

Machine Learning for Official Statistics & SDGs

Support Vector Machine



[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier
- ▶ Lots of refinements

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier
- ▶ Lots of refinements
- ▶ Based on two very different ideas

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier
- ▶ Lots of refinements
- ▶ Based on two very different ideas
 - ▶ Dimension augmentation

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier
- ▶ Lots of refinements
- ▶ Based on two very different ideas
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier
- ▶ Lots of refinements
- ▶ Based on two very different ideas
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
- ▶ Technical details require advanced mathematics

[UNDERSTANDING SVM]

Support Vector Machine (SVM) is a classifier

- ▶ One of the best "out of the box" classifier
 - ▶ Lots of refinements
 - ▶ Based on two very different ideas
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
 - ▶ Technical details require advanced mathematics
- ↪ Focus on seminal ideas and intuitions

[THE VAPNIK–CHERVONENKIS DIMENSION]

[THE VAPNIK–CHERVONENKIS DIMENSION]

- ▶ The *Vapnik–Chervonenkis* (VC) dimension measures the complexity that can be learned by a classifier

[THE VAPNIK–CHERVONENKIS DIMENSION]

- ▶ The *Vapnik–Chervonenkis* (VC) dimension measures the complexity that can be learned by a classifier
- ▶ The VC dimension is defined as the maximum number of points allowing a separation by a linear classifier

[THE VAPNIK–CHERVONENKIS DIMENSION]

- ▶ The *Vapnik–Chervonenkis* (VC) dimension measures the complexity that can be learned by a classifier
 - ▶ The VC dimension is defined as the maximum number of points allowing a separation by a linear classifier
- ↪ $VC(f) = \# \text{ points always separable in the current space}$

[THE VAPNIK–CHERVONENKIS DIMENSION]

- ▶ The *Vapnik–Chervonenkis* (VC) dimension measures the complexity that can be learned by a classifier
 - ▶ The VC dimension is defined as the maximum number of points allowing a separation by a linear classifier
- ↪ $VC(f) = \# \text{ points always separable in the current space}$
- ▶ Useful to measure the complexity of a particular model in a given space

[THE VAPNIK–CHERVONENKIS DIMENSION]

[THE VAPNIK–CHERVONENKIS DIMENSION]

In 1-D, there are 4 situations with 2 points:



Both orange

[THE VAPNIK–CHERVONENKIS DIMENSION]

In 1-D, there are 4 situations with 2 points:



Both blue

[THE VAPNIK–CHERVONENKIS DIMENSION]

In 1-D, there are 4 situations with 2 points:



One orange, one blue

[THE VAPNIK–CHERVONENKIS DIMENSION]

In 1-D, there are 4 situations with 2 points:



One blue, one orange

[THE VAPNIK–CHERVONENKIS DIMENSION]

What happens with 3 points?



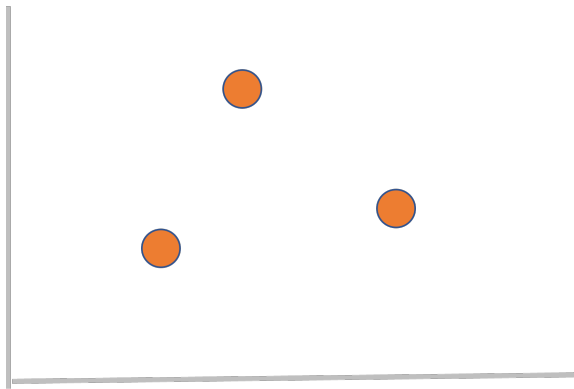
A situation where no linear classifier works!

$$\hookrightarrow VC(f_{lin}) = 2$$

[THE VAPNIK–CHERVONENKIS DIMENSION]

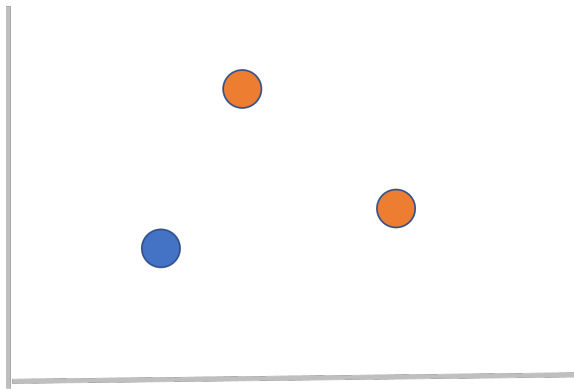
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



