

Introduction to Support Vector Machines

L FLEX, Machine Learning, Paul Sabatier University

Romain Thoreau
`romain.thoreau@cnes.fr`

February 2025

Binary Classification: Introductory Examples

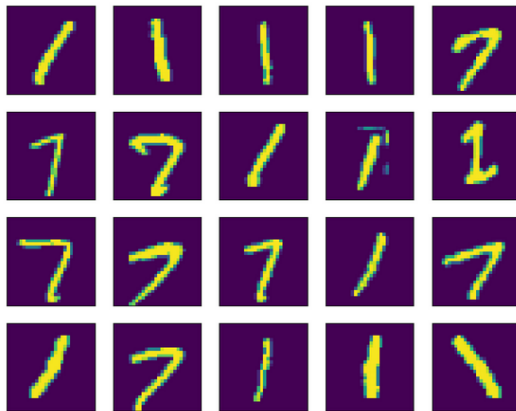


Figure: Images ($28\text{px} \times 28\text{px}$) of hand-written digits (1 & 7)

Binary Classification: Introductory Examples

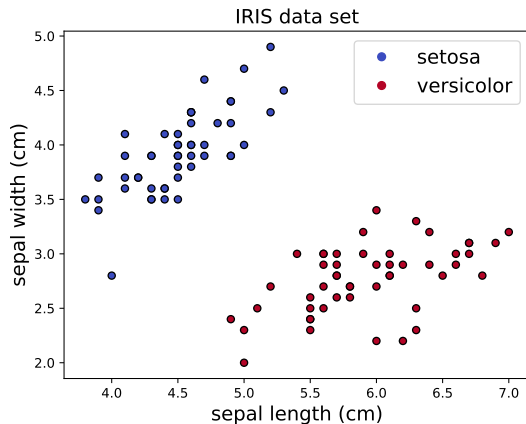
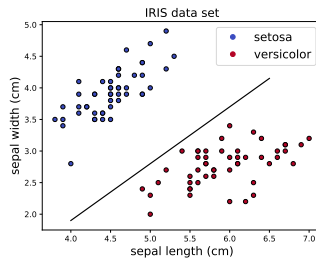
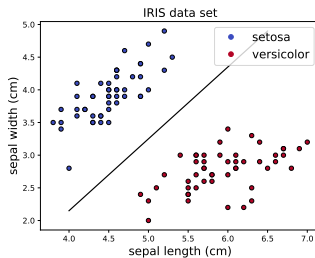
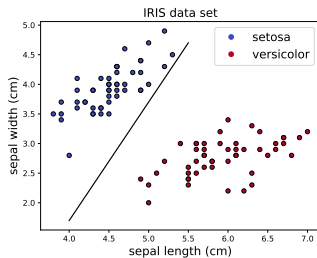


Figure: Samples of the IRIS data set [Fisher, 1936]

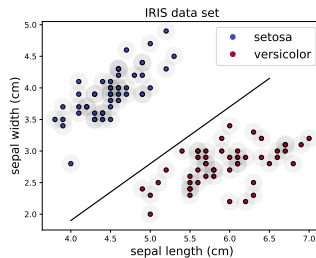
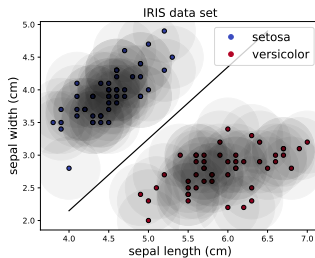
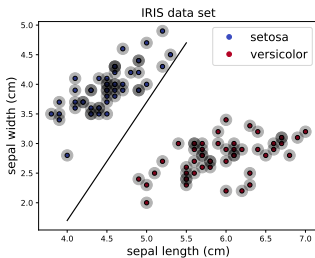
Linear Classification with a Hyperplane

- A priori, which hyperplane generalizes best?



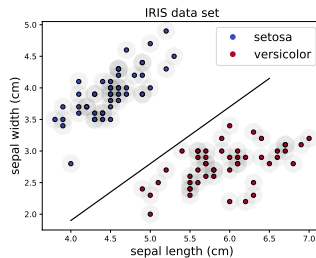
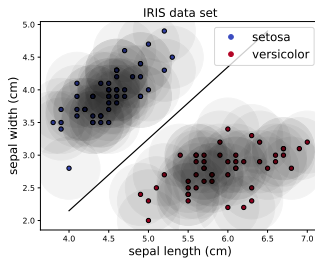
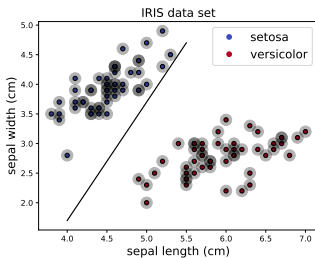
Linear Classification with a Hyperplane

- ▶ A priori, which hyperplane generalizes best?
- ▶ Which hyperplane is the most robust to gaussian noise in the data set?



