

# Pattern Recognition - Summary

Friedrich-Alexander-Universität Erlangen-Nürnberg

Winter Term 2025/26

November 26, 2025

## Contents

<b>1</b>	<b>Optimization</b>	<b>2</b>
1.1	Convexity . . . . .	2
1.2	Unconstrained Optimization . . . . .	2
1.2.1	Finding a Suitable Step Size (Line Search) . . . . .	2
1.2.2	Gradient Descent . . . . .	3
1.2.3	Normalized Steepest Descent (General Norms) . . . . .	3
1.2.4	Newton's Method . . . . .	4
1.3	Constrained Optimization . . . . .	5
1.3.1	The Lagrangian & Dual Function . . . . .	5
1.3.2	Strong Duality & Slater's Condition . . . . .	6
1.3.3	Karush-Kuhn-Tucker (KKT) Conditions . . . . .	6
<b>2</b>	<b>Support Vector Machines (SVM)</b>	<b>7</b>
2.1	Hard Margin Problems . . . . .	7
2.2	Soft Margin Problems . . . . .	8

# 1 Optimization

Optimization is crucial for many solutions in pattern recognition (e.g., training classifiers). We distinguish between unconstrained and constrained optimization problems.

## 1.1 Convexity

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is called **convex** if for all  $\vec{x}, \vec{y} \in \mathbb{R}^n$  and for all  $\theta \in [0, 1]$ , the following condition holds:

### Definition: Convexity & Concavity

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is **convex** if the domain  $\text{dom}(f)$  is a convex set and if  $\forall \vec{x}, \vec{y} \in \text{dom}(f)$ , and  $\theta$  with  $0 \leq \theta \leq 1$ , we have:

$$f(\theta \vec{x} + (1 - \theta) \vec{y}) \leq \theta f(\vec{x}) + (1 - \theta) f(\vec{y}) \quad (1.1)$$

A function is **concave** if  $-f$  is convex.

**Implication:** For convex functions, any local minimum is also a global minimum. This property is particularly useful because gradient-based methods won't get stuck in suboptimal local minima.

## 1.2 Unconstrained Optimization

Here, we aim to find the minimum of a function  $f(\vec{x})$  without any restrictions on  $\vec{x}$ . Typically, we assume  $f$  is twice differentiable and convex.

$$\vec{x}^* = \text{argmin}_{\vec{x}} f(\vec{x}) \quad (1.2)$$

A necessary and sufficient condition for the minimum is the zero-crossing of the gradient:

$$\nabla f(\vec{x}^*) = 0 \quad (1.3)$$

Since a closed-form solution is often impossible, we use iterative approaches:

initialization:  $\vec{x}^{(0)}$

iteration step:  $\vec{x}^{(k+1)} = \vec{x}^{(k)} + t^{(k)} \Delta \vec{x}^{(k)}$

where  $\Delta \vec{x}^{(k)}$  is the **search direction** and  $t^{(k)}$  is the **step size**.

### 1.2.1 Finding a Suitable Step Size (Line Search)

Choosing the correct step size  $t^{(k)}$  is crucial:

- **Too small:** Convergence is extremely slow.
- **Too large:** The algorithm might overshoot the minimum or diverge.

Instead of finding the exact optimal  $t$  (which is computationally expensive), we use **inexact line search** methods like **Backtracking Line Search** (Armijo-Goldstein).

The goal is to find a step size  $t$  that, that puts us below the red line, ensuring that the function value decreases sufficiently (not just barely) relative to the step size.