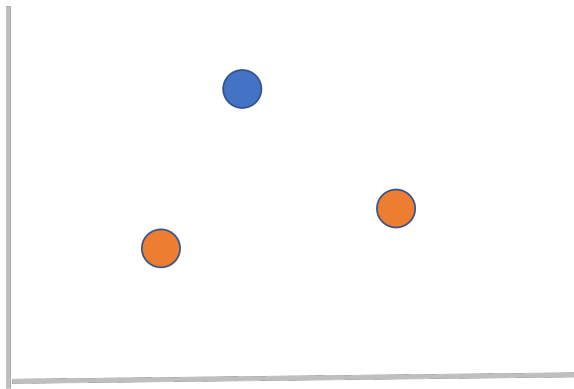
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



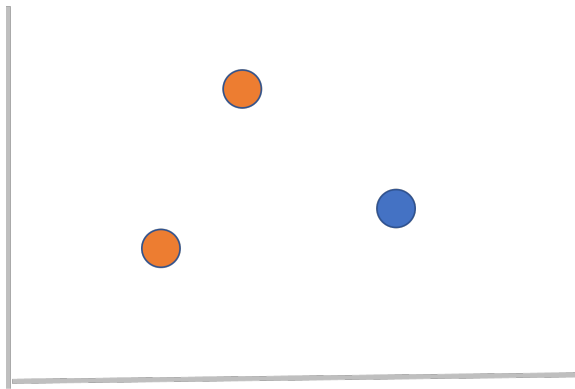
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



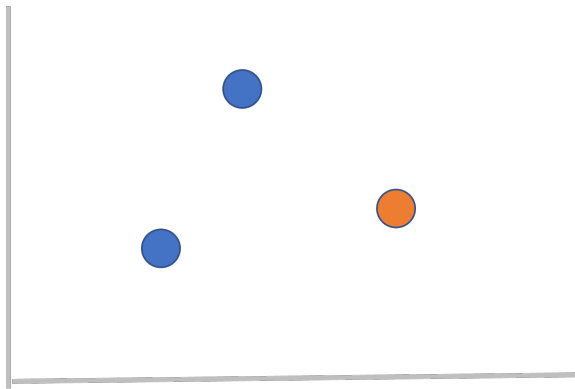
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



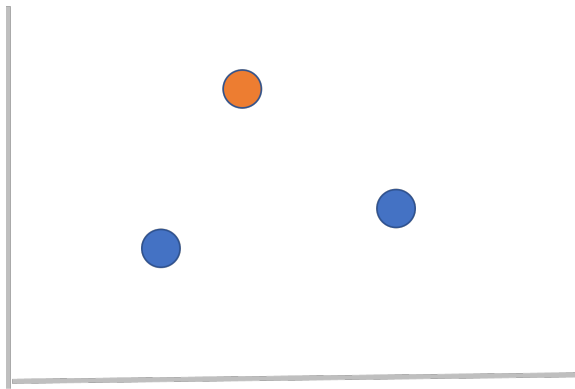
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



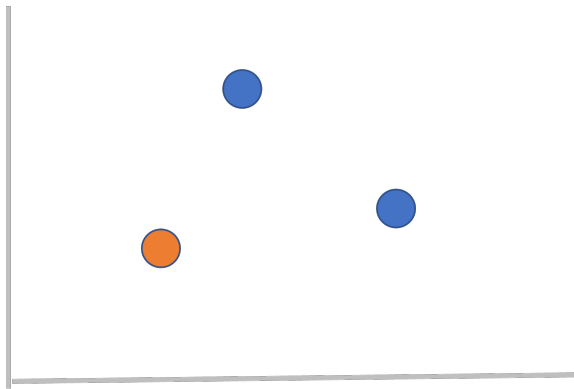
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



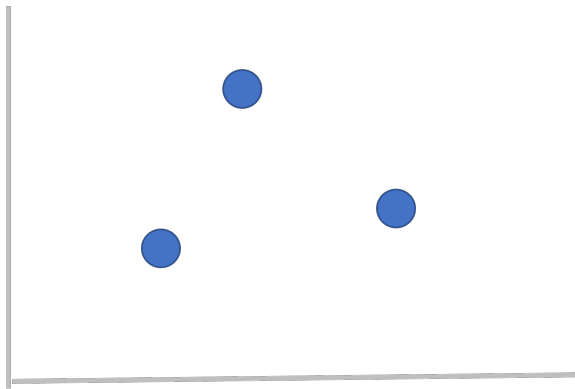
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



