

Machine Learning



Amin Golzari Oskouei

a.golzari@azaruniv.ac.ir

a.golzari@tabrizu.ac.ir

<https://github.com/Amin-Golzari-Oskouei>

Azərbaycan Şahid Mədani Universiteti
2023

Support Vector Machines : SVM



Table of Contents

2

- ❑ **Motivation** : Optimal decision boundary
- ❑ **Basic concepts** : Support vectors and margin maximization.
- ❑ **Objective function** : the primal problem, and the dual problem.
- ❑ **Linear and Non-linear classification** : Soft margin
- ❑ **Non-linear classification** : Kernel trick
- ❑ **Multi-Class classification** : Multi-Class Support Vector Machine

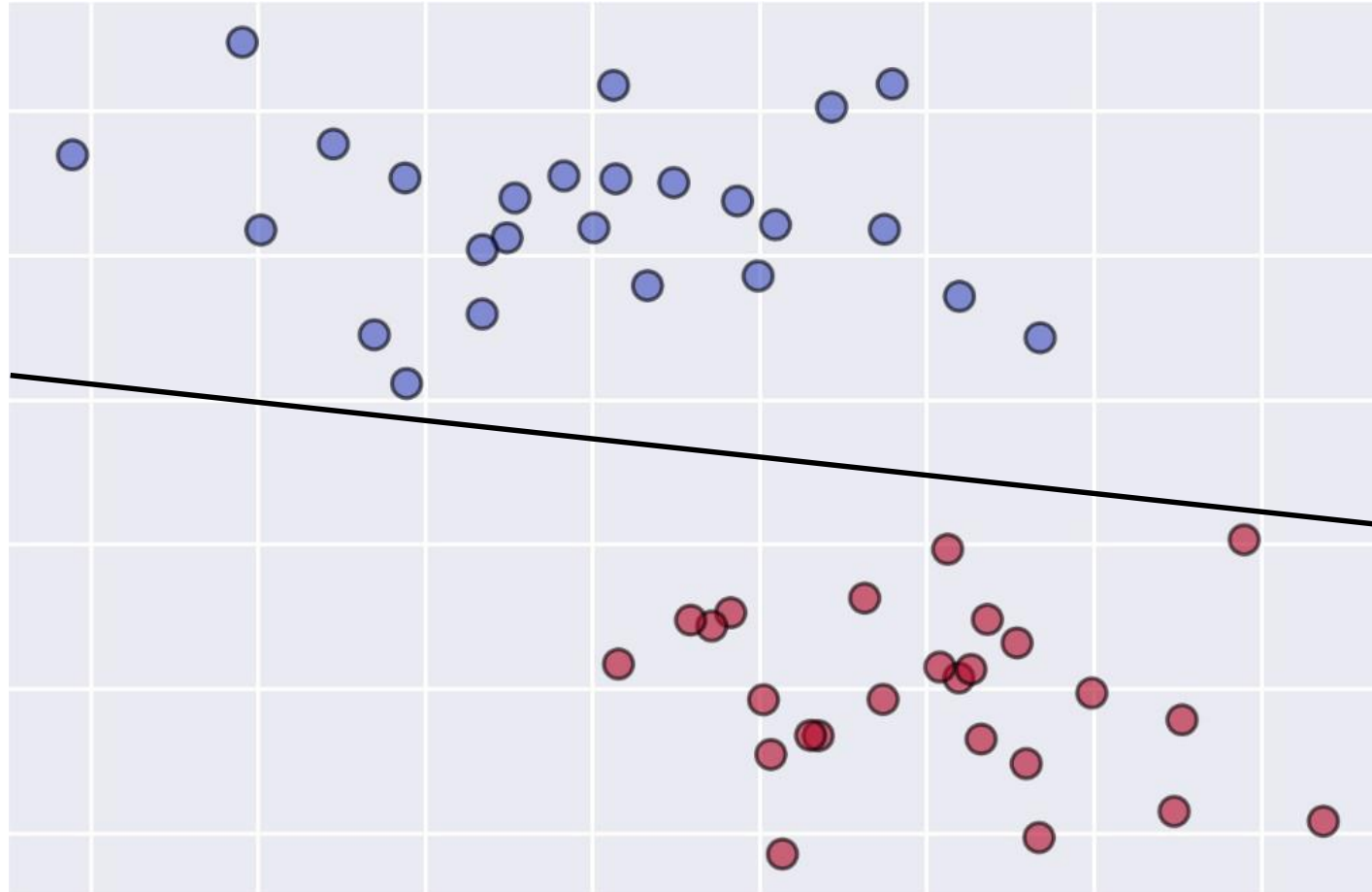
Introduction

3

- ❑ Support Vector Machines [Vanpik, 1990's]
 - ❑ One of the most popular machine learning algorithms!
 - ❑ Better data separation compared to other machine learning methods (classification problems)!
 - ❑ Relatively easy to use!
 - ❑ Using kernel tricks :
 - classification, regression, distribution estimation, one-class classification, and more.

Motive: Linearly separable data

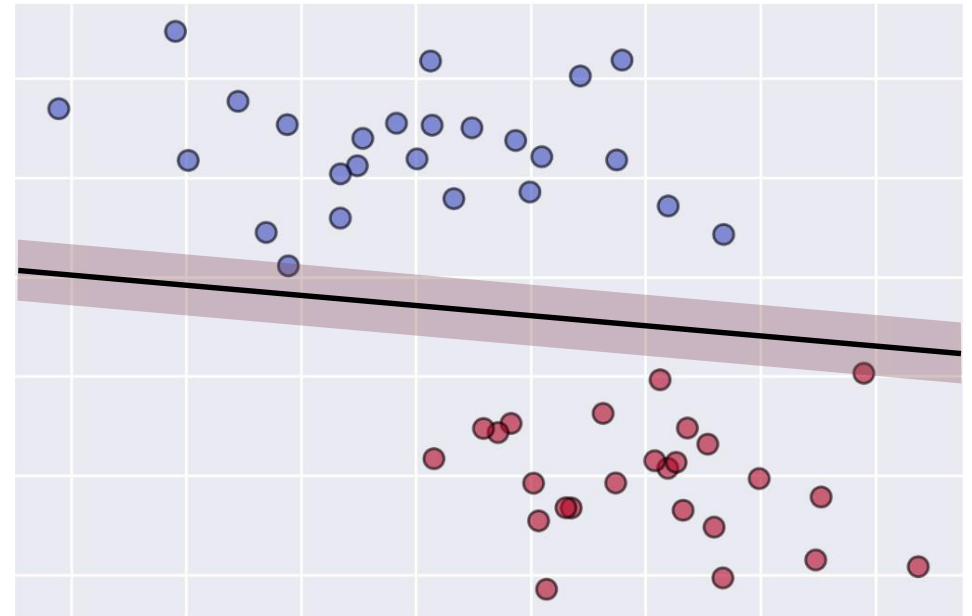
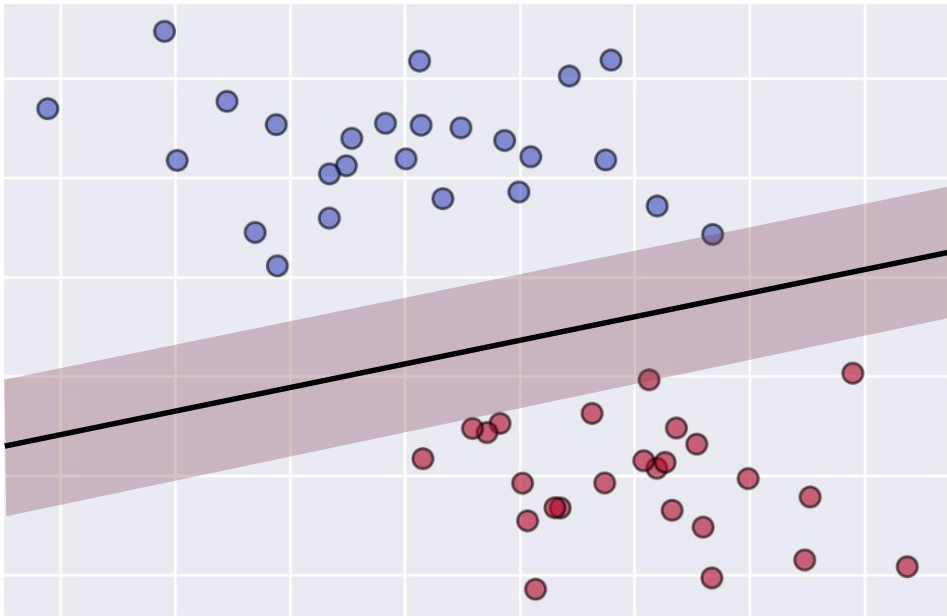
4



Motive: Optimal decision boundary

5

❑ Question : Which decision boundary is better?

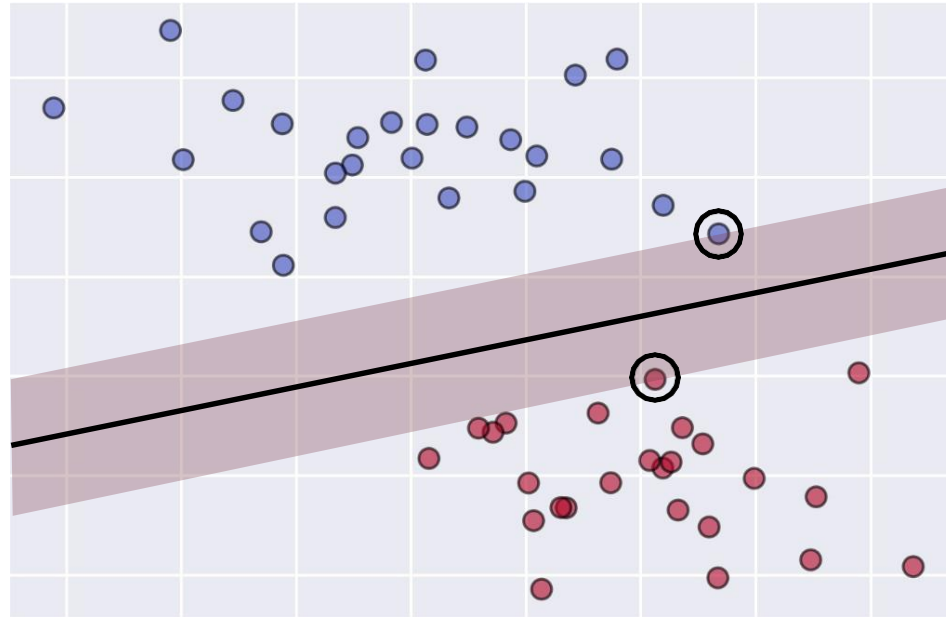


❑ Solution for maximum margin : Maximum robustness against data corruption. [Increasing generalization capability]

Motive: Support Vectors

6

- **Support Vector** : Nearest data points to the decision boundary.