Linear Classification with a Hyperplane

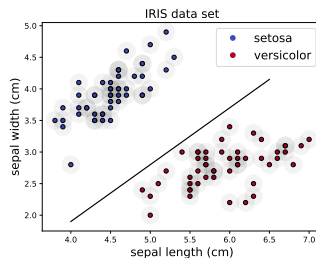
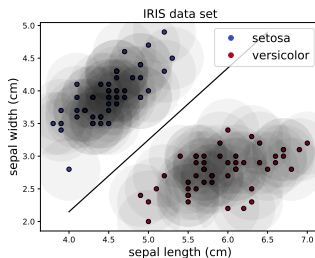
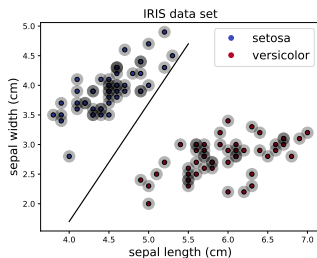
- ▶ A priori, which hyperplane generalizes best?
- ▶ Which hyperplane is the most robust to gaussian noise in the data set?



- ▶ An infinity of hyperplanes perfectly separates the data.

Linear Classification with a Hyperplane

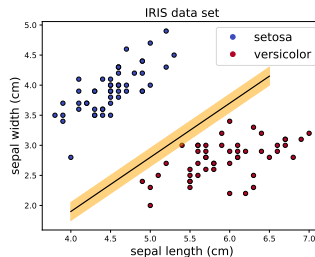
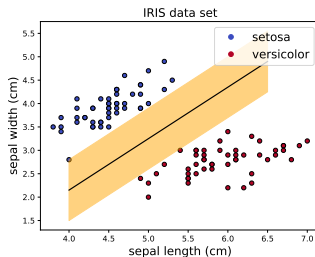
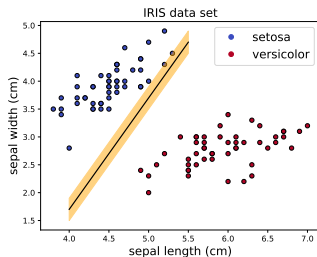
- ▶ A priori, which hyperplane generalizes best?
- ▶ Which hyperplane is the most robust to gaussian noise in the data set?



- ▶ An infinity of hyperplanes perfectly separates the data.
- ▶ We need **inductive biases**: assumptions about the data / model / training algorithm in order to select a way to generalize from the training data [Mitchell, 1980, Zhao et al., 2018].

Linear Classification with a Hyperplane

- ▶ A priori, which hyperplane generalizes best?
- ▶ Which hyperplane is the most robust to gaussian noise in the data set?



- ▶ An infinity of hyperplanes perfectly separates the data.
- ▶ We need **inductive biases**: assumptions about the data / model / training algorithm in order to select a way to generalize from the training data [Mitchell, 1980, Zhao et al., 2018].
- ▶ **Support Vector Machines**: select the hyperplane the maximizes the **margin**.

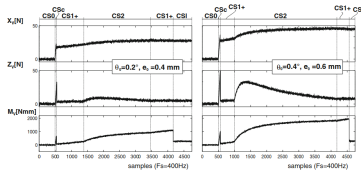
A few applications of SVM



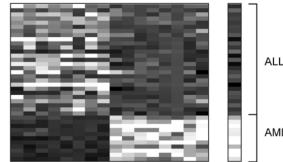
(a) Medical image classification,
e.g. [Camlica et al., 2015]



(b) Hyperspectral image segmentation,
e.g. [Mercier and Lennon, 2003]



(c) Contact states recognition in robotics,
e.g. [Jakovljevic et al., 2012]



(d) Gene selection, e.g. [Guyon et al., 2002]

Objectives of the lecture

- ▶ Develop a geometric intuition about classification problems,
- ▶ Understand the fundamentals of **supervised max-margin** classification models,
- ▶ Understand the need for model **regularization**,
- ▶ Understand the need for **inductive biases** in machine learning.

Bibliography

Those slides are inspired by numerous great resources including:

- ▶ Andrew Ng's lecture notes:
<https://see.stanford.edu/materials/aimlcs229/cs229-notes3.pdf>
- ▶ Pattern Recognition and Machine Learning, Christopher M. Bishop,
https://www.cs.uoi.gr/~arly/courses/ml/tmp/Bishop_book.pdf
- ▶ Kévin Bailly's lectures on SVM, <https://sites.google.com/view/bailly/>

Overview

- ▶ Linearly separable data
- ▶ Non-separable Data
- ▶ Non-Linear Classification

Support Vector Machines: Binary Linear Classification

We consider a labeled training data set $\mathcal{D}_{train} = \{(\mathbf{x}^{(i)}, y^{(i)}) | i \in \{1, \dots, N\}\}$ where:

$$\forall i \in \{1, \dots, N\}, \begin{cases} \mathbf{x}^{(i)} \in \mathcal{X} \subset \mathbb{R}^{D \times 1} \\ y^{(i)} \in \mathcal{Y} = \{-1, 1\} \end{cases}$$