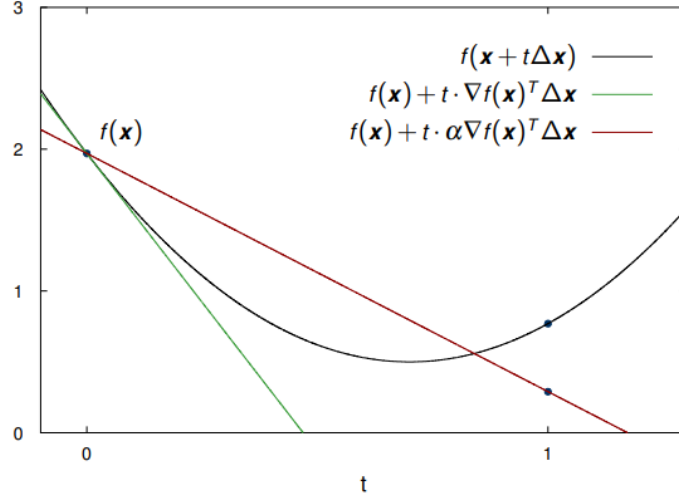


Figure 1: Backtracking line search

### Armijo-Goldstein Condition (Backtracking)

We start with a large step size ( $t = 1$ ) and iteratively reduce it ( $t := \beta t$ ) until the function value decreases sufficiently. The condition is:

$$f(\vec{x} + t\Delta\vec{x}) \leq f(\vec{x}) + \alpha t \nabla f(\vec{x})^T \Delta\vec{x} \quad (1.4)$$

Where  $\alpha \in (0, 0.5)$  defines the required "steepness" of the descent.

### 1.2.2 Gradient Descent

The most natural choice for the search direction is the direction of steepest descent, which is the negative gradient.

$$\Delta\vec{x}^{(k)} = -\nabla f(\vec{x}^{(k)}) \quad (1.5)$$

#### Algorithm: Gradient Descent

1. Set direction:  $\Delta\vec{x}^{(k)} = -\nabla f(\vec{x}^{(k)})$
2. Line search (find optimal step size  $t$ ):

$$t^{(k)} = \operatorname{argmin}_{t \geq 0} f(\vec{x}^{(k)} + t\Delta\vec{x}^{(k)})$$

(Usually approximated via **Backtracking Line Search** / **Armijo-Goldstein**).

3. Update:  $\vec{x}^{(k+1)} = \vec{x}^{(k)} + t^{(k)}\Delta\vec{x}^{(k)}$
4. Repeat until convergence ( $\|\vec{x}^{(k)} - \vec{x}^{(k-1)}\| < \epsilon$ ).

### 1.2.3 Normalized Steepest Descent (General Norms)

Ideally, we want the direction that gives the largest decrease in the linear approximation of  $f$ . This depends on the chosen norm  $\|\cdot\|$ .

$$\Delta\vec{x} = \operatorname{argmin}_{\vec{u}} \{\nabla f(\vec{x})^T \vec{u} \mid \|\vec{u}\| = 1\} \quad (1.6)$$

- **$L_2$ -Norm:** Direction is  $-\nabla f(\vec{x})$  (Standard Gradient Descent). Unit ball is a sphere.