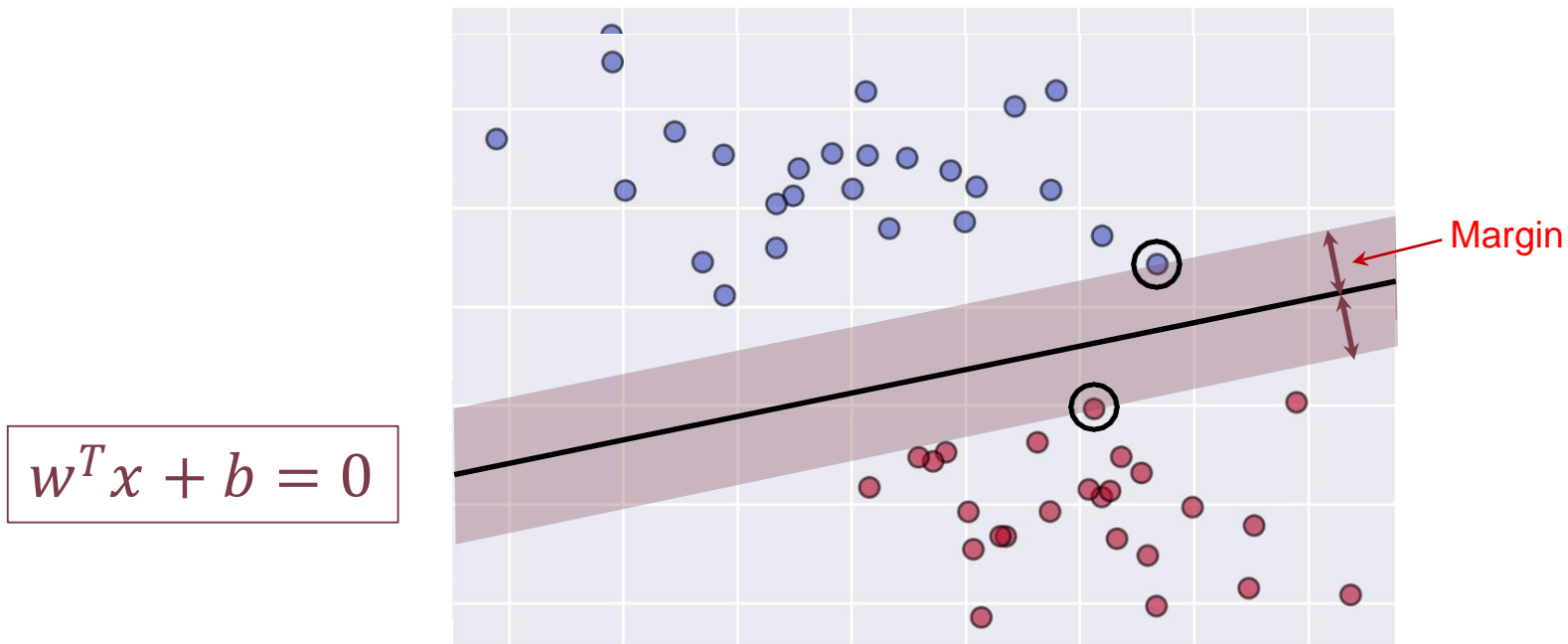


- **Objective**: Maximizing the distance between support vectors and the decision boundary.

Motive: Support Vectors

7

- **Support Vector** : Distance of support vectors from decision boundary.



- **Objective**: Maximizing the distance between support vectors and the decision boundary.

Optimal Decision Boundary: Symbols

8

- Training examples

$$X = (x^t, y^t), \quad y^t = \begin{cases} +1 & \text{if } x^t \in C_1 \\ -1 & \text{if } x^t \in C_2 \end{cases}$$

- Objective: Finding the vector w and the value b in a way that :

$$w^T x^t + b \geq +1 \quad \text{for } y^t = +1$$

$$w^T x^t + b \leq -1 \quad \text{for } y^t = -1$$



$$y^t(w^T x^t + b) \geq +1$$

$$\max(0, 1 - y^t(w^T x^t + b))$$

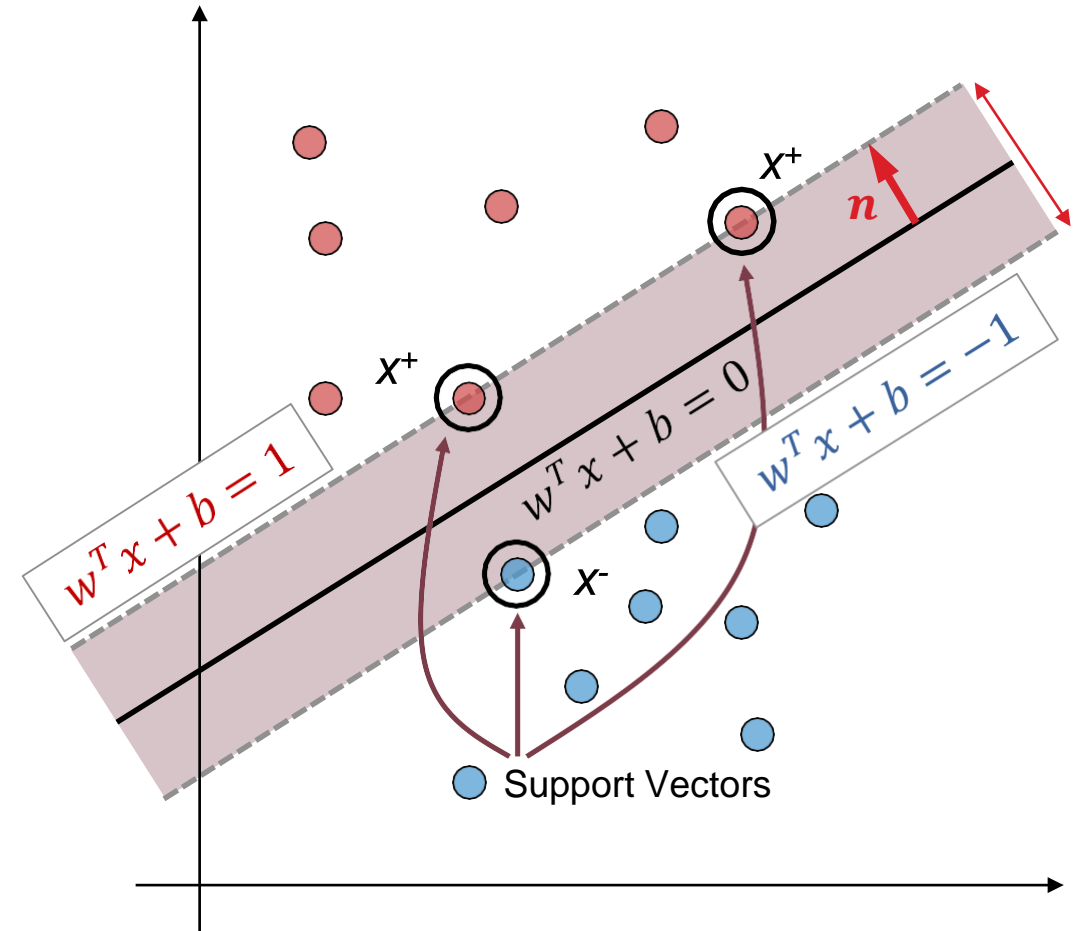
Objective function : Calculate margin

9

□ We already know :

$$w^T x^+ + b = +1$$

$$w^T x^- + b = -1$$



Objective function : Calculate margin

10

□ We already know :

$$w^T x^+ + b = +1$$

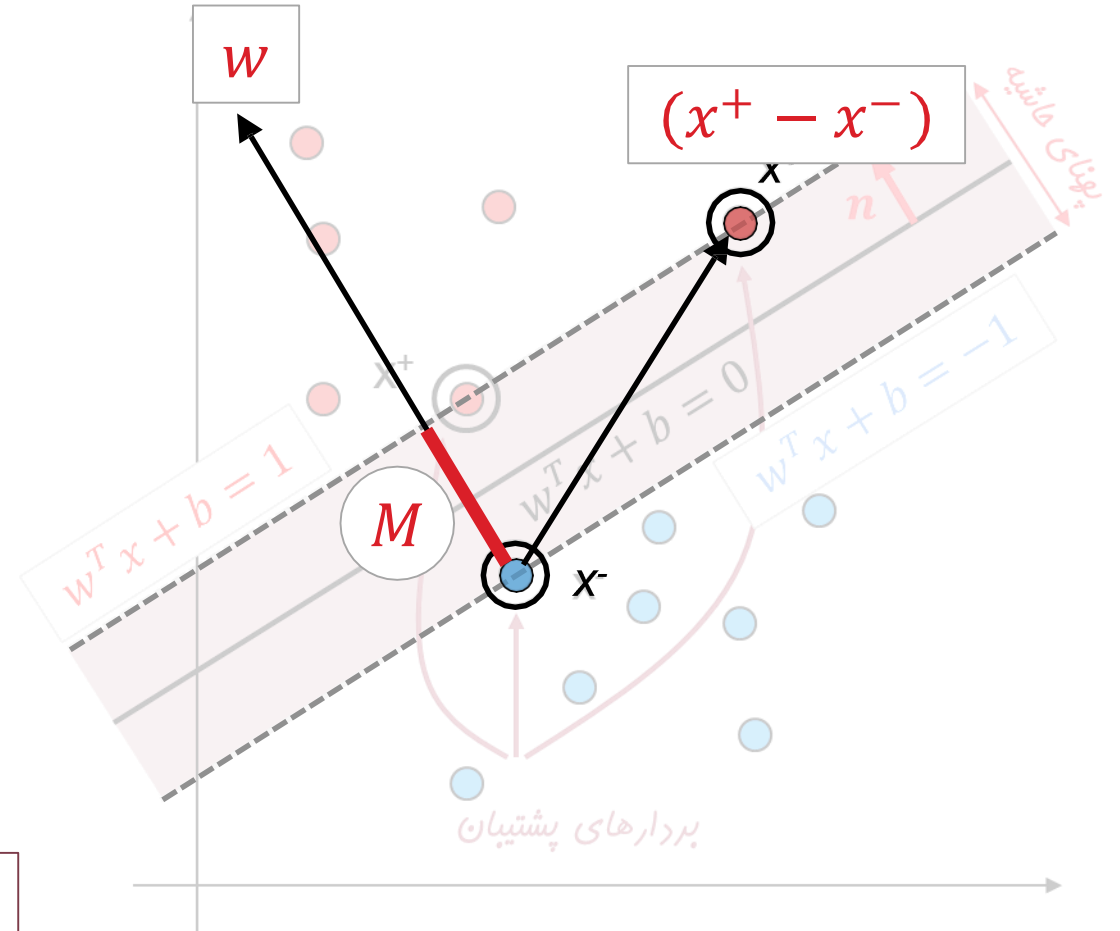
$$w^T x^- + b = -1$$

□ Then :

$$w^T (x^+ - x^-) = 2$$

$$\Rightarrow \|w\| \cdot \|x^+ - x^-\| \cos \alpha = 2$$

$$\Rightarrow \|w\| \cdot M = 2 \Rightarrow \boxed{M = \frac{2}{\|w\|}}$$



Objective function

11

- ❑ **Objective** : Maximizing the margin size [the distance between support vectors and the decision boundary].

$$M = \frac{2}{\|w\|}$$

- ❑ **Attention** : To maximize the margin, one can minimize the size of the vector w .
- ❑ **Constraints**: The decision boundary must correctly separate the data points of both classes.