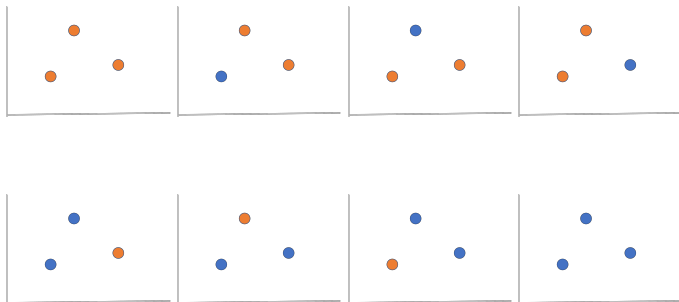
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



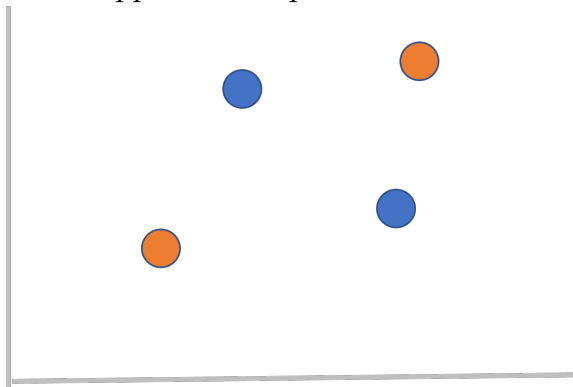
[THE VAPNIK–CHERVONENKIS DIMENSION]

In 2-D, there are $2^3 = 8$ different situations with 3 points:



[THE VAPNIK–CHERVONENKIS DIMENSION]

What happens with 4 points?



A situation where no linear classifier works!

$$\hookrightarrow VC(f_{lin}) = 3$$

[INTERPRETATION OF THE VC DIMENSION]

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points
It can be shown that:

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points

It can be shown that:

- ▶ Classifiers with a **high** VC dimension will be better at learning and fitting the training data set

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points

It can be shown that:

- ▶ Classifiers with a **high** VC dimension will be better at learning and fitting the training data set
- ↪ **High** risk of over-fitting

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points

It can be shown that:

- ▶ Classifiers with a **high** VC dimension will be better at learning and fitting the training data set
↳ **High** risk of over-fitting
- ▶ Classifiers with a **low** VC dimension make more errors on the training data set, but may be better in prediction

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points

It can be shown that:

- ▶ Classifiers with a **high** VC dimension will be better at learning and fitting the training data set
 - ↪ **High** risk of over-fitting
- ▶ Classifiers with a **low** VC dimension make more errors on the training data set, but may be better in prediction
 - ↪ **Low** risk of over-fitting

[INTERPRETATION OF THE VC DIMENSION]

- ▶ $VC(f)$ measures the complexity of the classifier f , linear or not
- ▶ A highly non-linear classifier will always be able to separate a high number of points

It can be shown that:

- ▶ Classifiers with a **high** VC dimension will be better at learning and fitting the training data set
 - ↪ **High** risk of over-fitting
- ▶ Classifiers with a **low** VC dimension make more errors on the training data set, but may be better in prediction
 - ↪ **Low** risk of over-fitting
- ▶ New version of the **Bias-Variance** trade off!

[IMPLICATIONS OF THE VC DIMENSION]

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

- $\epsilon(\cdot)$ is the error rate,

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

- $\epsilon(\cdot)$ is the error rate,
- m is the nb of training sets and

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

- $\epsilon(\cdot)$ is the error rate,
- m is the nb of training sets and
- $G(\cdot)$ a function **decreasing** with m (obviously)

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

- $\epsilon(\cdot)$ is the error rate,
- m is the nb of training sets and
- $G(\cdot)$ a function **decreasing** with m (obviously)
increasing with VC !!

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

- $\epsilon(\cdot)$ is the error rate,
- m is the nb of training sets and
- $G(\cdot)$ a function **decreasing** with m (obviously)
increasing with VC !!