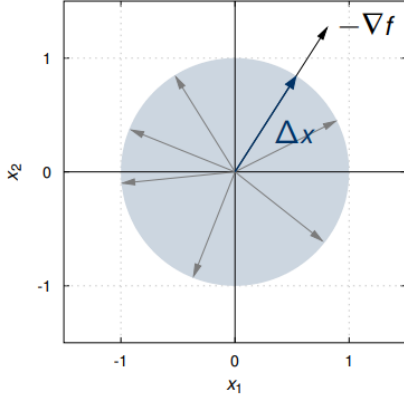


Figure 2: Unit ball in  $L_2$  norm

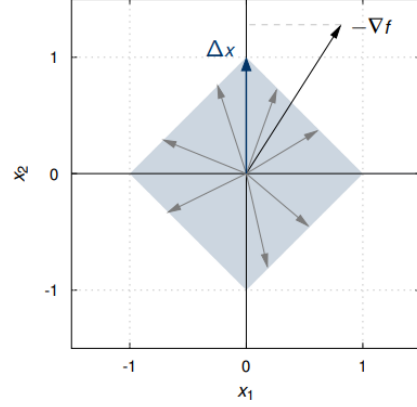


Figure 3: Unit ball in  $L_1$  norm

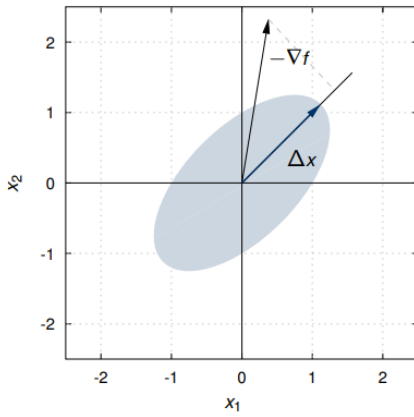


Figure 4: Unit ball in  $L_P$  norm

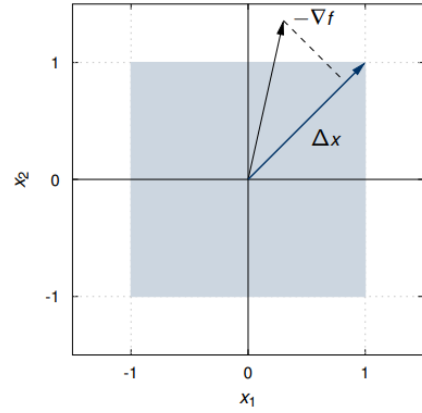


Figure 5: Unit ball in  $L_\infty$  norm

- **$L_1$ -Norm:** Direction is along the coordinate axis with the largest absolute gradient value (Coordinate Descent). Unit ball is a diamond.
- **$L_\infty$ -Norm:** Direction points towards the corners of the unit hypercube (e.g., vector with entries  $\pm 1$ ). Unit ball is a square/cube.
- **$L_P$ -Norm (Quadratic Norm):** Defined by a positive definite matrix  $\mathbf{P}$  as  $\|\vec{u}\|_{\mathbf{P}} = \sqrt{\vec{u}^T \mathbf{P} \vec{u}}$ .

$$\Delta \vec{x} = -\mathbf{P}^{-1} \nabla f(\vec{x}) \quad (1.7)$$

This transforms the space to make the contours spherical before taking the gradient step (similar to whitening in LDA). The update direction aligns with the largest principal component in the transformed space.

**Crucial Connection:** If we set  $\mathbf{P} = \nabla^2 f(\vec{x})$  (the Hessian), we effectively perform steepest descent in the local curvature norm, which yields **Newton's Method**.

#### 1.2.4 Newton's Method

Newton's method finds the search direction by approximating the function  $f(\vec{x})$  locally using a **second-order Taylor polynomial**:

$$f(\vec{x} + \Delta \vec{x}) \approx f(\vec{x}) + \nabla f(\vec{x})^T \Delta \vec{x} + \frac{1}{2} \Delta \vec{x}^T \nabla^2 f(\vec{x}) \Delta \vec{x} \quad (1.8)$$

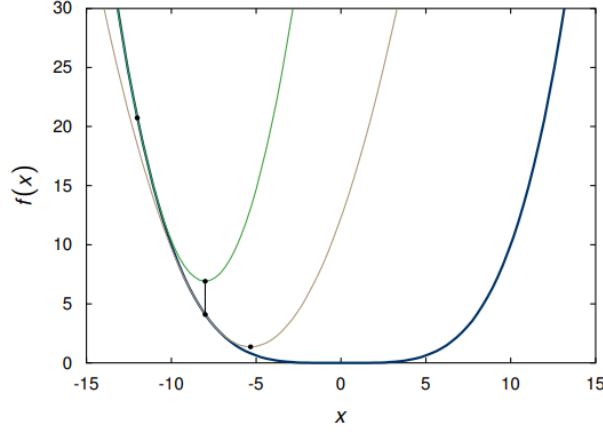


Figure 6: Newtons Method using 2nd order Taylor approximation

We find the optimal step  $\Delta\vec{x}$  by setting the derivative of this approximation to zero:

$$\begin{aligned}\nabla_{\Delta\vec{x}}(\dots) &= \nabla f(\vec{x}) + \nabla^2 f(\vec{x})\Delta\vec{x} = 0 \\ \Rightarrow \Delta\vec{x} &= -(\nabla^2 f(\vec{x}))^{-1}\nabla f(\vec{x})\end{aligned}$$

**Key Difference:** While Gradient Descent assumes the local geometry is a sphere (linear approximation), Newton's Method accounts for the local **curvature** (ellipsoidal geometry) given by the Hessian  $\nabla^2 f(\vec{x})$ .