Objective function : Formal expression

12

□ Objective function :

$$\begin{aligned} & \min \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & (\mathbf{w}^T \mathbf{x}^t + b) \geq +1 \quad \text{if } y^t = +1 \\ & (\mathbf{w}^T \mathbf{x}^t + b) \leq -1 \quad \text{if } y^t = -1 \end{aligned}$$

□ Simplifying :

$$\begin{aligned} & \min \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^t (\mathbf{w}^T \mathbf{x}^t + b) \geq +1 \end{aligned}$$

Objective function : Formal expression

13

□ Objective function :

$$\begin{aligned} & \min \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^t (\mathbf{w}^T \mathbf{x}^t + b) \geq +1 \end{aligned}$$

□ Solving the problem using Lagrange Multipliers :



$$\begin{aligned} L_p &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^m \alpha^t [y^t (\mathbf{w}^T \mathbf{x}^t + b) - 1] \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^m \alpha^t y^t (\mathbf{w}^T \mathbf{x}^t + b) + \sum_{t=1}^m \alpha^t \end{aligned}$$

Objective function : Formal expression

14

□ Objective function :

$$\begin{aligned} L_p &= \frac{1}{2} \|w\|^2 - \sum_{t=1}^m \alpha^t [y^t (w^T x^t + b) - 1] \\ &= \frac{1}{2} \|w\|^2 - \sum_{t=1}^m \alpha^t y^t (w^T x^t + b) + \sum_{t=1}^m \alpha^t \end{aligned}$$

The decision boundary is a linear combination of the training data.

$$\begin{aligned} \frac{\partial L_p}{\partial w} = 0 &\Rightarrow w = \sum_{t=1}^m \alpha^t y^t x^t \\ \frac{\partial L_p}{\partial b} = 0 &\Rightarrow \sum_{t=1}^m \alpha^t y^t = 0 \end{aligned}$$

Objective function : Dual form

15

□ Objective function :

$$\begin{aligned} L_d &= \frac{1}{2} (w^T w) - w^T \sum_{t=1}^m \alpha^t y^t x^t - b \sum_{z=1}^m \alpha^t y^t + \sum_{t=1}^m \alpha^t \\ &= \frac{1}{2} (w^T w) + \sum_{t=1}^m \alpha^t \\ &= \frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s (x^t)^T x^s + \sum_{t=1}^m \alpha^t \end{aligned}$$

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $\alpha^t \geq 0 \forall t$

- Many of the alpha coefficients are equal to zero, and only a few have values greater than zero.
- The data points for which the alpha values are greater than zero are **the support vectors**.

Objective function : Vector form

16

□ Objective function :

$$L_d = \frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m \alpha^t \alpha^s y^t y^s (x^t)^T x^s + \sum_{t=1}^m \alpha^t$$
$$= -\frac{1}{2} \alpha^T Q \alpha + e^T \alpha$$

$$Q_{ts} = y^t y^s (x^t)^T x^s, \quad e = [1 \quad 1 \quad \dots \quad 1]^T \in \mathbb{R}^m$$

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $\alpha^t \geq 0 \forall t$

- Many of the alpha coefficients are equal to zero, and only a few have values greater than zero.
- The data points for which the alpha values are greater than zero are **the support vectors**.

Non-linearly Separable Data: Soft Margin

17

- ❑ If the data is not linearly separable, what happens?
- ❑ **Soft Margin:** Allowing for a small margin of error in separation
- ❑ **Soft Error**
- ❑ **New objective function**

$$y^t(\mathbf{w}^T \mathbf{x}^t + b) \geq 1 - \varepsilon^t$$

$$\text{soft error} : \sum_{t=1}^m \varepsilon^t$$

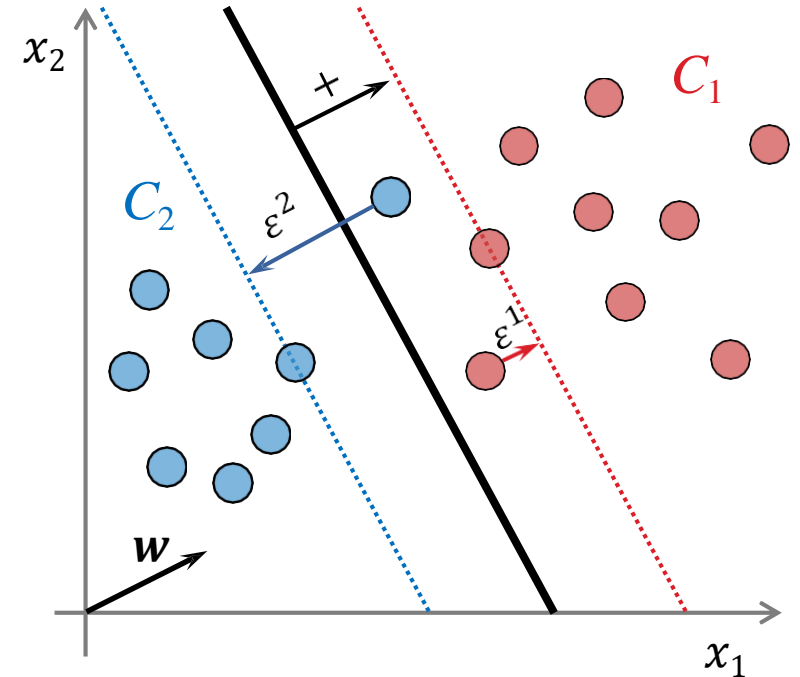
Penalty coefficient

$$\begin{aligned} & \min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^m \varepsilon^t \\ \text{s.t.} \quad & y^t(\mathbf{w}^T \mathbf{x}^t + b) \geq 1 - \varepsilon^t \\ & \varepsilon^t \geq 0 \end{aligned}$$

Non-linearly Separable Data: Soft Margin

18

$$\begin{aligned} \min & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^m \varepsilon^t \\ \text{s.t.} & y^t (\mathbf{w}^T \mathbf{x}^t + b) \geq 1 - \varepsilon^t \\ & \varepsilon^t \geq 0 \end{aligned}$$



$$\begin{aligned} L_p = & \frac{1}{2} \|\mathbf{w}\|^2 + c \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t [y^t (\mathbf{w}^T \mathbf{x}^t + b) - 1 + \varepsilon^t] \\ & - \sum_{t=1}^m \mu^t \varepsilon^t \end{aligned}$$

Non-linearly Separable Data: Soft Margin

19

$$L_p = \frac{1}{2} \|w\|^2 + c \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t [y^t (w^T x^t + b) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

$$\frac{\partial L_p}{\partial w} = 0 \Rightarrow w = \sum_{t=1}^m \alpha^t y^t x^t$$

$$\frac{\partial L_p}{\partial b} = 0 \Rightarrow \sum_{t=1}^m \alpha^t y^t = 0$$

$$\frac{\partial L_p}{\partial \varepsilon^t} = 0 \Rightarrow C - \alpha^t - \mu^t = 0 \Rightarrow 0 \leq \alpha^t \leq C$$

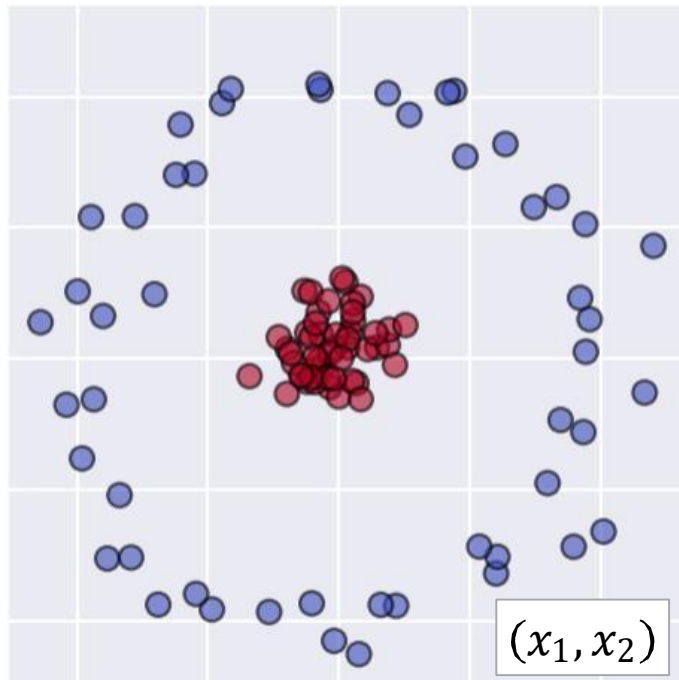
A horizontal decorative bar at the top of the slide, consisting of a red rectangular section on the left and a teal rectangular section on the right.

Kernel Trick and Non-Linear Classification

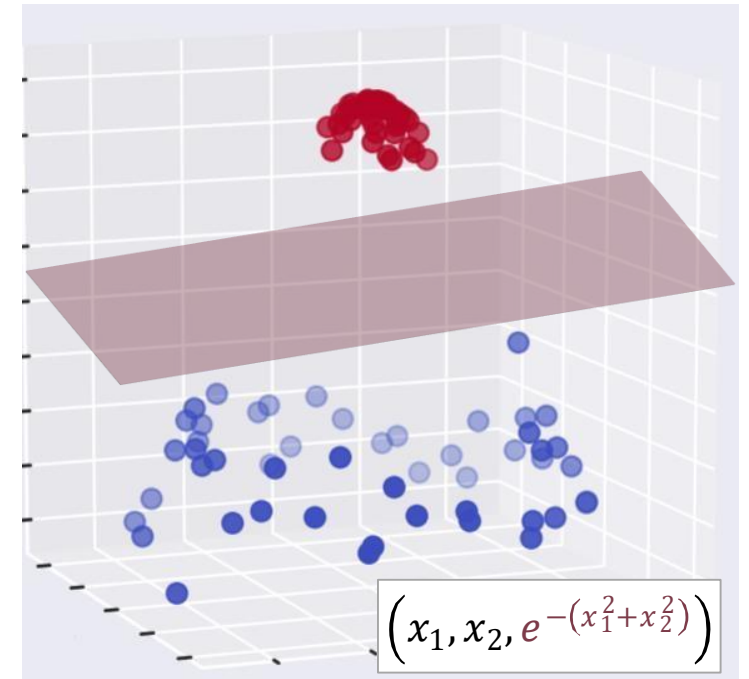
Kernel Functions

21

- ❑ **Idea** : Mapping the problem to a new feature space using non-linear transformations
 - ❑ Using a linear model in the new feature space for data classification.



