

# Optimization Techniques for Machine Learning

AMLZC326 · #14 Machine Learning II

Anshid Aboobacker

# TABLE OF CONTENTS

- 1 Linear Separable Data
- 2 Geometry of Hyperplanes
- 3 Margin and Supporting Hyperplanes
- 4 Hard Margin SVM
- 5 Soft Margin SVM
- 6 Hinge Loss

# THE CENTRAL QUESTION

Among all separating lines,  
**why choose one over another?**

# LEARNING OBJECTIVES

By the end of this lecture you should be able to:

- Motivate the maximum-margin classifier geometrically
- Derive the hard margin SVM as a constrained quadratic program
- Extend to soft margin SVM with slack variables and the regularisation parameter  $C$
- Interpret hinge loss as the convex surrogate for 0–1 misclassification loss

# LINEAR CLASSIFICATION

- We want to separate data into two classes:  $y \in \{+1, -1\}$

$$f(x) = \text{sign}(w^T x + b)$$

- Decision boundary:

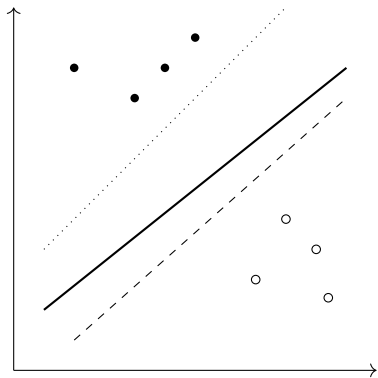
$$w^T x + b = 0$$

# GEOMETRIC INTERPRETATION

- $w^T x + b > 0 \Rightarrow +1$
- $w^T x + b < 0 \Rightarrow -1$

**A hyperplane separates the data**

# LINEARLY SEPARABLE DATA



Many separating lines exist

## KEY OBSERVATION

All these classifiers perfectly separate the data

**Which one should we choose?**

# SCALE INVARIANCE

- The classifier depends on:

$$\text{sign}(w^T x + b)$$

- If we scale:

$$(w, b) \rightarrow (\lambda w, \lambda b)$$

- Then:

$$\text{sign}(w^T x + b) = \text{sign}(\lambda w^T x + \lambda b)$$

**Same classifier, different parameters**

# NORMALIZATION IDEA

- Since scaling does not change classification,
- we can fix a scale:

$$y_i(w^T x_i + b) \geq 1$$

- This removes ambiguity

# WHAT JUST HAPPENED?

- Classification depends only on sign
- But geometry depends on magnitude

**We are now ready to define margin**

# TABLE OF CONTENTS

- 1 Linear Separable Data
- 2 Geometry of Hyperplanes**
- 3 Margin and Supporting Hyperplanes
- 4 Hard Margin SVM
- 5 Soft Margin SVM
- 6 Hinge Loss

# WHAT IS A HYPERPLANE?

- Decision boundary:

$$w^T x + b = 0$$

- In  $\mathbb{R}^2$ : a line
- In  $\mathbb{R}^3$ : a plane
- In  $\mathbb{R}^d$ : a hyperplane

**This is the set of points where classification changes**

# WHY IS $w$ PERPENDICULAR?

Let  $x_a, x_b$  lie on the hyperplane:

$$w^T x_a + b = 0, \quad w^T x_b + b = 0$$

Subtract:

$$w^T (x_a - x_b) = 0$$

- $(x_a - x_b)$  lies along the hyperplane
- Dot product = 0  $\Rightarrow$  orthogonal

**$w$  is normal to the hyperplane**

# GEOMETRIC MEANING OF $w$

- Direction of  $w$ :
  - ▶ Points toward positive class
- Magnitude of  $w$ :
  - ▶ Controls how fast values change