### 1.3 Constrained Optimization

We consider the primal optimization problem with constraints:

$$\begin{aligned}\text{minimize} \quad & f_0(\vec{x}) \\ \text{subject to} \quad & f_i(\vec{x}) \leq 0, \quad i = 1, \dots, m \\ & h_i(\vec{x}) = 0, \quad i = 1, \dots, p\end{aligned}$$

*Note:  $f_0$  is not required to be convex yet.*

#### 1.3.1 The Lagrangian & Dual Function

To handle constraints, we define the **Lagrangian**  $L$ :

$$L(\vec{x}, \vec{\lambda}, \vec{\nu}) = f_0(\vec{x}) + \sum_{i=1}^m \lambda_i f_i(\vec{x}) + \sum_{i=1}^p \nu_i h_i(\vec{x}) \quad (1.9)$$

We formulate the **Lagrange dual function**  $g(\vec{\lambda}, \vec{\nu})$  by minimizing  $L$  over  $\vec{x}$ :

$$g(\vec{\lambda}, \vec{\nu}) = \inf_{\vec{x}} L(\vec{x}, \vec{\lambda}, \vec{\nu}) \quad (1.10)$$

#### Properties of the Dual Function

- $g(\vec{\lambda}, \vec{\nu})$  is **always concave**, even if the primal problem is not convex.
- It provides a **lower bound** on the optimal primal value  $p^*$ :  $g(\vec{\lambda}, \vec{\nu}) \leq p^*$ .

This leads to the **Dual Problem**: Find the best lower bound by maximizing  $g$ .

$$\text{maximize } g(\vec{\lambda}, \vec{\nu}) \quad \text{subject to } \vec{\lambda} \succeq 0$$

### 1.3.2 Strong Duality & Slater's Condition

- **Weak Duality:**  $d^* \leq p^*$  (Always true). The difference  $p^* - d^*$  is the *duality gap*.
- **Strong Duality:**  $d^* = p^*$  (Gap is zero). This allows solving the dual problem to find the primal solution.

Strong duality is guaranteed if **Slater's Condition** holds:

#### Theorem: Slater's Condition

For a **convex** optimization problem, strong duality holds if there exists a strictly feasible point  $\vec{x}$  such that:

$$f_i(\vec{x}) < 0 \quad \forall i = 1, \dots, m \quad \text{and} \quad A\vec{x} = \vec{b} \quad (1.11)$$

**Refinement:** If constraints are affine (linear), touching the boundary ( $f_i(\vec{x}) \leq 0$ ) is allowed. Only non-linear constraints require strict inequality.

### 1.3.3 Karush-Kuhn-Tucker (KKT) Conditions

If strong duality holds, any optimal pair  $(\vec{x}^*, \vec{\lambda}^*, \nu^*)$  **must** satisfy the KKT conditions:

1. **Primal Feasibility:**

$$f_i(\vec{x}^*) \leq 0, \quad h_i(\vec{x}^*) = 0$$

2. **Dual Feasibility:**

$$\vec{\lambda}^* \succeq 0$$

3. **Complementary Slackness:** (Crucial for SVMs!)

$$\lambda_i^* \cdot f_i(\vec{x}^*) = 0$$

*Meaning: Either a constraint is active ( $f_i(\vec{x}) = 0$ ) or its multiplier is zero ( $\lambda_i = 0$ ).*

4. **Stationarity:** Gradient of Lagrangian is zero.

$$\nabla f_0(\vec{x}^*) + \sum \lambda_i^* \nabla f_i(\vec{x}^*) + \sum \nu_i^* \nabla h_i(\vec{x}^*) = 0$$

**Conclusion:** For convex problems with strong duality, KKT conditions are necessary and sufficient for optimality. Finding a point that satisfies them means we found the global optimum.

## 2 Support Vector Machines (SVM)

Just like other classifiers (Neural Nets, Nearest Neighbor, etc.), the goal of SVMs is to draw a linear line (decision boundary) to separate classes. But instead of drawing any sufficient line to separate the classes, SVMs aim to find a unique decision boundary that **maximizes the margin (distance)** between each class. The solution to this problem is unique and depends only on the features that are close to the decision boundary.