↪ The higher $\text{VC}(f)$ the greater the difference between performance on the *training* and *validation* sample

[IMPLICATIONS OF THE VC DIMENSION]

- The Vapnik–Chervonenkis inequality (*simplified*):

$$\epsilon_{\text{validation}}(f) \leq \epsilon_{\text{train}}(f) + G(\text{VC}(f), m)$$

Where:

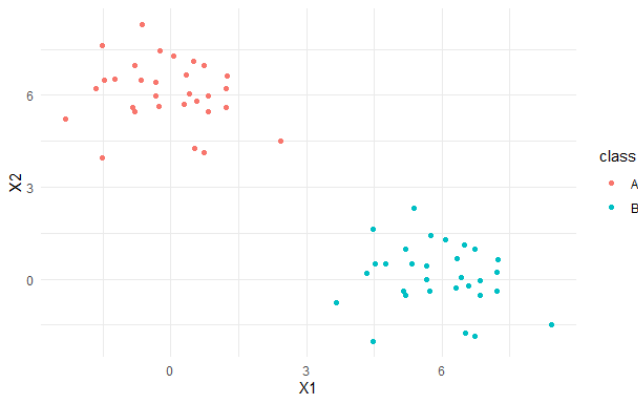
- $\epsilon(\cdot)$ is the error rate,
- m is the nb of training sets and
- $G(\cdot)$ a function **decreasing** with m (obviously)
increasing with VC !!

- ↪ The higher $\text{VC}(f)$ the greater the difference between performance on the *training* and *validation* sample
- ↪ Better to use simple models (low VC)

[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

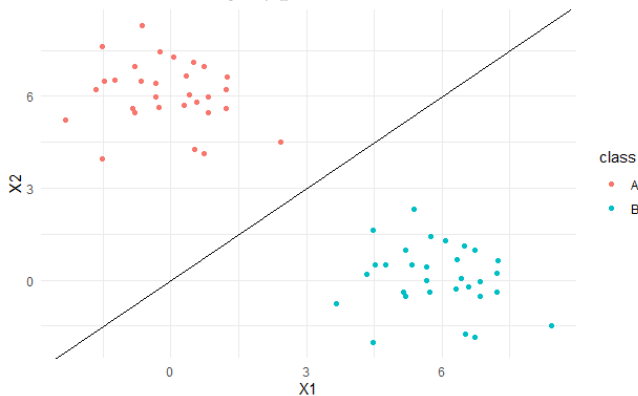
[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

Take the following hypothetical situation



[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

Take the following hypothetical situation

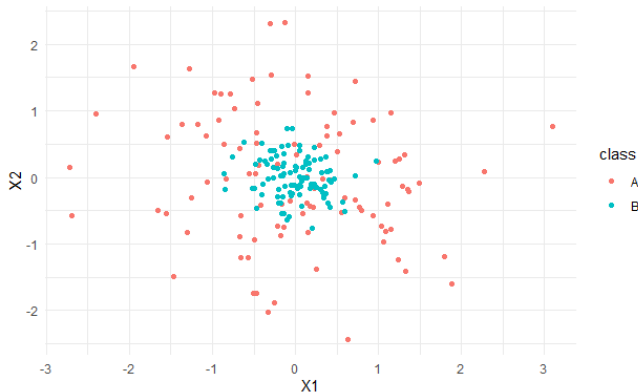


↪ It is easy to *separate* the two classes with a linear classifier

[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

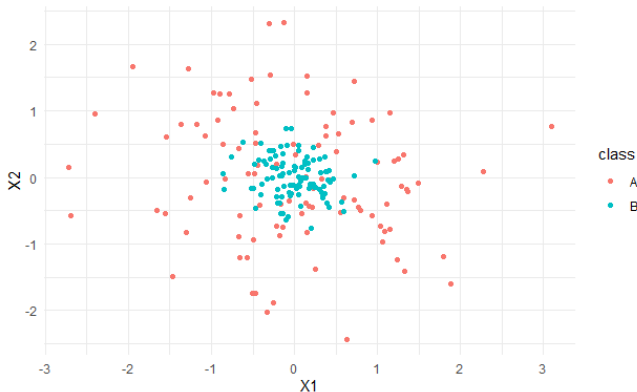
[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

But in this situation?



[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

But in this situation?



↪ There is no way to *separate* the two classes with a linear classifier!!

[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

[LINEAR BOUNDARIES IN A SIMPLE EXAMPLE]

Unless...