Assuming that the data is linearly separable, how can we define a function $f_\theta : \mathcal{X} \longrightarrow \mathcal{Y}$, parameterized by θ , such that:

$$\forall i \in \{1, \dots, N\}, f_\theta(\mathbf{x}^{(i)}) = y^{(i)} \quad \text{for an optimal } \theta?$$

Support Vector Machines: Binary Linear Classification

We consider a labeled training data set $\mathcal{D}_{train} = \{(\mathbf{x}^{(i)}, y^{(i)}) | i \in \{1, \dots, N\}\}$ where:

$$\forall i \in \{1, \dots, N\}, \begin{cases} \mathbf{x}^{(i)} \in \mathcal{X} \subset \mathbb{R}^{D \times 1} \\ y^{(i)} \in \mathcal{Y} = \{-1, 1\} \end{cases}$$

Assuming that the data is linearly separable, how can we define a function $f_\theta : \mathcal{X} \longrightarrow \mathcal{Y}$, parameterized by θ , such that:

$$\forall i \in \{1, \dots, N\}, f_\theta(\mathbf{x}^{(i)}) = y^{(i)} \quad \text{for an optimal } \theta?$$

Linear Classifier

Let's denote $\theta = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$ where $\mathbf{w} \in \mathbb{R}^{D \times 1}$ and $b \in \mathbb{R}$. Besides, let's define the hyperplane $\mathcal{H} : \mathbf{w}^T \mathbf{x} + b = 0, \mathbf{x} \in \mathbb{R}^{D \times 1}$. We define a linear classifier f_θ as follows:

$$f_\theta : \mathcal{X} \longrightarrow \mathcal{Y}$$

$$f_\theta(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) \quad \text{where } \forall v \in \mathbb{R}, \begin{cases} \sigma(v) = 1 \text{ if } v \geq 0 \\ \sigma(v) = -1 \text{ otherwise} \end{cases}$$

Support Vector Machines: Margin Definition

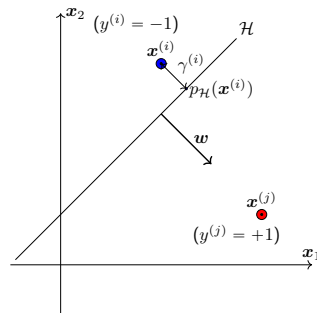
We aim to find the model parameters that maximize the margin between the hyperplane \mathcal{H} and its closest training data points.

- First, let us show that the distance $\gamma^{(i)} := d(\mathbf{x}^{(i)}, \mathcal{H}) = d(\mathbf{x}^{(i)}, p_{\mathcal{H}}(\mathbf{x}^{(i)}))$ can be expressed as follows:

$$\gamma^{(i)} = \frac{1}{\|\mathbf{w}\|} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

Support Vector Machines: Margin Definition

We aim to find the model parameters that maximize the margin between the hyperplane \mathcal{H} and its closest training data points.



- First, let us show that the distance $\gamma^{(i)} := d(\mathbf{x}^{(i)}, \mathcal{H}) = d(\mathbf{x}^{(i)}, p_{\mathcal{H}}(\mathbf{x}^{(i)}))$ can be expressed as follows:

$$\gamma^{(i)} = \frac{1}{\|\mathbf{w}\|} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

Support Vector Machines: Functional and Geometric Margins

- Let's define the **geometric margin** γ :

$$\gamma = \min_{i \in \{1, \dots, N\}} \gamma^{(i)} = \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)} = \frac{1}{\|\mathbf{w}\|} \hat{\gamma}$$

Support Vector Machines: Functional and Geometric Margins

- Let's define the **geometric margin** γ :

$$\gamma = \min_{i \in \{1, \dots, N\}} \gamma^{(i)} = \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)} = \frac{1}{\|\mathbf{w}\|} \hat{\gamma}$$

where $\hat{\gamma}$, called the **functional margin**, is defined as follows:

$$\hat{\gamma} = \min_{i \in \{1, \dots, N\}} \hat{\gamma}^{(i)} = \min_{i \in \{1, \dots, N\}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

Support Vector Machines: Functional and Geometric Margins

- Let's define the **geometric margin** γ :

$$\gamma = \min_{i \in \{1, \dots, N\}} \gamma^{(i)} = \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)} = \frac{1}{\|\mathbf{w}\|} \hat{\gamma}$$

where $\hat{\gamma}$, called the **functional margin**, is defined as follows:

$$\hat{\gamma} = \min_{i \in \{1, \dots, N\}} \hat{\gamma}^{(i)} = \min_{i \in \{1, \dots, N\}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

What does the functional margin represent? Is it a good metric for model uncertainty? Should we maximize the functional margin or the geometric margin?