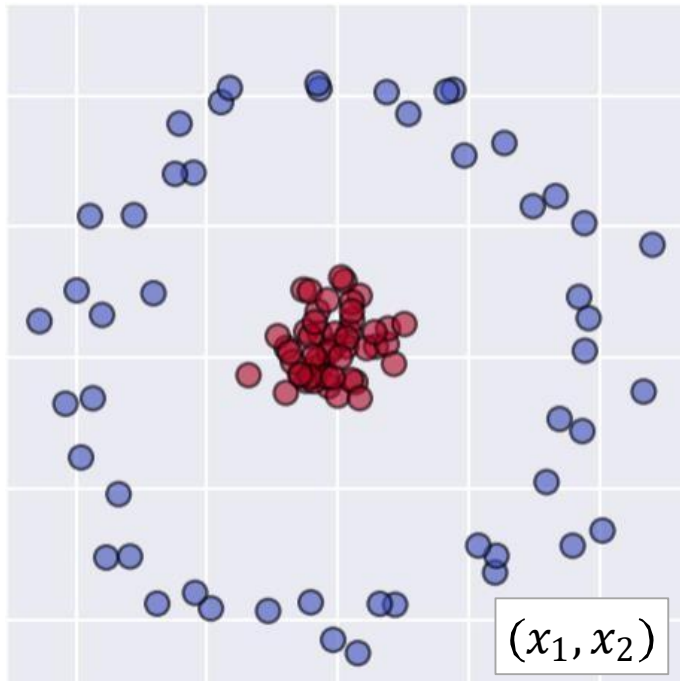
$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$



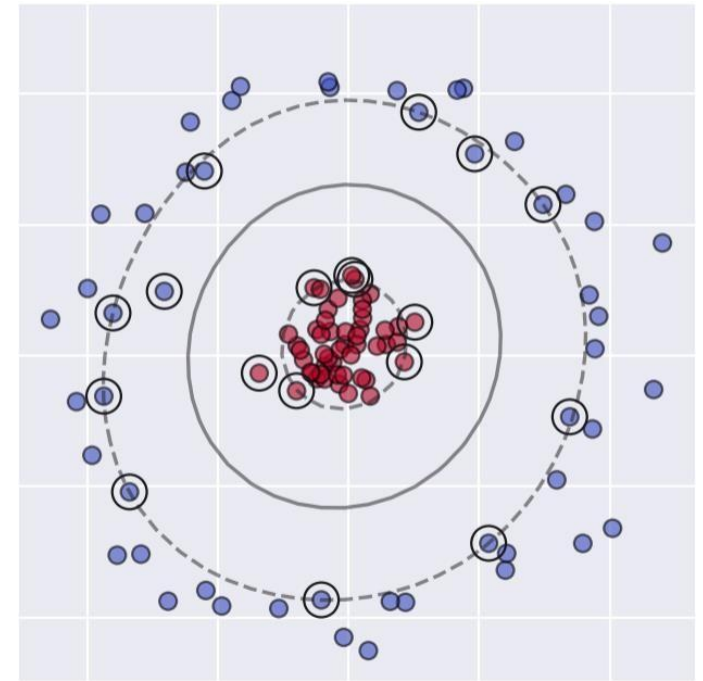
Kernel Functions

22

- ❑ **Idea** : Mapping the problem to a new feature space using non-linear transformations
 - ❑ Using a linear model in the new feature space for data classification.
 - ❑ A linear model in the new feature space corresponds to a non-linear model in the original space.



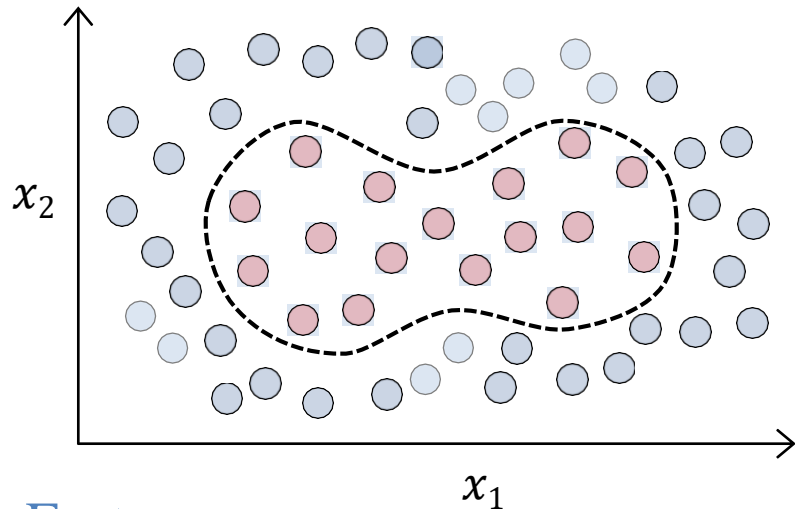
$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$



Non-Linear Decision Boundary

23

❑ Prediction : $y = 1$ if :



$$h(x) = b + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + \dots \geq 0$$

❑ Features :

$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1^2, \quad f_4 = x_2^2, \quad f_5 = x_1x_2, \quad \dots$$

$$h(f) = b + w_1f_1 + w_2f_2 + w_3f_3 + w_4f_4 + w_5f_5 + \dots$$

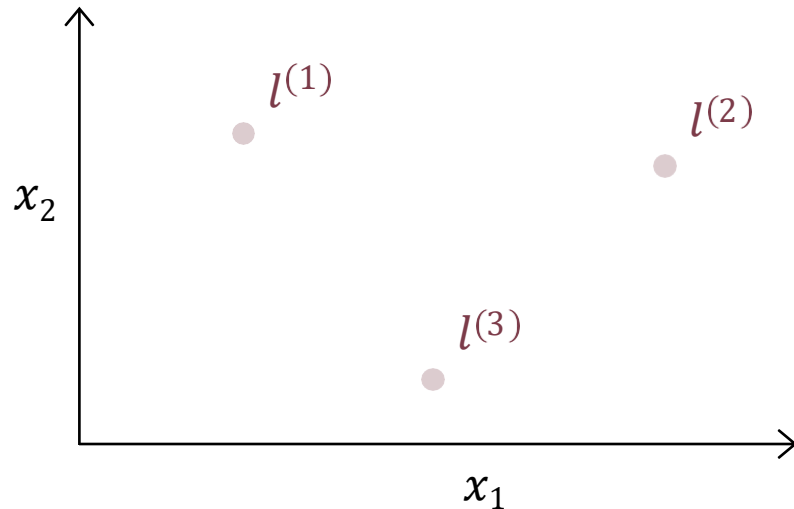
← Linear decision boundary

❑ Question : Is there a better method for feature selection?

Kernels as Similarity Measures

24

- ❑ **Idea** : Given x , select a new set of features based on its similarity to reference points using l^2 , l^1 , and l^3 norms.



$$f_1 = \text{sim}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$f_2 = \text{sim}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{sim}(x, l^{(3)}) = \exp\left(-\frac{\|x - l^{(3)}\|^2}{2\sigma^2}\right)$$

Kernel (Gaussian Kernel)

- ❑ **Kernel Function**: A measure for calculating the similarity between data x and y .

Kernels as Similarity Measures

25

□ Kernel function :

$$f_i = \text{sim}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$

□ First case : $x \approx l^{(i)}$

$$f_i \approx \exp\left(-\frac{0}{2\sigma^2}\right) = \exp(0) = 1$$

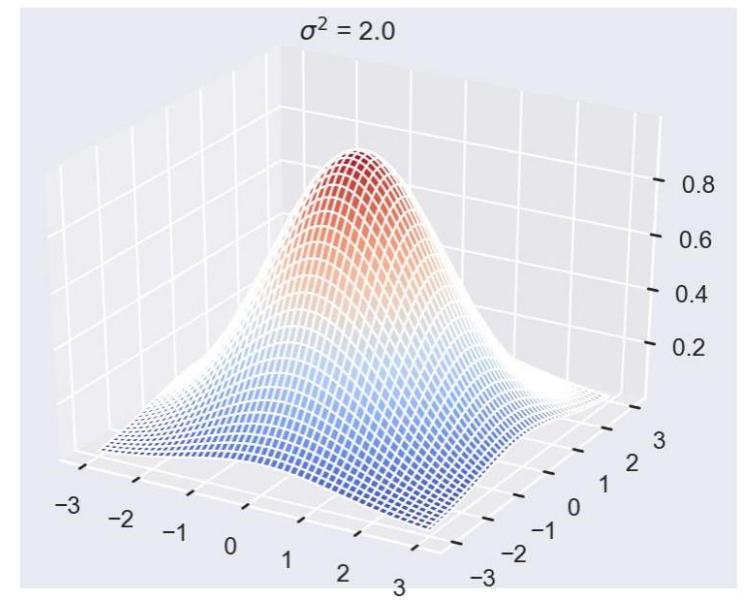
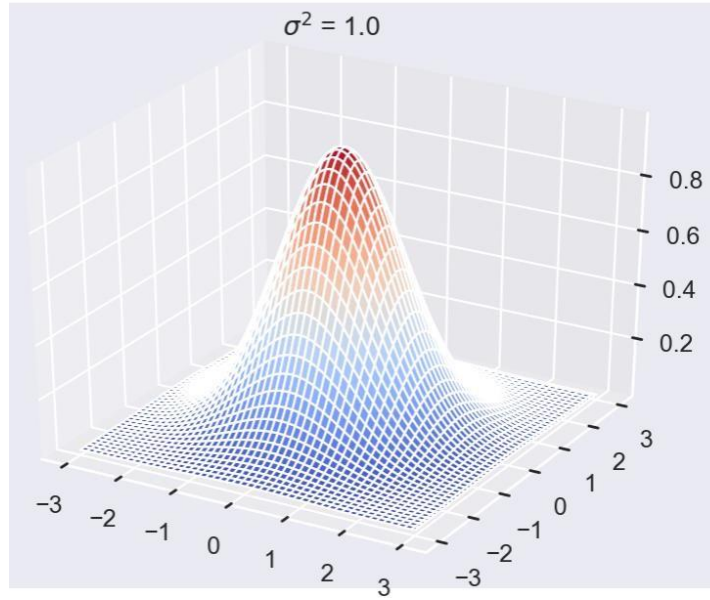
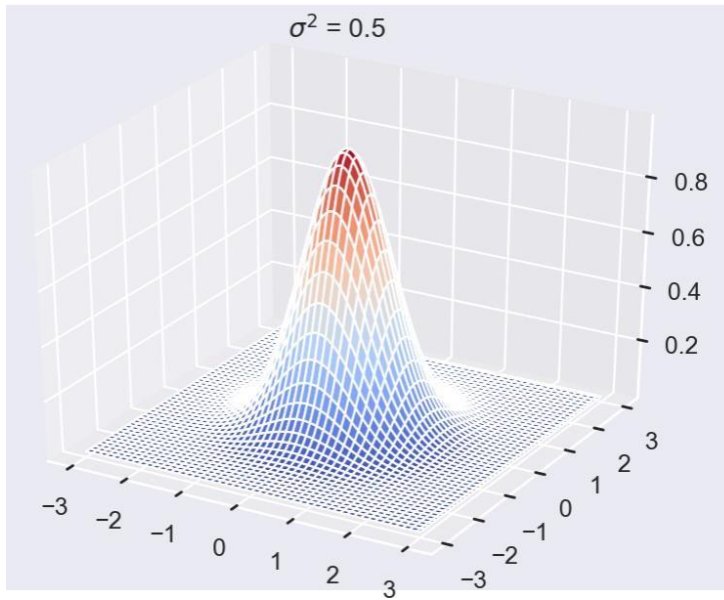
□ Second case : x is very far from $l^{(i)}$.