Introduction
○

The VC dimension
○○○○○

Dimension augmentation
○○●

Soft margin
○○○

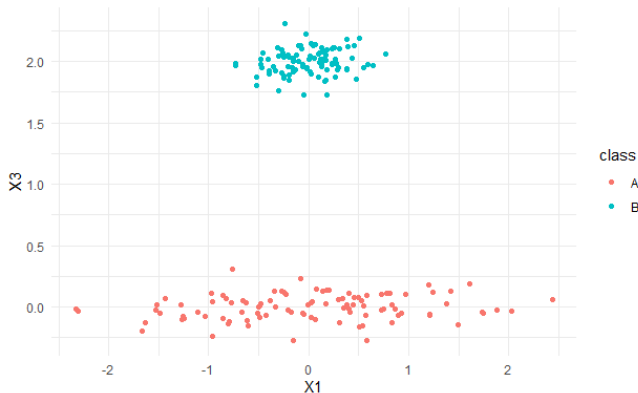
Kernel trick
○○○○

Wrap-up
○

[DIMENSION AUGMENTATION]

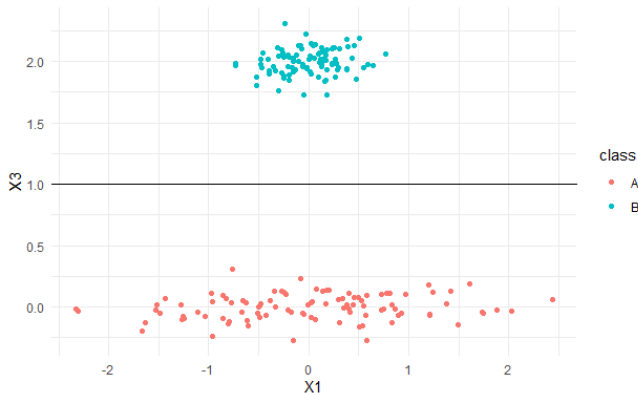
[DIMENSION AUGMENTATION]

Unless one adds a third dimension!



[DIMENSION AUGMENTATION]

Unless one adds a third dimension!

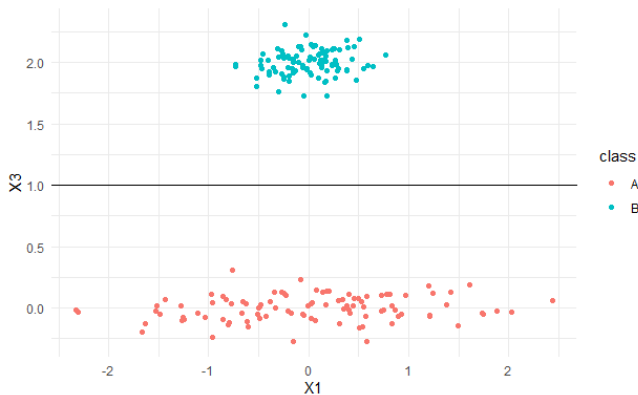


↪ A linear classifier can *separate* the two classes in an *augmented* space (third dimension)

[MAXIMAL MARGIN CLASSIFIER]

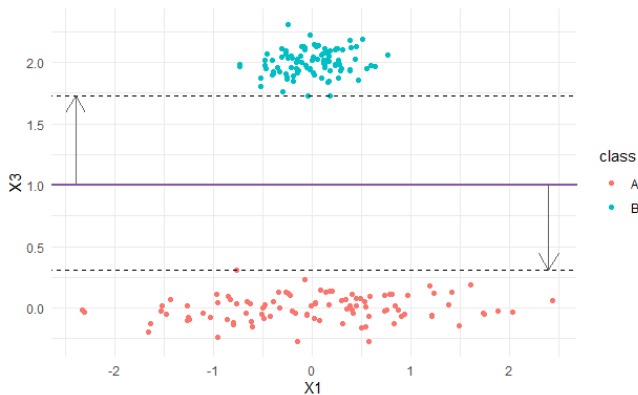
[MAXIMAL MARGIN CLASSIFIER]

The maximum margin classifier is defined as:



[MAXIMAL MARGIN CLASSIFIER]

The maximum margin classifier is defined as:



↪ M is the margin.