Support Vector Machines: Functional and Geometric Margins

- Let's define the **geometric margin** γ :

$$\gamma = \min_{i \in \{1, \dots, N\}} \gamma^{(i)} = \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)} = \frac{1}{\|\mathbf{w}\|} \hat{\gamma}$$

where $\hat{\gamma}$, called the **functional margin**, is defined as follows:

$$\hat{\gamma} = \min_{i \in \{1, \dots, N\}} \hat{\gamma}^{(i)} = \min_{i \in \{1, \dots, N\}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

What does the functional margin represent? Is it a good metric for model uncertainty? Should we maximize the functional margin or the geometric margin?

Notice that the functional margin can be arbitrarily increased by a factor $\alpha > 0$ by substituting \mathbf{w} and b with $\alpha\mathbf{w}$ and αb while f_θ is kept unchanged:

$$\begin{cases} f_{\theta'}(\mathbf{x}) &= \sigma(\mathbf{w}'^T \mathbf{x} + b') = \sigma(\alpha(\mathbf{w}^T \mathbf{x} + b)) = \sigma(\mathbf{w}^T \mathbf{x} + b) = f_\theta(\mathbf{x}) \\ \hat{\gamma}'^{(i)} &= \alpha \hat{\gamma}^{(i)} \end{cases}$$

Support Vector Machines: Objective Function

We aim to maximize the geometric margin: $\max_{\mathbf{w}, b} \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)}$.

► *Can we turn this optimization problem into a quadratic optimization problem with linear constraints, such that standard solvers could solve it?*

SVM as a Quadratic Optimization Problem

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \quad & \frac{1}{2} \boldsymbol{\theta}^T Q \boldsymbol{\theta} + \mathbf{p}^T \boldsymbol{\theta} \\ \text{s.t.} \quad & A \boldsymbol{\theta} \geq \mathbf{c} \end{aligned}$$

Support Vector Machines: Objective Function

We aim to maximize the geometric margin: $\max_{\mathbf{w}, b} \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)}$.

► *Can we turn this optimization problem into a quadratic optimization problem with linear constraints, such that standard solvers could solve it?*

SVM as a Quadratic Optimization Problem

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \quad & \frac{1}{2} \boldsymbol{\theta}^T Q \boldsymbol{\theta} + \mathbf{p}^T \boldsymbol{\theta} \\ \text{s.t.} \quad & A \boldsymbol{\theta} \geq \mathbf{c} \end{aligned}$$

SVM Formulation

Since f_{θ} is invariant to the scale of the functional margin, we can arbitrarily set $\hat{\gamma} = 1$:

$$\begin{aligned} \max_{\mathbf{w}, b} \quad & \frac{1}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & \min_{i \in \{1, \dots, N\}} \underbrace{y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)}_{\hat{\gamma}^{(i)}} = 1 \end{aligned}$$

Support Vector Machines: Objective Function

We aim to maximize the geometric margin: $\max_{\mathbf{w}, b} \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)}$.

► *Can we turn this optimization problem into a quadratic optimization problem with linear constraints, such that standard solvers could solve it?*

SVM as a Quadratic Optimization Problem

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \quad & \frac{1}{2} \boldsymbol{\theta}^T Q \boldsymbol{\theta} + \mathbf{p}^T \boldsymbol{\theta} \\ \text{s.t.} \quad & A \boldsymbol{\theta} \geqslant \mathbf{c} \end{aligned}$$

SVM Formulation

Since f_{θ} is invariant to the scale of the functional margin, we can arbitrarily set $\hat{\gamma} = 1$:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & \min_{i \in \{1, \dots, N\}} \underbrace{y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b)}_{\hat{\gamma}^{(i)}} = 1 \end{aligned}$$

Support Vector Machines: Objective Function

We aim to maximize the geometric margin: $\max_{\mathbf{w}, b} \min_{i \in \{1, \dots, N\}} \frac{1}{\|\mathbf{w}\|} \hat{\gamma}^{(i)}$.

► *Can we turn this optimization problem into a quadratic optimization problem with linear constraints, such that standard solvers could solve it?*

SVM as a Quadratic Optimization Problem

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \quad & \frac{1}{2} \boldsymbol{\theta}^T Q \boldsymbol{\theta} + \mathbf{p}^T \boldsymbol{\theta} \\ \text{s.t.} \quad & A \boldsymbol{\theta} \geq \mathbf{c} \end{aligned}$$

SVM Formulation

Since f_{θ} is invariant to the scale of the functional margin, we can arbitrarily set $\hat{\gamma} = 1$:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & \forall i \in \{1, \dots, N\}, y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \end{aligned}$$

SVM as a Quadratic Optimization Problem

- Write the corresponding values of θ , Q , p , A and c .