### 2.1 Hard Margin Problems

The hard margin SVM needs linearly separable classes. Lets assume there is an affine function

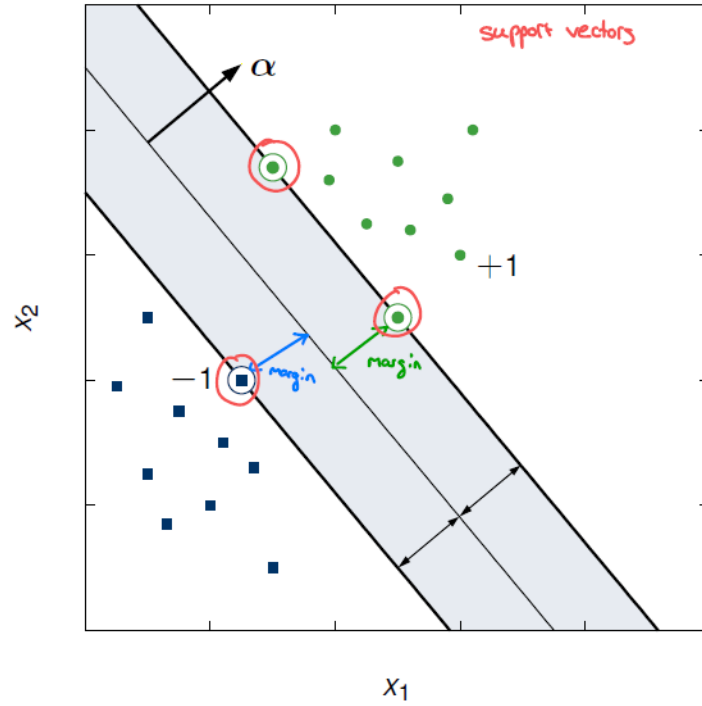


Figure 7: Hard margin SVM

defined as:

$$f(x) = \vec{\alpha}^T x + \alpha_0 \quad (2.1)$$

Where  $\vec{\alpha}$  is the normal vector to the decision boundary and  $\alpha_0$  is some sort of bias. For any point  $x$  on the decision boundary, it holds that  $f(x) = 0$ .

There are three major points we need to think about to end up with a nice optimization problem:

#### 1. Introduce margin constraints:

We need to ensure that all points are classified correctly and therefore lie outside the margin. Therefore, we introduce the following constraints:

- For points of class +1:  $f(x) \geq 1$
- For points of class -1:  $f(x) \leq -1$

This means that for a sample point  $x_i$  with the label  $y_i = +1$ , it holds that  $f(x_i) \geq 1$ . Similarly, for a sample point  $x_j$  with the label  $y_j = -1$ , it holds that  $f(x_j) \leq -1$ . These two constraints can be combined into one single constraint ( $y \in \{+1, -1\}$ ):

$$y_i f(x_i) - 1 = y_i (\vec{\alpha}^T x_i + \alpha_0) - 1 \geq 0 \quad \forall i \quad (2.2)$$

## 2. Define the margin:

The margin is defined as the distance between the decision boundary and the closest points from either class. To compute the margin width, we take a sample from each class that lies exactly on the margin (i.e., satisfies  $y_i f(x_i) - 1 = 0 \quad \forall i$ ) and subtract them from each other. When we project the resulting vector onto the normalized normal vector of the hyperplane, we get the margin width:

$$\text{width} = \frac{\vec{\alpha}}{\|\vec{\alpha}\|_2} \cdot (\vec{x}_{y=+1} - \vec{x}_{y=-1}) \quad (2.3)$$

Now we multiply this out:

$$\text{width} = \frac{1}{\|\vec{\alpha}\|_2} (\vec{\alpha}^T \vec{x}_{y=+1} - \vec{\alpha}^T \vec{x}_{y=-1}) \quad (2.4)$$

We know from our margin constraints defined in step 1 that for support vectors (points on the margin), the inequality becomes an equality:

- For the positive support vector  $\vec{x}_{y=+1}$ :

$$\vec{\alpha}^T \vec{x}_{y=+1} + \alpha_0 = 1 \quad \Rightarrow \quad \vec{\alpha}^T \vec{x}_{y=+1} = 1 - \alpha_0$$