Introduction
○

The VC dimension
○○○○○

Dimension augmentation
○○○○

Soft margin
○●○

Kernel trick
○○○○

Wrap-up
○

[MAXIMAL MARGIN CLASSIFIER]

[MAXIMAL MARGIN CLASSIFIER]

The optimization problem to find this margin is :
with $(\beta_0 + x'\beta)$ defining the boundary (line or hyperplane)

$$\begin{aligned} & \text{Max}_{\beta_0, \beta} \quad M \\ & \text{subject to : } y_i(\beta_0 + x'_i\beta) \geq M \end{aligned}$$

and $y_i = 1$ or -1 depending on the class of observation i

[MAXIMAL MARGIN CLASSIFIER]

The optimization problem to find this margin is :
↪ This can be relaxed

$$\begin{aligned} & \text{Max}_{\beta_0, \beta} \quad M \\ \text{subject to : } & y_i(\beta_0 + x'_i \beta) \geq M - \xi \end{aligned}$$

where ξ is the *slack* parameter

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

The slack parameter will play an important role

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

The slack parameter will play an important role

- ▶ A small *slack* allows few exceptions

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

The slack parameter will play an important role

- ▶ A small *slack* allows few exceptions

↪ Boundary (classifier) sensitive to outliers

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

The slack parameter will play an important role

- ▶ A small *slack* allows few exceptions
- ↪ Boundary (classifier) sensitive to outliers
- ▶ A large *slack* allows lots of errors and misclassifications

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

The slack parameter will play an important role

- ▶ A small *slack* allows few exceptions
- ↪ Boundary (classifier) sensitive to outliers
- ▶ A large *slack* allows lots of errors and misclassifications

The optimization problem is equivalent to:

$$\begin{aligned} \underset{\beta_0, \beta}{\text{Min}} \quad & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to :} \quad & y_i(\beta_0 + x' \beta) \geq M - \xi_i \end{aligned} \tag{1}$$

[SOFT MARGIN]

Introducing *soft margins* is one idea implemented in SVM

The slack parameter will play an important role

- ▶ A small *slack* allows few exceptions
- ↪ Boundary (classifier) sensitive to outliers
- ▶ A large *slack* allows lots of errors and misclassifications

The optimization problem is equivalent to:

$$\begin{aligned} \text{Min}_{\beta_0, \beta} \quad & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to :} \quad & y_i(\beta_0 + x' \beta) \geq M - \xi_i \end{aligned} \quad (1)$$

- ↪ SVM is a trade-off between maximal M and cost C of errors due to ξ

[THE *kernel trick*]

The previous optimization method applies to an (*augmented*) space with separable classes

[THE *kernel trick*]

The previous optimization method applies to an (*augmented*) space with separable classes

- How do we determine this *augmented* space?

[THE *kernel trick*]

The previous optimization method applies to an (*augmented*) space with separable classes

- How do we determine this *augmented* space?

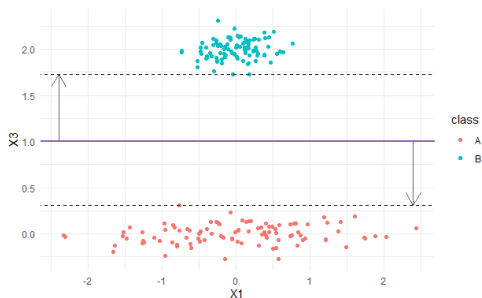
How to find $\phi(\cdot)$ so that $X_3 = \phi(X_1, X_2)$ separates classes?

[THE *kernel trick*]

The previous optimization method applies to an (*augmented*) space with separable classes

- How do we determine this *augmented* space?

How to find $\phi(\cdot)$ so that $X_3 = \phi(X_1, X_2)$ separates classes?

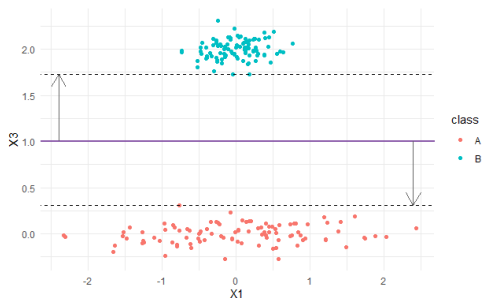


[THE *kernel trick*]

The previous optimization method applies to an (*augmented*) space with separable classes

- How do we determine this *augmented* space?

How to find $\phi(\cdot)$ so that $X_3 = \phi(X_1, X_2)$ separates classes?



↪ The *kernel trick*!