**$w$  controls geometry, not just classification**

# DISTANCE FROM A POINT TO HYPERPLANE

$$D = \frac{|w^T x + b|}{\|w\|}$$

- Numerator: raw score
- Denominator: normalization by scale

**Distance depends inversely on  $\|w\|$**

# SIGNED DISTANCE

$$\text{Signed distance} = \frac{w^T x + b}{\|w\|}$$

- Positive  $\Rightarrow$  one side
- Negative  $\Rightarrow$  other side

**This connects geometry to classification**

# KEY INSIGHT

- Classification uses:

$$\text{sign}(w^T x + b)$$

- Distance uses:

$$\frac{w^T x + b}{\|w\|}$$

**Only  $\|w\|$  controls margin**

# TRANSITION TO MARGIN

We now know:

- $w$  defines direction
- $\|w\|$  defines scale of distance

**So what is the best hyperplane?**

# TABLE OF CONTENTS

- 1 Linear Separable Data
- 2 Geometry of Hyperplanes
- 3 Margin and Supporting Hyperplanes**
- 4 Hard Margin SVM
- 5 Soft Margin SVM
- 6 Hinge Loss

# FROM CLASSIFICATION TO GEOMETRY

We don't just want to classify correctly

**We want the boundary to be as far as possible from all points**

# NORMALIZATION RECALL

We fixed scale using:

$$y_i(w^T x_i + b) \geq 1$$

- Closest points satisfy:

$$y_i(w^T x_i + b) = 1$$

**These define the margin boundaries**

# SUPPORTING HYPERPLANES

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

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

- Parallel to decision boundary
- Touch closest points

**These are the supporting hyperplanes**

# WHAT IS THE MARGIN?

Distance between the two supporting  
hyperplanes

**We want to maximize this distance**

# DERIVING THE MARGIN

Take one point on each hyperplane:

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

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

Subtract:

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

Project onto direction of  $w$ :

$$\text{distance} = \frac{2}{\|w\|}$$

# MARGIN FORMULA

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

Maximize margin  $\iff$  minimize  $\|w\|$

# OPTIMIZATION INSIGHT

Instead of:

$$\max \frac{2}{\|w\|}$$

We solve:

$$\min \frac{1}{2} \|w\|^2$$

- Same solution
- Easier to optimize

# WHO DEFINES THE MARGIN?

Only the closest points matter

**These points are called support vectors**

## KEY INSIGHT

- Far away points do not affect the boundary
- Only boundary-touching points matter

**SVM depends only on a few critical points**

We now know:

- What to optimize
- What constraints to enforce

**Let's turn this into an optimization  
problem**

# TABLE OF CONTENTS

- 1 Linear Separable Data
- 2 Geometry of Hyperplanes
- 3 Margin and Supporting Hyperplanes
- 4 Hard Margin SVM**
- 5 Soft Margin SVM
- 6 Hinge Loss

# WHAT HAVE WE BUILT SO FAR?

- Linear classifier:

$$f(x) = \text{sign}(w^T x + b)$$

- Margin:

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

- Goal:

Maximize margin

# REWRITING THE GOAL

$$\max \frac{2}{\|w\|}$$

Equivalent to:

$$\min \frac{1}{2} \|w\|^2$$

**Cleaner, convex objective**

# WHAT ABOUT CORRECT CLASSIFICATION?

We must ensure all points are classified correctly:

$$y_i(w^T x_i + b) > 0$$

After normalization:

$$y_i(w^T x_i + b) \geq 1$$

**Encodes both correctness + margin**

# HARD MARGIN SVM FORMULATION

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

subject to  $y_i(w^T x_i + b) \geq 1$

**This is the hard margin SVM**

# UNDERSTANDING THE CONSTRAINTS

- If  $y_i = +1$ :

$$w^T x_i + b \geq 1$$

- If  $y_i = -1$ :

$$w^T x_i + b \leq -1$$

**Points lie outside margin boundaries**

# GEOMETRIC MEANING

- Decision boundary:

$$w^T x + b = 0$$

- Margin boundaries:

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

**All points lie outside the margin**

## IMPORTANT OBSERVATION

Only points satisfying

$$y_i(w^T x_i + b) = 1$$

matter

**These are the support vectors**

# NATURE OF THE PROBLEM

- Objective: quadratic
- Constraints: linear

## **Convex optimization problem**

**Unique global solution**

## LIMITATION OF HARD MARGIN

What if data is not perfectly separable?