- For the negative support vector  $\vec{x}_{y=-1}$ :

$$\vec{\alpha}^T \vec{x}_{y=-1} + \alpha_0 = -1 \quad \Rightarrow \quad \vec{\alpha}^T \vec{x}_{y=-1} = -1 - \alpha_0$$

Substituting these expressions back into the width equation:

$$\text{width} = \frac{1}{\|\vec{\alpha}\|_2} ((1 - \alpha_0) - (-1 - \alpha_0)) \quad (2.5)$$

$$= \frac{1}{\|\vec{\alpha}\|_2} (1 - \alpha_0 + 1 + \alpha_0) \quad (2.6)$$

$$= \frac{2}{\|\vec{\alpha}\|_2} \quad (2.7)$$

## 3. Minimize the norm:

Since we want to **maximize** the margin width  $\frac{2}{\|\vec{\alpha}\|_2}$ , this is mathematically equivalent to **minimizing** the length of the normal vector  $\|\vec{\alpha}\|_2$ . For mathematical convenience (to make derivatives easier later), we minimize the squared norm:

### Primal Optimization Problem (Hard Margin)

$$\text{minimize} \quad \frac{1}{2} \|\vec{\alpha}\|_2^2 \quad \text{subject to} \quad y_i (\vec{\alpha}^T x_i + \alpha_0) \geq 1 \quad \forall i$$

## 2.2 Soft Margin Problems

In real world applications, data is often not perfectly linearly separable. The Soft Margin SVM relaxes the hard margin constraints by allowing some points to violate the margin or even be misclassified.

Therefore, we introduce slack variables  $\xi_i \geq 0$  for each training sample  $x_i$ , which measure the degree of misclassification:

- $\xi_i = 0$ : point is correctly classified and outside or on the margin
- $0 < \xi_i \leq 1$ : point is inside the margin but still correctly classified



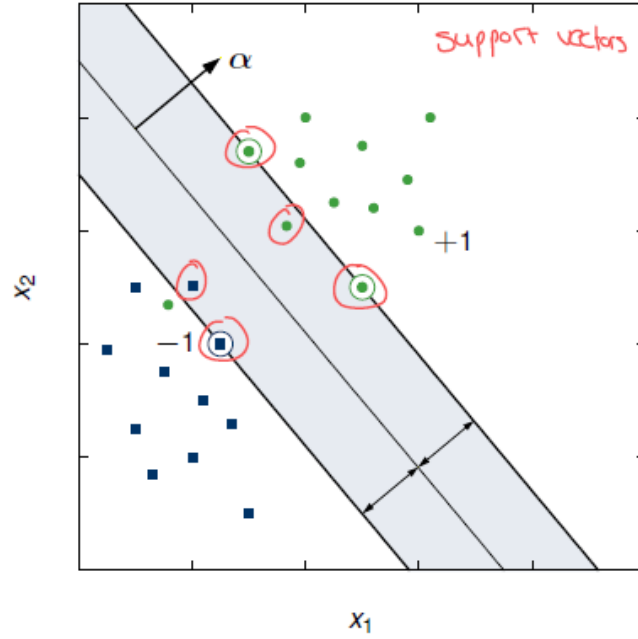


Figure 8: Soft margin SVM

- $\xi_i > 1$ : point is misclassified

The margin constraints are now relaxed to:

$$y_i(\vec{\alpha}^T x_i + \alpha_0) \geq 1 - \xi_i \quad \forall i \quad (2.8)$$

The primal optimization problem becomes:

**Primal Optimization Problem (Soft Margin)**

$$\text{minimize} \quad \frac{1}{2} \|\vec{\alpha}\|_2^2 + \mu \sum_{i=1}^n \xi_i \quad \text{subject to} \quad -(y_i(\vec{\alpha}^T x_i + \alpha_0) - 1 + \xi_i) \leq 0, \quad -\xi_i \leq 0 \quad \forall i$$

where  $\mu > 0$  is a hyperparameter that controls the trade-off between maximizing the margin and minimizing the classification error.