$$f_i \approx \exp\left(-\frac{\infty}{2\sigma^2}\right) = \exp(-\infty) = 0$$

Kernels as Similarity Measures

26

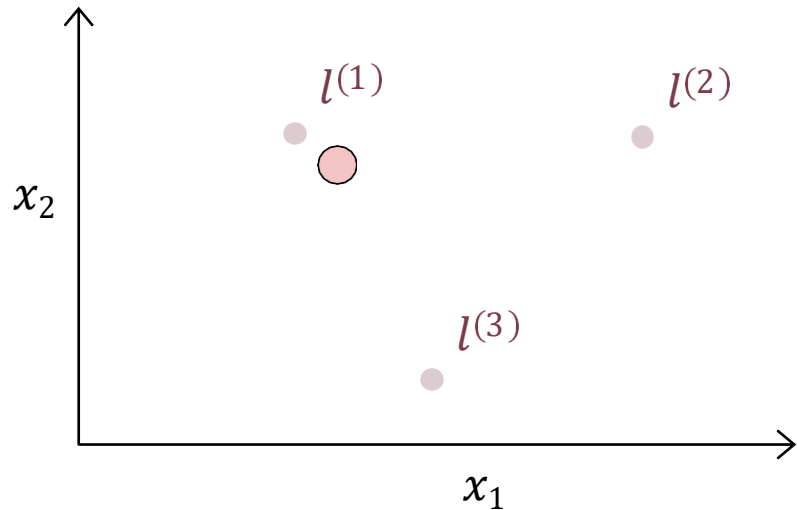
$$f_i = \text{sim}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$



Kernels as Similarity Measures

27

□ Prediction : $y = 1$ if :



$$b + w_1 f_1 + w_2 f_2 + w_3 f_3 \geq 0$$

$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

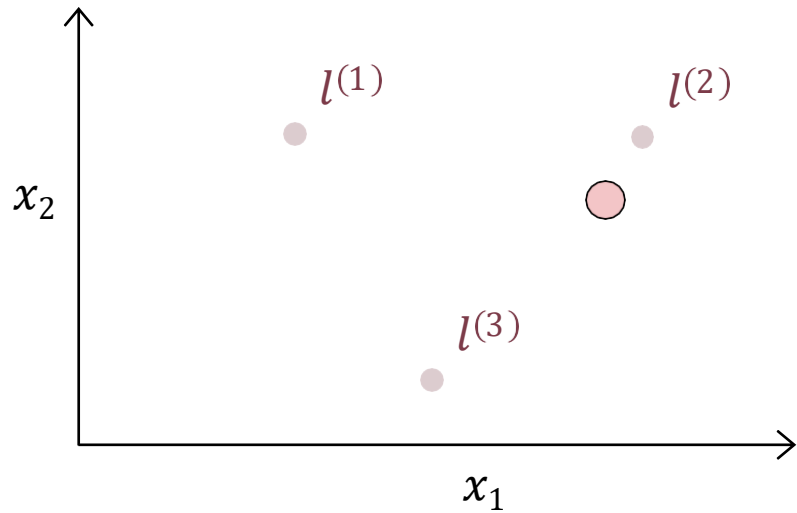
$$f_1 \approx 1, f_2 \approx f_3 \approx 0$$

$$h(f) \approx -0.5 + (1.0)(1.0) + (1.0)(0.0) + (0.0)(0.0) = 0.5 \geq 0 \Rightarrow \boxed{y = 1}$$

Kernels as Similarity Measures

28

□ Prediction : $y = 1$ if :



$$b + w_1 f_1 + w_2 f_2 + w_3 f_3 \geq 0$$

$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

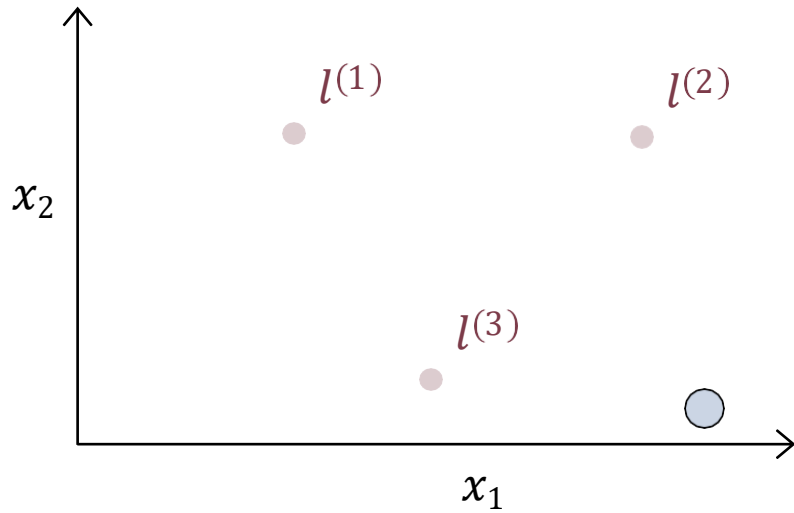
$$f_1 \approx f_3 \approx 0, f_2 \approx 1$$

$$h(f) \approx -0.5 + (1.0)(0.0) + (1.0)(1.0) + (0.0)(0.0) = 0.5 \geq 0 \Rightarrow \boxed{y = 1}$$

Kernels as Similarity Measures

29

□ Prediction : $y = 1$ if :



$$b + w_1 f_1 + w_2 f_2 + w_3 f_3 \geq 0$$

$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

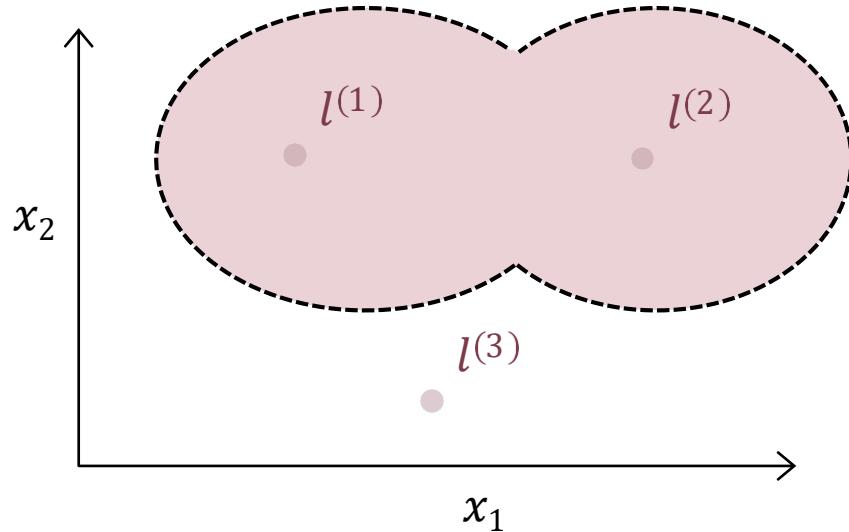
$$f_1 \approx f_2 \approx f_3 \approx 0$$

$$h(f) \approx -0.5 + (1.0)(0.0) + (1.0)(0.0) + (0.0)(0.0) = -0.5 < 0 \Rightarrow \boxed{y = 0}$$

Kernels as Similarity Measures

30

□ Prediction : $y = 1$ if :



$$b + w_1 f_1 + w_2 f_2 + w_3 f_3 \geq 0$$

$$\begin{array}{cccc} \uparrow & \uparrow & \uparrow & \uparrow \\ -0.5 & 1.0 & 1.0 & 0.0 \end{array}$$

□ Decision Boundary : Classifies points near $l^{(1)}$ and $l^{(2)}$ as Class 1, and other points as Class 0.

Several Questions

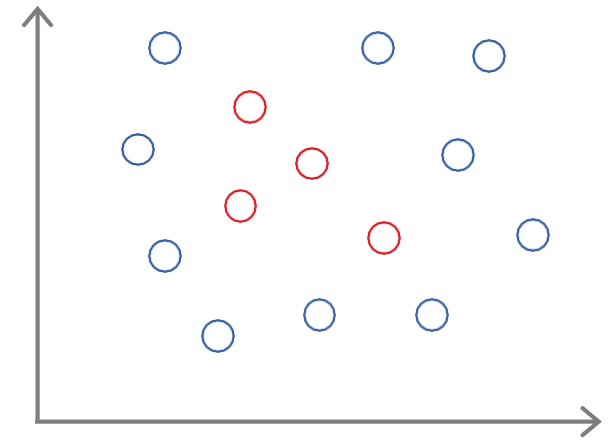
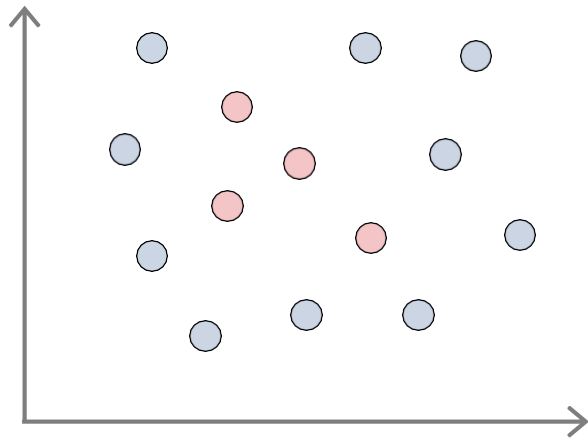
31

- ❑ How does the learning algorithm automatically select reference points?
- ❑ How are appropriate values for kernel function parameters determined?
- ❑ Are there other types of kernels?