[THE *kernel trick*]

Finding $\phi(\cdot)$ is difficult and sometimes infeasible!

[THE *kernel trick*]

Finding $\phi(\cdot)$ is difficult and sometimes infeasible!

- One can rewrite the optimization problem

[THE *kernel trick*]

Finding $\phi(\cdot)$ is difficult and sometimes infeasible!

- One can rewrite the optimization problem
- ↪ Transform the X s through kernels $K(\cdot)$

[THE *kernel trick*]

Finding $\phi(\cdot)$ is difficult and sometimes infeasible!

- ▶ One can rewrite the optimization problem
- \hookrightarrow Transform the X s through kernels $K(\cdot)$
- ▶ A *kernel* quantifies the similarity of two observations (*inner product*)

[THE *kernel trick*]

Finding $\phi(\cdot)$ is difficult and sometimes infeasible!

- ▶ One can rewrite the optimization problem
 \hookrightarrow Transform the X s through kernels $K(\cdot)$
- ▶ A *kernel* quantifies the similarity of two observations (*inner product*)
- ▶ The optimization can be written **only** in terms of $K(X_i, X_j)$!

[THE *kernel trick*]

Finding $\phi(\cdot)$ is difficult and sometimes infeasible!

- ▶ One can rewrite the optimization problem
 - ↪ Transform the X s through kernels $K(\cdot)$
- ▶ A *kernel* quantifies the similarity of two observations (*inner product*)
- ▶ The optimization can be written **only** in terms of $K(X_i, X_j)$!
 - ↪ The solution is non linear in the original space

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

- ▶ The slack parameter

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

- ▶ The slack parameter
- ▶ The type of kernel (*linear, radial, sigmoid, ..*)

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

- ▶ The slack parameter
- ▶ The type of kernel (*linear, radial, sigmoid, ..*)
- ▶ In an ML framework, one trains the SVM on several samples

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

- ▶ The slack parameter
- ▶ The type of kernel (*linear, radial, sigmoid, ..*)
- ▶ In an ML framework, one trains the SVM on several samples
 - ▶ Number of CV sample should be reduced (100)

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

- ▶ The slack parameter
- ▶ The type of kernel (*linear, radial, sigmoid, ..*)
- ▶ In an ML framework, one trains the SVM on several samples
 - ▶ Number of CV sample should be reduced (100)
 - ▶ Linear kernel are less computationally intensive

[HYPERPARAMETERS IN SVM]

Several hyperparameters have to be selected:

- ▶ The slack parameter
- ▶ The type of kernel (*linear, radial, sigmoid, ..*)
- ▶ In an ML framework, one trains the SVM on several samples
 - ▶ Number of CV sample should be reduced (100)
 - ▶ Linear kernel are less computationally intensive
- ▶ Proceed with parsimony and increase complexity

Introduction
○

The VC dimension
○○○○○

Dimension augmentation
○○○○

Soft margin
○○○

Kernel trick
○○○●

Wrap-up
○

[QUIZ TIME]

Introduction
○

The VC dimension
○○○○○

Dimension augmentation
○○○○

Soft margin
○○○

Kernel trick
○○○●

Wrap-up
○

[QUIZ TIME]

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
- ↪ Soft margins with a *slack* parameter

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
- ↪ Soft margins with a *slack* parameter
 - ▶ A trick to avoid costly features transformations

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
 - ↪ Soft margins with a *slack* parameter
 - ▶ A trick to avoid costly features transformations
 - ↪ The *Kernel Trick*

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
 - ↪ Soft margins with a *slack* parameter
 - ▶ A trick to avoid costly features transformations
 - ↪ The *Kernel Trick*
- ▶ Due to the VC theory, SVM use simple models

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
 - ↪ Soft margins with a *slack* parameter
 - ▶ A trick to avoid costly features transformations
 - ↪ The *Kernel Trick*
- ▶ Due to the VC theory, SVM use simple models
- ▶ SVM requires optimization and CPU

[TAKEAWAYS]

SVM is a complex but popular technique

- ▶ Based on several clever ideas:
 - ▶ Dimension augmentation
 - ▶ Maximal margin classifier
 - ↪ Soft margins with a *slack* parameter
 - ▶ A trick to avoid costly features transformations
 - ↪ The *Kernel Trick*
- ▶ Due to the VC theory, SVM use simple models
- ▶ SVM requires optimization and CPU
- ▶ Implemented in many software!