SVM as a Quadratic Optimization Problem

- Write the corresponding values of θ , Q , p , A and c .

$$\theta = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$Q = \begin{bmatrix} 0 & \mathbf{0}_{1 \times D} \\ \mathbf{0}_{D \times 1} & \mathbf{I}_D \end{bmatrix}$$

$$p = \mathbf{0}_{(1+D) \times 1}$$

$$A = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_N^T \end{bmatrix} \text{ with:}$$

$$\mathbf{a}_i^T = [y^{(i)} \ y^{(i)} \mathbf{x}^{(i)T}]$$

$$\mathbf{c} = \mathbf{1}_{N \times 1}$$

Support Vector Machines: Example

- Let's consider the following data set:

$$\mathcal{D}_{train} = \{([0 \ 0]^T, -1), ([2 \ 2]^T, -1), ([2 \ 0]^T, 1), ([3 \ 0]^T, 1)\}$$

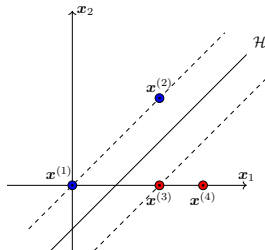
1. Write the linear constraints for this data set,
2. Compute the optimal hyperplane parameters,
3. Compute the optimal margin,
4. Draw the hyperplane, margin, samples.
Which samples are on the margin?

Support Vector Machines: Example

- Let's consider the following data set:

$$\mathcal{D}_{train} = \{([0 \ 0]^T, -1), ([2 \ 2]^T, -1), ([2 \ 0]^T, 1), ([3 \ 0]^T, 1)\}$$

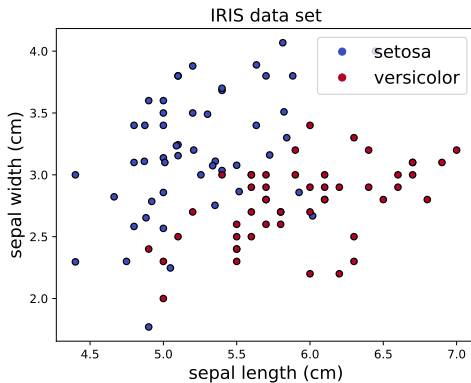
- Write the linear constraints for this data set,
- Compute the optimal hyperplane parameters,
- Compute the optimal margin,
- Draw the hyperplane, margin, samples.
Which samples are on the margin?



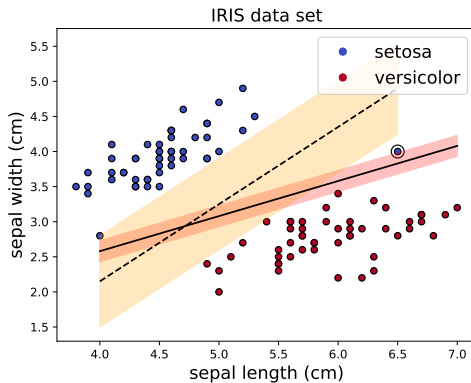
What happens if we remove $x^{(4)}$ from the data set? $x^{(3)}$?

- ▶ Linearly separable data
- ▶ **Non-separable Data**
- ▶ Non-Linear Classification

Non-separable data: examples of data with noise and outliers



(a) Samples of a noisy version of the IRIS data set



(b) Samples of the IRIS data set with an outlier

Regularization of the SVM problem

- ▶ With the current formulation of the SVM problem, the optimization algorithm will not converge.

SVM Formulation

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & \forall i \in \{1, \dots, N\}, y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \end{aligned}$$

- ▶ *Can we modify the optimization problem to handle non-separable data?*

Regularization of the SVM problem

- ▶ With the current formulation of the SVM problem, the optimization algorithm will not converge.

SVM Formulation

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & \forall i \in \{1, \dots, N\}, y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \end{aligned}$$

- ▶ *Can we modify the optimization problem to handle non-separable data?*

Relaxed SVM Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \forall i \in \{1, \dots, N\}, y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i \\ & \forall i \in \{1, \dots, N\}, \xi_i \geq 0 \end{aligned}$$

Regularization of the SVM problem

Relaxed SVM Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \forall i \in \{1, \dots, N\}, y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i \\ & \forall i \in \{1, \dots, N\}, \xi_i \geq 0 \end{aligned}$$

Regularization of the SVM problem

Relaxed SVM Formulation

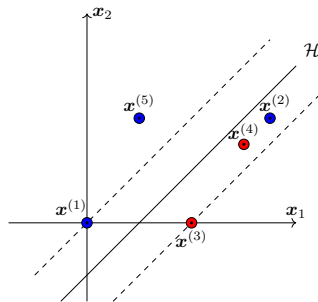
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \forall i \in \{1, \dots, N\}, y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i \\ & \forall i \in \{1, \dots, N\}, \xi_i \geq 0 \end{aligned}$$

- ▶ What happens if $C = 0$? $C = \infty$?
- ▶ Give the formulation of the relaxed problem to match the quadratic formulation:

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \quad & \frac{1}{2} \boldsymbol{\theta}^T Q \boldsymbol{\theta} + \mathbf{p}^T \boldsymbol{\theta} \\ \text{s.t.} \quad & A \boldsymbol{\theta} \geq \mathbf{c} \end{aligned}$$

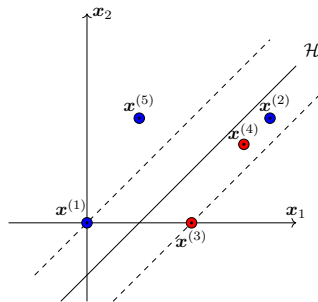
Regularization of the SVM problem

- What are the values of $\xi_i, i \in \{1, \dots, 5\}$?