Select Reference Points

32

- ❑ How does the learning algorithm automatically select reference points?
 - ❑ For each sample in the training set, a reference point is chosen equal to that sample.



Feature Mapping

33

□ Training set :

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

□ Reference points :

$$l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$$

□ Mapping feature space :

$$x = \begin{bmatrix} x_0 = 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$



$$f = \begin{bmatrix} f_0 = 1 \\ f_1 = K(x, l^{(1)}) \\ f_2 = K(x, l^{(2)}) \\ \vdots \\ f_m = K(x, l^{(m)}) \end{bmatrix}$$

Kernel Trick

34

- Kernel Function : Preprocessing data x using kernel functions.

$$\begin{aligned} \mathbf{z} &= \varphi(\mathbf{x}) \\ &= (\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x}), \dots, \varphi_k(\mathbf{x})) \end{aligned}$$

$$g(\mathbf{z}) = \mathbf{w}^T \mathbf{z} + b$$

$$g(\mathbf{x}) = \mathbf{w}^T \varphi(\mathbf{x}) + b$$

It may be infinite!

- Decision boundary :

$$\mathbf{w} = \sum_{t=1}^m a^t y^t \mathbf{z}^t = \sum_{t=1}^m a^t y^t \varphi(\mathbf{x}^t)$$

- Classification of new data :

$$g(\mathbf{x}) = \mathbf{w}^T \varphi(\mathbf{x}) + b = \left(\sum_{t=1}^m a^t y^t \varphi(\mathbf{x}^t)^T \right) \varphi(\mathbf{x}) + b = \left(\sum_{t=1}^m a^t y^t \varphi(\mathbf{x}^t)^T \varphi(\mathbf{x}) \right) + b = \left(\sum_{t=1}^m a^t y^t k(\mathbf{x}^t, \mathbf{x}) \right) + b$$

Kernel Function

35

$$L_p = \frac{1}{2} \|W\|^2 + C \sum_{t=1}^m \varepsilon^t$$

s.t. $y^t W^t \varphi(x^t) \geq 1 - \varepsilon^t$

$$\varepsilon^t \geq 0$$

$$L_p = \frac{1}{2} \|W\|^2 + C \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t y^t [W^T \varphi(x^t) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

Lagrange Coefficients

Lagrange Coefficients

Kernel Functions: The Primary Problem

36

$$L_p = \frac{1}{2} \|W\|^2 + C \sum_{t=1}^m \varepsilon^t - \sum_{t=1}^m \alpha^t y^t [W^T \varphi(x^t) - 1 + \varepsilon^t] - \sum_{t=1}^m \mu^t \varepsilon^t$$

$$\frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{t=1}^m \alpha^t y^t \varphi(\mathbf{x}^t)$$

$$\frac{\partial L_p}{\partial \varepsilon^t} = 0 \Rightarrow C - \alpha^t - \mu^t = 0 \Rightarrow 0 \leq \alpha^t \leq C$$

Kernel Functions: The Dual Problem

37

$$L_d = -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m a^t a^s y^t y^s \varphi(x^t)^T \varphi(x^s) + \sum_{t=1}^m a^t$$

subject to $\sum_{t=1}^m \alpha^t y^t = 0$ and $0 \leq \alpha^t \leq C \forall t$

□ Idea of Kernel Machines : [Kernel Trick]

- Replacing the dot product of basis functions with a kernel function as $K(x^t, x^s)$

$$L_d = -\frac{1}{2} \sum_{t=1}^m \sum_{s=1}^m a^t a^s y^t y^s K(x^t, x^s) + \sum_{t=1}^m a^t$$

Gram Matrix: A symmetric positive definite matrix (for linear separability).

Kernel Functions: Polynomial Kernel

38

□ **Polynomial Kernel** : A polynomial of degree q

$$K(x^t, x) = (x^T x^t + 1)^q$$

□ **Example** : $[q = 2, d = 2]$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$K(x, y) = (x^T y + 1)^2$$

$$= (x_1 y_1 + x_2 y_2 + 1)^2$$

$$= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2$$

3 multiplications,
2 additions

6 multiplications,
5 additions

$$\varphi(x) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T$$

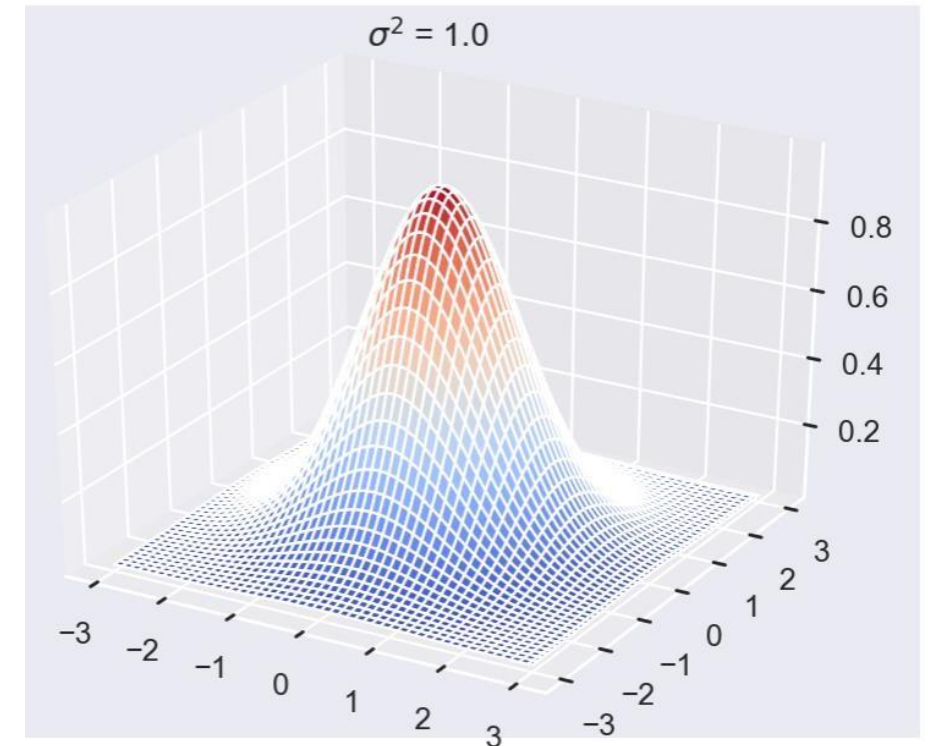
$$\varphi(y) = [1, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_1 y_2, y_1^2, y_2^2]^T$$

Kernel Functions: Gaussian Kernel

39

- Gaussian kernel :
- Finding an appropriate value for σ
 - Using a validation set [Model selection]
 - Larger values : Smoother decision boundary

$$K(\mathbf{x}^t, \mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}^t - \mathbf{x}\|^2}{2\sigma^2}\right)$$

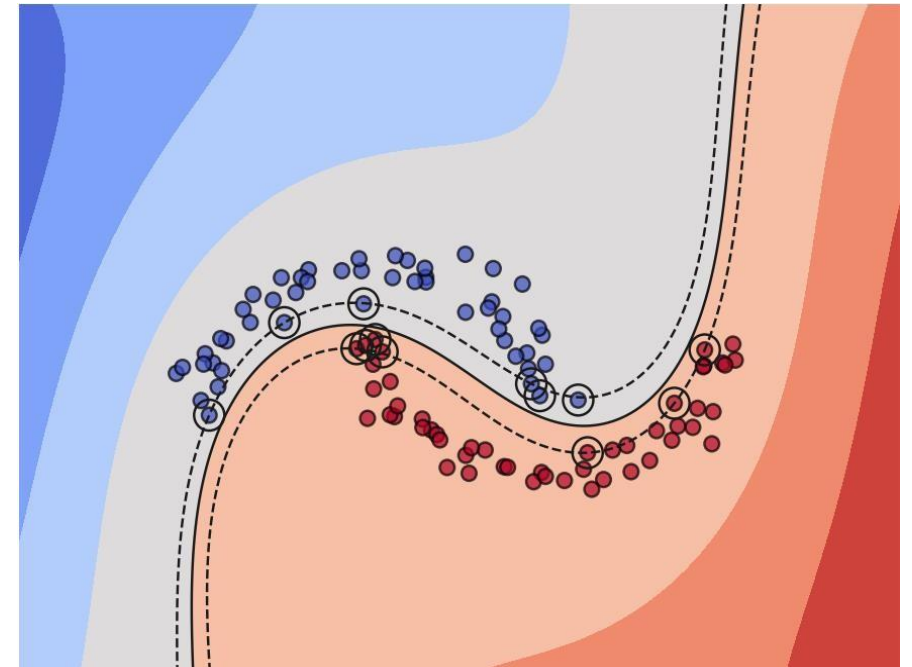


Kernel Functions: Gaussian Kernel

40

- Gaussian kernel :
- Finding an appropriate value for σ
 - Using a validation set [Model selection]
 - Larger values : Smoother decision boundary

$$K(\mathbf{x}^t, \mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}^t - \mathbf{x}\|^2}{2\sigma^2}\right)$$



Hyperparameters: Model Selection

41

- ❑ **First Solution:** [A very bad idea]

- ❑ Selecting a value that results in the highest classification accuracy on the **test dataset**