**Hard margin fails completely**

We need flexibility

**Allow some violations of constraints**

# TABLE OF CONTENTS

- 1 Linear Separable Data
- 2 Geometry of Hyperplanes
- 3 Margin and Supporting Hyperplanes
- 4 Hard Margin SVM
- 5 Soft Margin SVM**
- 6 Hinge Loss

# WHY HARD MARGIN FAILS

Real-world data is noisy

**Perfect separation may not exist**

# THE PROBLEM

- Even one outlier can break hard margin SVM

**We need to allow mistakes**

# INTRODUCING SLACK VARIABLES

We relax constraints:

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

$\xi_i$  **measures violation**

# MEANING OF $\xi_i$

- $\xi_i = 0 \rightarrow$  correctly classified, outside margin
- $0 < \xi_i < 1 \rightarrow$  inside margin
- $\xi_i > 1 \rightarrow$  misclassified

**Slack = how bad the violation is**

# NEW OBJECTIVE

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

- First term  $\rightarrow$  maximize margin
- Second term  $\rightarrow$  penalize violations

# ROLE OF $C$

- Large  $C$ :
  - ▶ Penalize mistakes heavily
  - ▶ Smaller margin
- Small  $C$ :
  - ▶ Allow more violations
  - ▶ Larger margin

**$C$  controls the trade-off**

# GEOMETRIC PICTURE

- Some points allowed inside margin
- Some points even misclassified

**We trade perfection for robustness**

# SOFT MARGIN SVM FORMULATION

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

subject to  $y_i(w^T x_i + b) \geq 1 - \xi_i$

$$\xi_i \geq 0$$

We are balancing:

- Margin size
- Classification errors

**This is a trade-off problem**

Can we write this without slack variables?

**Yes — using a loss function**

# TABLE OF CONTENTS

- 1 Linear Separable Data
- 2 Geometry of Hyperplanes
- 3 Margin and Supporting Hyperplanes
- 4 Hard Margin SVM
- 5 Soft Margin SVM
- 6 Hinge Loss**

# FROM CONSTRAINTS TO LOSS

We had:

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

Can we eliminate  $\xi_i$ ?

# KEY OBSERVATION

Constraint implies:

$$\xi_i \geq 1 - y_i(w^T x_i + b)$$

**Minimum  $\xi_i$  is:**

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b))$$

# HINGE LOSS DEFINITION

$$L_i = \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

- $0 \rightarrow$  correct and outside margin
- $> 0 \rightarrow$  violation

# INTERPRETATION

- If  $y_i(w^T x_i + b) \geq 1$ :  
 $L_i = 0$
- If  $y_i(w^T x_i + b) < 1$ :  
 $L_i > 0$

**Loss depends on margin, not just correctness**

# FINAL OBJECTIVE

$$J(w, b) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b))$$

- First term  $\rightarrow$  regularization
- Second term  $\rightarrow$  hinge loss

SVM = Regularization + Loss

Minimize Complexity + Error

# WHY HINGE LOSS?

- Penalizes margin violations
- Encourages large margin
- Convex  $\rightarrow$  easy to optimize

**Designed for maximum margin classification**

Classification is not enough

**Confidence (margin) matters**

# CONNECTION TO MODERN ML

- Linear regression  $\rightarrow$  squared loss
- Logistic regression  $\rightarrow$  log loss
- SVM  $\rightarrow$  hinge loss

**Same framework, different loss**

## FINAL TAKEAWAY

SVM is not a special trick

**It is a principled optimization problem**

Thank you :)