Regularization of the SVM problem

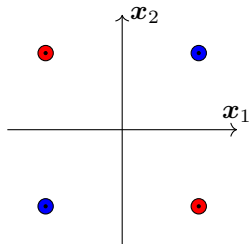
► What are the values of $\xi_i, i \in \{1, \dots, 5\}$?



$$\xi_1 = 0, \quad \xi_2 = 0, \quad \xi_3 = 0, \quad 0 < \xi_4 < 1, \quad \xi_5 = 0$$

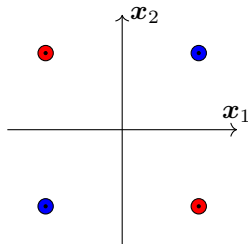
- ▶ Linearly separable data
- ▶ Non-separable Data
- ▶ Non-Linear Classification

Non-Linear Classification: the XOR example



Can you find a separating hyperplane?

Non-Linear Classification: the XOR example



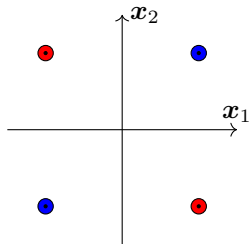
$$\phi : \mathbb{R}^2 \longrightarrow ?$$

$$\phi(\mathbf{x}) = ?$$

Can you find a separating hyperplane?

- *Could we map the data in a higher dimensional space where it would be linearly separable?*

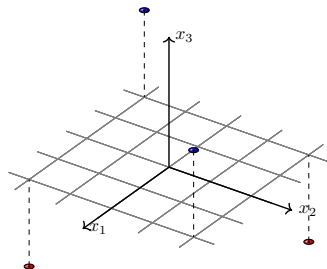
Non-Linear Classification: the XOR example



Can you find a separating hyperplane?

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$\phi(\mathbf{x}) = [x_1 \ x_2 \ x_1 x_2]^T$$



- *Could we map the data in a higher dimensional space where it would be linearly separable?*

Non-Linear Mapping before SVM

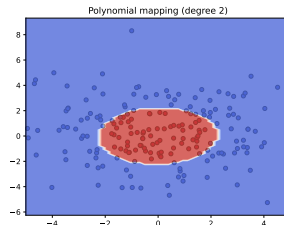
- ▶ So, when data is not linearly separable, we can find a mapping $\phi : \mathcal{X} \longrightarrow \mathcal{Z}$ from the input space \mathcal{X} to a feature space \mathcal{Z} of higher dimension, where features are linearly separable.
- ▶ Then, we apply the SVM on $z = \phi(x)$ instead of x .

Another example is a polynomial mapping:

$$\phi : \mathbb{R}^2 \longrightarrow \mathbb{R}^5$$

$$\phi(x) = [x_1 \ x_2 \ x_1 x_2 \ x_1^2 \ x_2^2]^T$$

- ▶ However, computing $\phi(x)$ can be computationally intensive.
- ▶ In order to circumvent this issue, we need to consider another formulation of the SVM optimization problem and to digress for a moment on Lagrange duality!



Lagrangian Duality

Primal Optimization Problem

The SVM optimization problem can be put in the following form, called the primal problem:

$$\begin{aligned} \min_{\boldsymbol{\theta}} \quad & f(\boldsymbol{\theta}) \\ \text{s.t.} \quad & g_i(\boldsymbol{\theta}) \leq 0, i \in \{1, \dots, N\} \end{aligned}$$

Lagrangian of the Primal Optimization Problem

$$\mathcal{L}(\boldsymbol{\theta}, \alpha) = f(\boldsymbol{\theta}) + \sum_{i=1}^N \alpha_i g_i(\boldsymbol{\theta})$$

► *Prove that the following optimization problem has the same solution than the primal problem:*

$$\min_{\boldsymbol{\theta}} \max_{\alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \alpha)$$

Lagrangian Duality

Primal Optimization problem

$$\min_{\boldsymbol{\theta}} \max_{\alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \alpha)$$

Dual Optimization problem

$$\max_{\alpha \geq 0} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \alpha)$$

Strong Duality

If the optimization problem is convex, then the solutions of the **primal** and **dual** problems are equal:

$$\min_{\boldsymbol{\theta}} \max_{\alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \alpha) = \max_{\alpha \geq 0} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \alpha)$$

Karush-Kuhn-Tucker (KKT) Conditions

$(\boldsymbol{\theta}^*, \alpha^*)$ are the optimal solutions of the primal / dual problem if and only if they satisfy the **Karush-Kuhn-Tucker conditions**:

Stationarity	$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \alpha) _{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \alpha=\alpha^*} = 0$	(1)
Primal feasibility	$\forall i \in \{1, \dots, N\}, g_i(\boldsymbol{\theta}^*) \leq 0$	(2)
Dual feasibility	$\forall i \in \{1, \dots, N\}, \alpha_i^* \geq 0$	(3)
Complementary slackness	$\forall i \in \{1, \dots, N\}, \alpha_i^* g_i(\boldsymbol{\theta}^*) = 0$	(4)

Condition (4) states that if $\alpha_i^* > 0$, then $g_i(\boldsymbol{\theta}^*) = 0$, meaning that the constraint $g_i(\boldsymbol{\theta}^*) \leq 0$ is active.