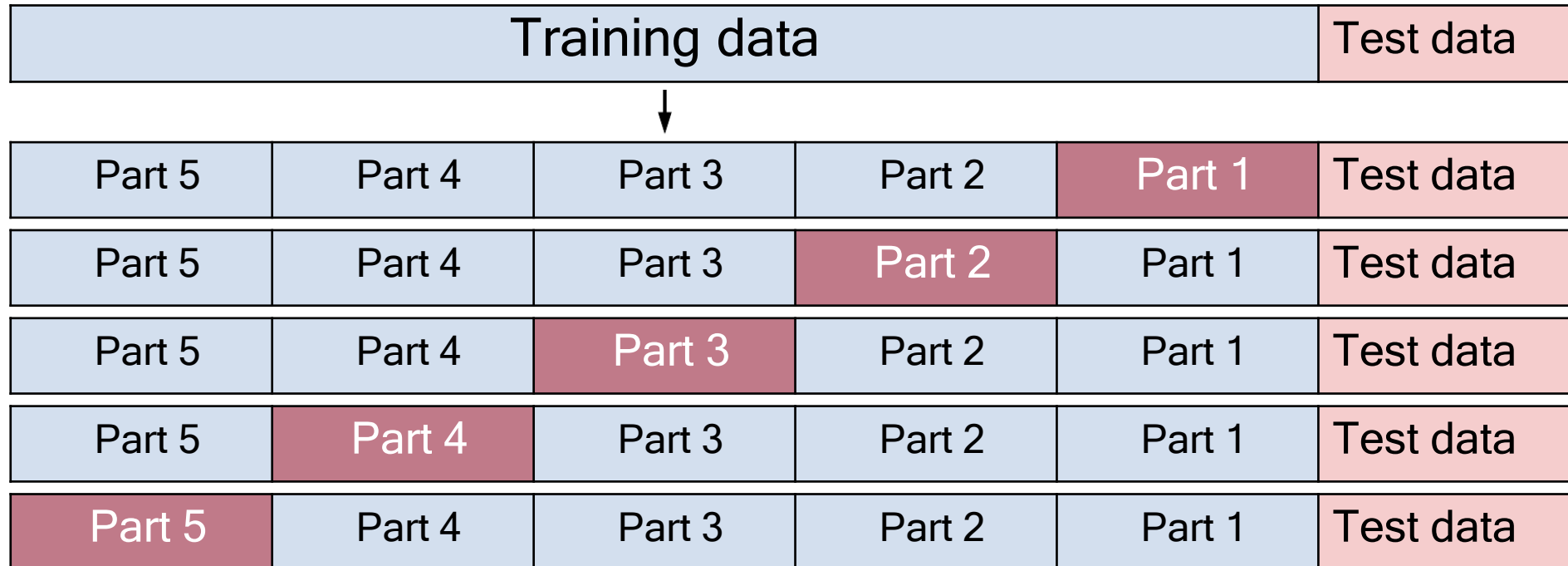
- ❑ **Attention :** [Very important]

- ❑ Use the test set at the end of the process and only for **estimating generalization capability** of the classifier

Hyperparameters: Model Selection

42

❑ Second Solution: [Cross-Validation]

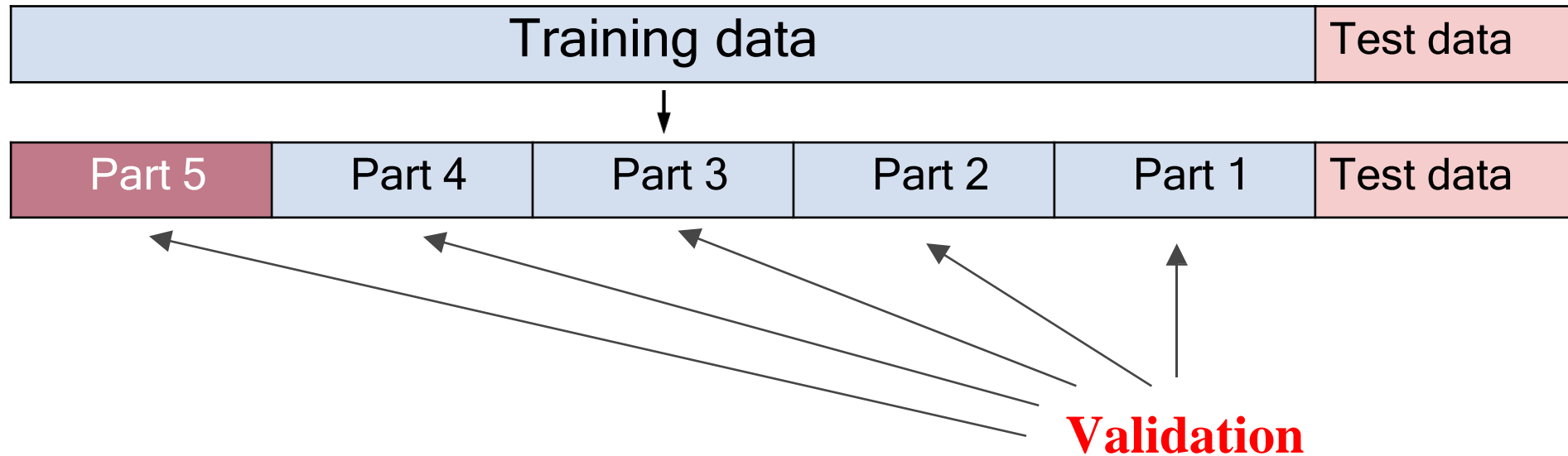


Validation Data [for determining hyperparameter values]

Hyperparameters: Model Selection

43

❑ Second Solution: [Cross-Validation]




- ❑ Select one portion as validation data each time and then take the average of the obtained results.

Hyperparameters: Model Selection

44

Splitting training data into training and validation data.



```
cv = StratifiedKFold(n_splits=5, shuffle=True)

C_range = np.logspace(-3, 5, 9)
gamma_range = np.logspace(-3, 5, 9)


pgrid = dict(gamma=gamma_range, C=C_range)
grid = GridSearchCV(SVC(), param_grid=pgrid, cv=cv)

grid.fit(X, y)
```

Hyperparameters: Model Selection

45

Determining the search range for hyperparameter tuning.



```
cv = StratifiedKFold(n_splits=5, shuffle=True)
```

```
C_range = np.logspace(-3, 5, 9)
```

```
gamma_range = np.logspace(-3, 5, 9)
```

```
pgrid = dict(gamma=gamma_range, C=C_range)
```

```
grid = GridSearchCV(SVC(), param_grid=pgrid, cv=cv)
```

```
grid.fit(X, y)
```

$\{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4, 10^5\}$

Hyperparameters: Model Selection

46

Creating a classifier



```
cv = StratifiedKFold(n_splits=5, shuffle=True)

C_range = np.logspace(-3, 5, 9)
gamma_range = np.logspace(-3, 5, 9)

pgrid = dict(gamma=gamma_range, C=C_range)
grid = GridSearchCV(SVC(), param_grid=pgrid, cv=cv)

grid.fit(X, y)
```

Hyperparameters: Model Selection

47

```
cv = StratifiedKFold(n_splits=5, shuffle=True)

C_range = np.logspace(-3, 5, 9)
gamma_range = np.logspace(-3, 5, 9)

pgrid = dict(gamma=gamma_range, C=C_range)
grid = GridSearchCV(SVC(), param_grid=pgrid, cv=cv)
```

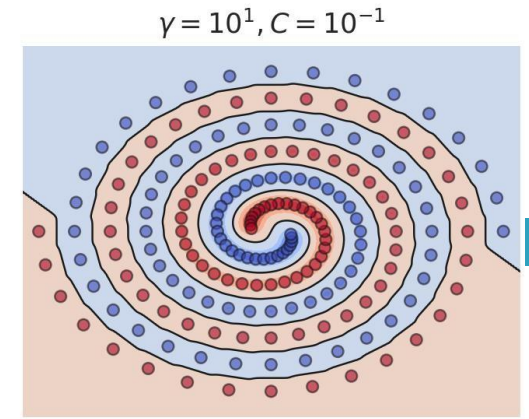
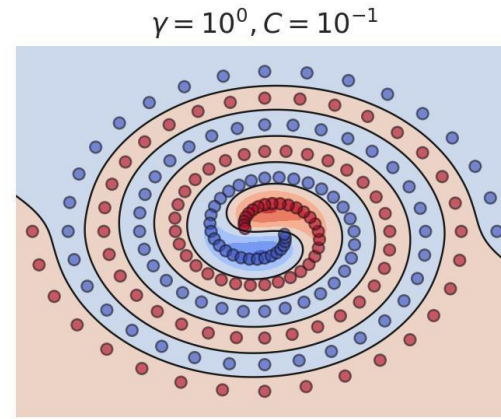
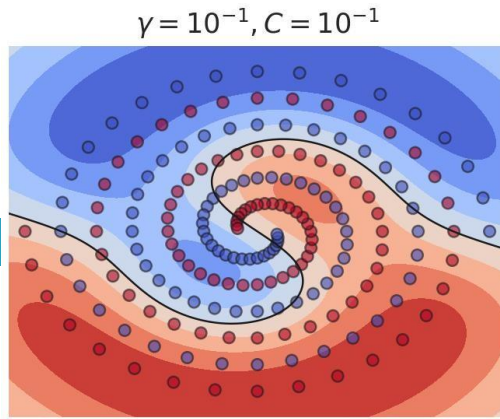
Training the
classifier



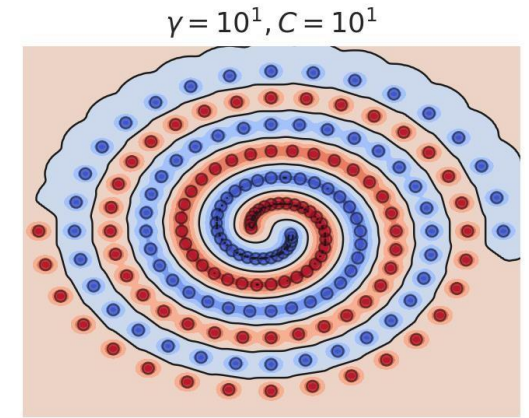
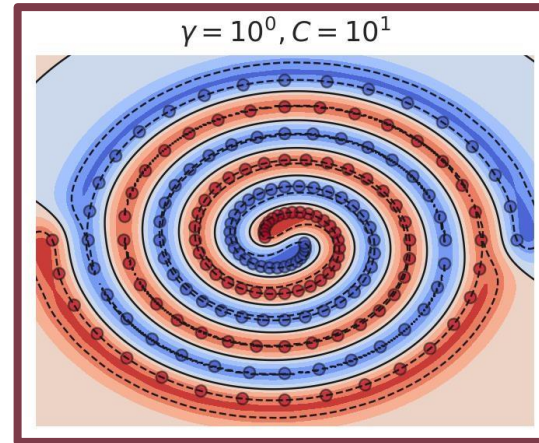
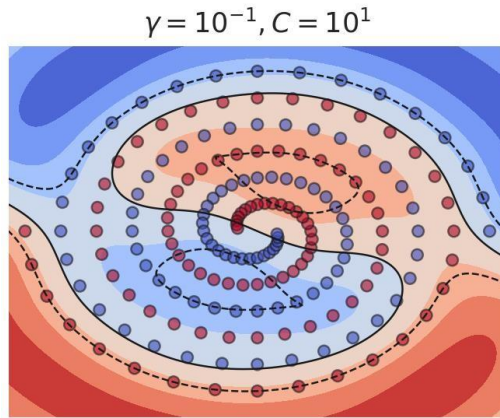
```
grid.fit(X, y)
```


