}x_{i1} + \dots + \beta_{p1}x_{ip} + \beta_{12}x_{i1}^2 + \beta_{p1}x_{ip}^2) \geq M(1 - s_i) \\ s_i \geq 0, \sum_{i=1}^n s_i \leq C \end{cases}$$



A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

- a. The Optimization Problem of SMC
- b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

- a. What if we enlarged the Feature Space?

b. The New Optimization Problem

2. Kernel Tricks

- a. Inner Product and the Kernel Idea
- b. Support Vector Machines

Kernels

! *More advanced*

A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

- a. The Optimization Problem of SMC
- b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

- a. What if we enlarged the Feature Space?
- b. The New Optimization Problem

2. Kernel Tricks

a. Inner Product and the Kernel Idea

- b. Support Vector Machines

We can represent the Support Vector Classification function with the classical inner product (dot product):

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle$$

We can generalize this by replacing the dot product above by a Kernel. A Kernel is the dot product between the image of two vectors \mathbf{x} and \mathbf{y} by φ , where $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$. $Ker(\mathbf{x}, \mathbf{y}) = \langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle$.

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i Ker(\mathbf{x}, \mathbf{x}_i).$$

Kernels

! *More advanced*

A- Linear Classification

1. Maximal Margin Classifier

- a. Hyperplanes
- b. Maximal Margin Hyperplane

2. Soft Margin Classifier

- a. The Optimization Problem of SMC
- b. The Bias-Variance tradeoff

B- Support Vector Machines

1. Non-linear Decision Boundaries

- a. What if we enlarged the Feature Space?
- b. The New Optimization Problem

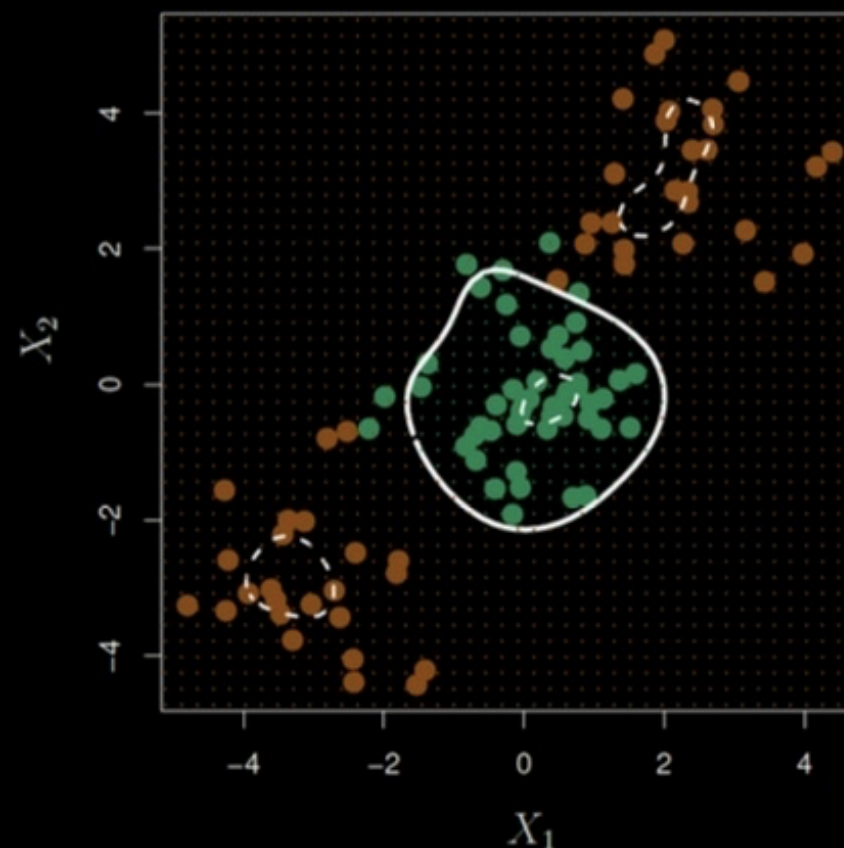
2. Kernel Tricks

- a. Inner Product and the Kernel Idea

b. Support Vector Machines

We can generalize this by replacing the dot product above by a Kernel. A Kernel is the dot product between the image of two vectors \mathbf{x} and \mathbf{y} by φ , where $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$. $Ker(\mathbf{x}, \mathbf{y}) = \langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle$.

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i Ker(\mathbf{x}, \mathbf{x}_i).$$



Example of Kernel: Radial Kernel

$$Ker(\mathbf{x}, \mathbf{x}_i) = \exp\left\{-\gamma \sum_{j=1}^p (x_j - x_{ij})^2\right\}$$

Thanks for coming!