Why do we care about the Dual Problem?

- *Write the Lagrangian of the SVM optimization problem:*

Why do we care about the Dual Problem?

- *Write the Lagrangian of the SVM optimization problem:*

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i (y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1)$$

Why do we care about the Dual Problem?

- *Write the Lagrangian of the SVM optimization problem:*

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i (y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1)$$

- *Find optimal conditions for \mathbf{w}^* and b^* (note that the Lagrangian is convex):*

Why do we care about the Dual Problem?

- *Write the Lagrangian of the SVM optimization problem:*

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i (y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1)$$

- *Find optimal conditions for \mathbf{w}^* and b^* (note that the Lagrangian is convex):*

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) = 0 \iff \mathbf{w} - \sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)} = 0 \iff \mathbf{w} = \sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

Why do we care about the Dual Problem?

- Write the Lagrangian of the SVM optimization problem:

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i (y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1)$$

- Find optimal conditions for \mathbf{w}^* and b^* (note that the Lagrangian is convex):

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \alpha) = 0 \iff \mathbf{w} - \sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)} = 0 \iff \mathbf{w} = \sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

$$\nabla_b \mathcal{L}(\mathbf{w}, b, \alpha) = 0 \iff \sum_{i=1}^N \alpha_i y^{(i)} = 0$$

Why do we care about the Dual Problem?

Hence, $\mathcal{L}(\mathbf{w}^*, b^*, \alpha) =$

Why do we care about the Dual Problem?

$$\begin{aligned} \text{Hence, } \mathcal{L}(\mathbf{w}^*, b^*, \alpha) &= \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \left(\sum_{j=1}^N \alpha_j y^{(j)} \mathbf{x}^{(j)} \right) - \sum_{i=1}^N \alpha_i y^{(i)} \left(\sum_{j=1}^N \alpha_j y^{(j)} \mathbf{x}^{(j)} \right)^T \mathbf{x}^{(i)} \\ &\quad - b \sum_{i=1}^N \alpha_i y^{(i)} + \sum_{i=1}^N \alpha_i \end{aligned}$$

Why do we care about the Dual Problem?

$$\begin{aligned}
 \text{Hence, } \mathcal{L}(\mathbf{w}^*, b^*, \alpha) &= \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \left(\sum_{j=1}^N \alpha_j y^{(j)} \mathbf{x}^{(j)} \right) - \sum_{i=1}^N \alpha_i y^{(i)} \left(\sum_{j=1}^N \alpha_j y^{(j)} \mathbf{x}^{(j)} \right)^T \mathbf{x}^{(i)} \\
 &\quad - b \sum_{i=1}^N \alpha_i y^{(i)} + \sum_{i=1}^N \alpha_i \\
 &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)}
 \end{aligned}$$

Why do we care about the Dual Problem?

$$\begin{aligned}
 \text{Hence, } \mathcal{L}(\mathbf{w}^*, b^*, \alpha) &= \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \left(\sum_{j=1}^N \alpha_j y^{(j)} \mathbf{x}^{(j)} \right) - \sum_{i=1}^N \alpha_i y^{(i)} \left(\sum_{j=1}^N \alpha_j y^{(j)} \mathbf{x}^{(j)} \right)^T \mathbf{x}^{(i)} \\
 &\quad - b \sum_{i=1}^N \alpha_i y^{(i)} + \sum_{i=1}^N \alpha_i \\
 &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)}
 \end{aligned}$$

SVM Dual Problem Formulation

$$\begin{aligned}
 \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} \\
 \text{s.t. } \quad & \forall i \in \{1, \dots, N\}, \alpha_i \geq 0, \\
 & \sum_{i=1}^N \alpha_i y^{(i)} = 0
 \end{aligned}$$

Why do we care about the Dual Problem?

- ▶ We can solve the SVM dual problem with a standard solver for quadratic optimization with linear constraints in order to find α^* .
- ▶ It follows that the optimal classifier is defined by:

$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y^{(i)} \mathbf{x}^{(i)} \quad b^* = ?$$

Why do we care about the Dual Problem?

- ▶ We can solve the SVM dual problem with a standard solver for quadratic optimization with linear constraints in order to find α^* .
- ▶ It follows that the optimal classifier is defined by:

$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y^{(i)} \mathbf{x}^{(i)} \quad b^* = ?$$

- ▶ Recall that, from the KKT conditions, we have that if $\alpha_i^* > 0$, then $y^{(i)}(\mathbf{w}^{*T} \mathbf{x}^{(i)} + b^*) = 1$. In other words, the only samples $\mathbf{x}^{(i)}$ that contribute to \mathbf{w}^* are the support vectors lying on the margin.