Grid Search

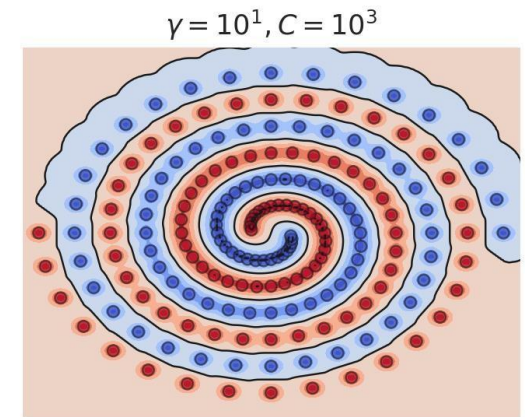
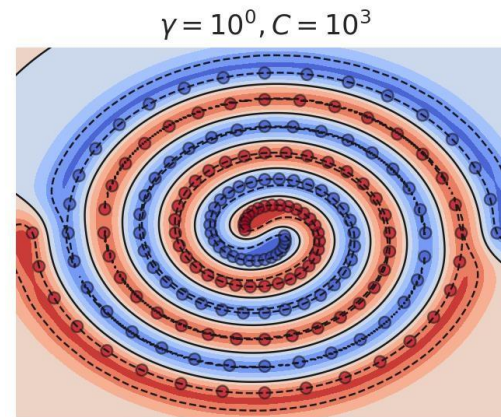
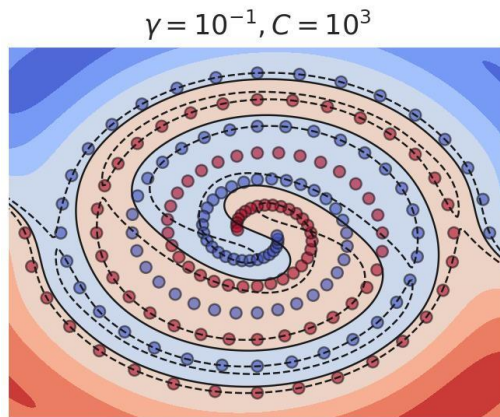
48



$C = 10^{-1}$



$C = 10^1$



$C = 10^3$

Support Vector Machine Parameters

49

❑ **Question :** How are the appropriate values for the kernel function parameters determined?

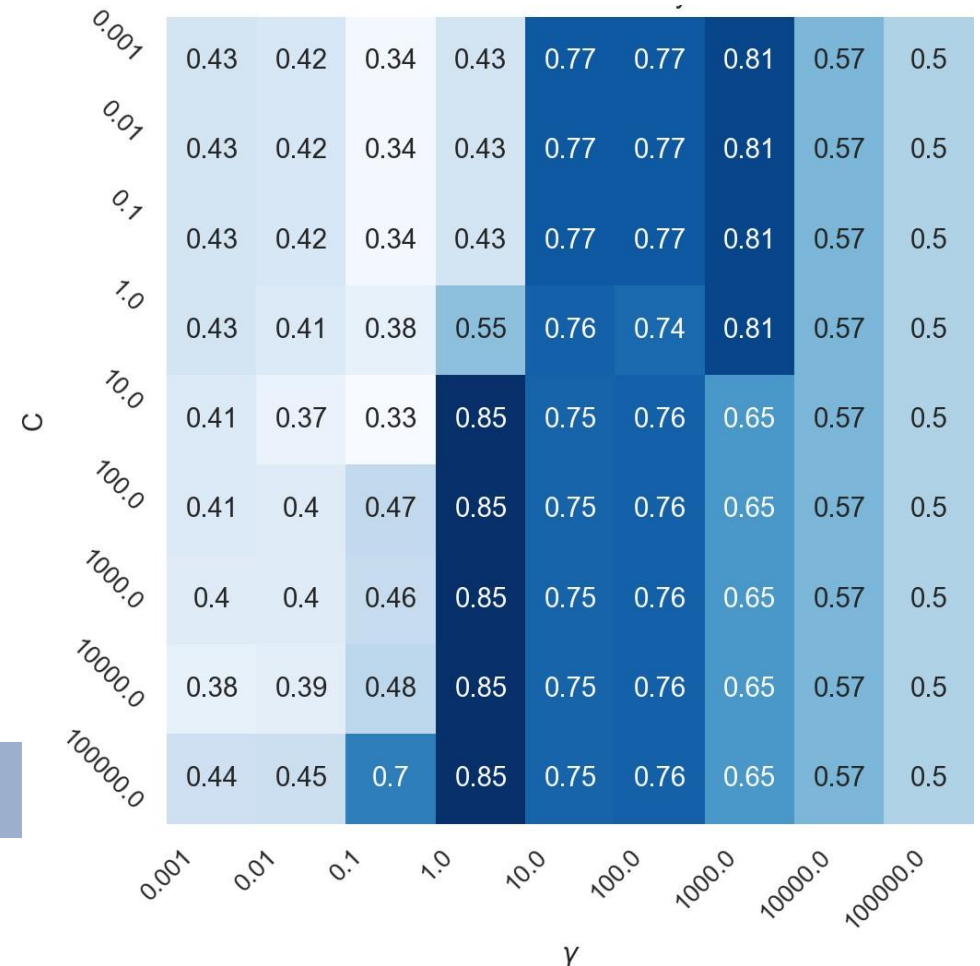
❑ **Parameter C :**

- Smaller values: Higher bias, lower variance.
- Larger values: Lower bias, higher variance.

❑ **Parameter σ :**

- Smaller values: Lower bias, higher variance.
- Larger values: Higher bias, lower variance.

Classification accuracy on validation data



Guide to Using Support Vector Machines

50

- ❑ **Implementation:** Using existing software packages like SVM^{light} and LIBSVM

- ❑ **Determining the kernel function :**
 - ❑ Linear kernel (no kernel): When $m \gg n$.
 - ❑ Gaussian, polynomial, string, and ...

- ❑ **Determining parameter values :** Grid search
 - ❑ Selecting a value for parameter C .
 - ❑ Selecting values for kernel function parameters (like σ).

Support Vector Machine, Neural Network, or Logistic Regression

51

❑ First Scenario : [$n \gg m$]

- ❑ Example: Spam detection (1000 training samples, 50,000 features).
- ❑ Logistic Regression or Linear Support Vector Machine

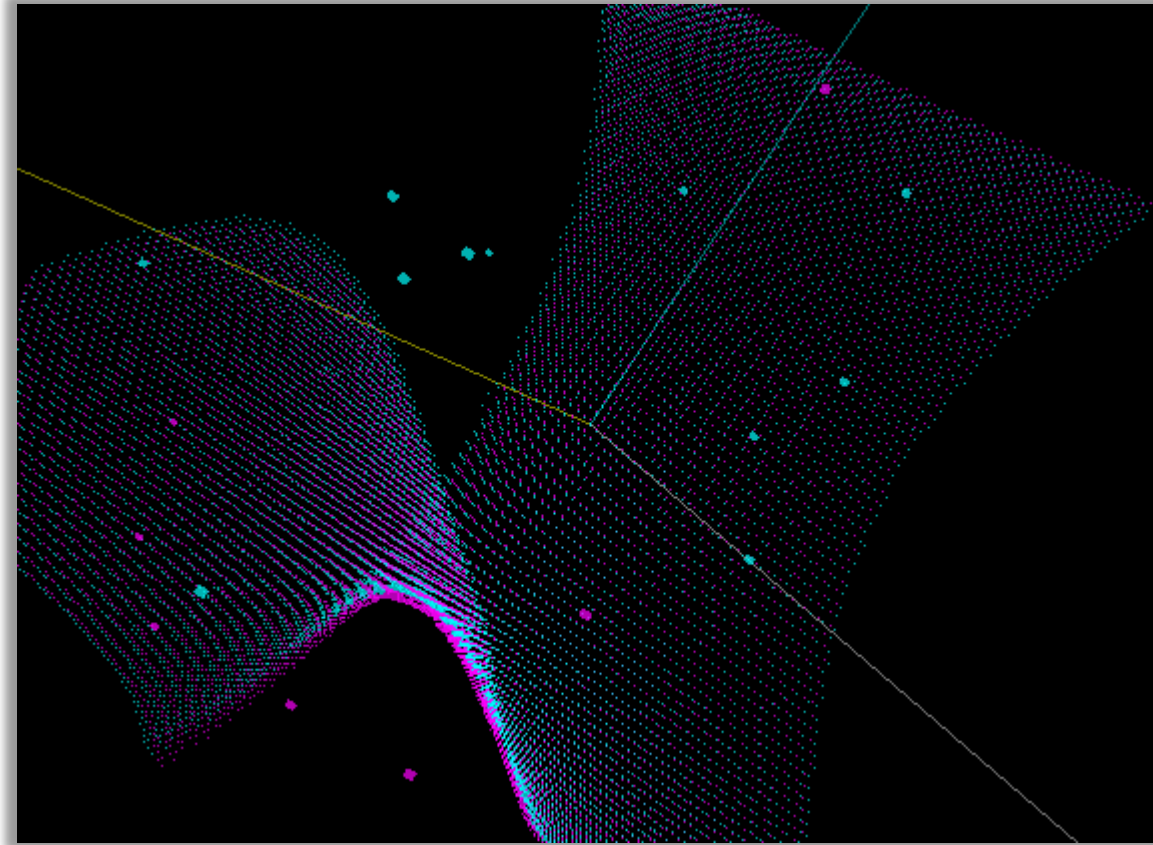
❑ Second Scenario : [Few features, a large number of training samples]

- ❑ Support Vector Machine with Gaussian Kernel.

Attention: Neural networks can be used in all the above scenarios, but they may require more training time.

Demo Execution

52



Appendix: More on Kernels

53

❑ **Question:** How do we know that using kernels helps us in data separation?

❑ In an n -dimensional space, any set of n independent vectors is linearly separable.

❑ If the matrix K is positive definite, then the data is linearly separable.

❑ **Theorem:** The matrix K is positive definite because $L^T L = K$.

❑ The i -th column in matrix L is equal to the vector $\phi(x^{(i)})$.

❑ **Proof:** Consider a non-zero vector v . In this case:

$$v^T K v = v^T L^T L v = (L v)^T (L v) = w^T w = \|w\|^2 \geq 0$$

And since both L and v are non-zero, the vector w is also non-zero. That is :

$$\|w\|^2 > 0 \Rightarrow v^T K v > 0 \Rightarrow K \text{ is positive definite}$$