Why do we care about the Dual Problem?

- ▶ We can solve the SVM dual problem with a standard solver for quadratic optimization with linear constraints in order to find α^* .
- ▶ It follows that the optimal classifier is defined by:

$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y^{(i)} \mathbf{x}^{(i)} \quad b^* = ?$$

- ▶ Recall that, from the KKT conditions, we have that if $\alpha_i^* > 0$, then $y^{(i)}(\mathbf{w}^{*T} \mathbf{x}^{(i)} + b^*) = 1$. In other words, the only samples $\mathbf{x}^{(i)}$ that contribute to \mathbf{w}^* are the support vectors lying on the margin.
- ▶ For a test sample \mathbf{x}^* , the prediction is given by:

$$\hat{y} = \sigma \left(\sum_{\alpha_i^* > 0} \alpha_i^* y^{(i)} \mathbf{x}^{(i)T} \mathbf{x}^* + b^* \right)$$

The Kernel Trick

- For non-linear classification problem, the prediction becomes:

$$\hat{y} = \sigma \left(\sum_{\alpha_i^* > 0} \alpha_i^* y^{(i)} \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^*) + b^* \right)$$

The Kernel Trick

- For non-linear classification problem, the prediction becomes:

$$\hat{y} = \sigma \left(\sum_{\alpha_i^* > 0} \alpha_i^* y^{(i)} K_{\phi}(\mathbf{x}^{(i)}, \mathbf{x}^*) + b^* \right)$$

Kernel Trick

We do not need to explicitly compute the mapping $\phi(\mathbf{x})$, that can be costly, if we know a formulation of kernel K_{ϕ} that does not depend on ϕ :

$$K_{\phi}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

The Kernel Trick

- For non-linear classification problem, the prediction becomes:

$$\hat{y} = \sigma \left(\sum_{\alpha_i^* > 0} \alpha_i^* y^{(i)} K_{\phi}(\mathbf{x}^{(i)}, \mathbf{x}^*) + b^* \right)$$

Kernel Trick

We do not need to explicitly compute the mapping $\phi(\mathbf{x})$, that can be costly, if we know a formulation of kernel K_{ϕ} that does not depend on ϕ :

$$K_{\phi}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

- *Find the kernel corresponding to the following polynomial mapping:*

$$\phi : \mathbb{R}^2 \longrightarrow \mathbb{R}^5$$

$$\phi(\mathbf{x}) = [\mathbf{x}_1 \ \mathbf{x}_2 \ \sqrt{2}\mathbf{x}_1\mathbf{x}_2 \ \mathbf{x}_1^2 \ \mathbf{x}_2^2]^T$$

The Kernel Trick

- ▶ The Kernel function can be thought as a similarity measure: the more similar \mathbf{x} and \mathbf{x}' , the larger $K_\phi(\mathbf{x}, \mathbf{x}')$.
- ▶ Radial Basis Function (RBF) Kernel:

$$K_\phi(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \exp\left(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2\right), \quad \gamma > 0$$

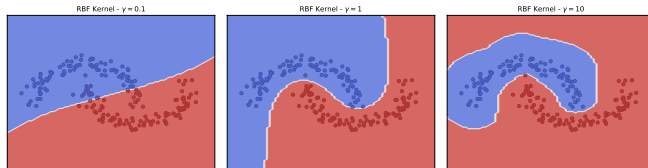









Figure: SVM decision boundaries on the Two Moons data set with various RBF kernels

- ▶ We don't need to explicitly know ϕ !

-  Camlica, Z., Tizhoosh, H. R., and Khalvati, F. (2015).
Medical image classification via svm using lbp features from saliency-based folded data.
In *2015 IEEE 14th international conference on machine learning and applications (ICMLA)*, pages 128–132. IEEE.
-  Fisher, R. A. (1936).
The use of multiple measurements in taxonomic problems.
Annals of eugenics, 7(2):179–188.
-  Guyon, I., Weston, J., Barnhill, S., and Vapnik, V. (2002).
Gene selection for cancer classification using support vector machines.
Machine learning, 46:389–422.
-  Jakovljevic, Z., Petrovic, P. B., and Hodolic, J. (2012).
Contact states recognition in robotic part mating based on support vector machines.
The International Journal of Advanced Manufacturing Technology, 59:377–395.
-  Mercier, G. and Lennon, M. (2003).
Support vector machines for hyperspectral image classification with spectral-based kernels.
In *IGARSS 2003. 2003 IEEE International Geoscience and Remote Sensing Symposium. Proceedings (IEEE Cat. No. 03CH37477)*, volume 1, pages 288–290. IEEE.
-  Mitchell, T. M. (1980).
The need for biases in learning generalizations (rutgers computer science tech. rept. cbm-tr-117).
Rutgers University.
-  Zhao, S., Ren, H., Yuan, A., Song, J., Goodman, N., and Ermon, S. (2018).
Bias and generalization in deep generative models: An empirical study.
Advances in Neural Information Processing Systems, 31.