# Week 11: Support Vector Machines

Maximum Margin Classification

2025

## Learning Objectives
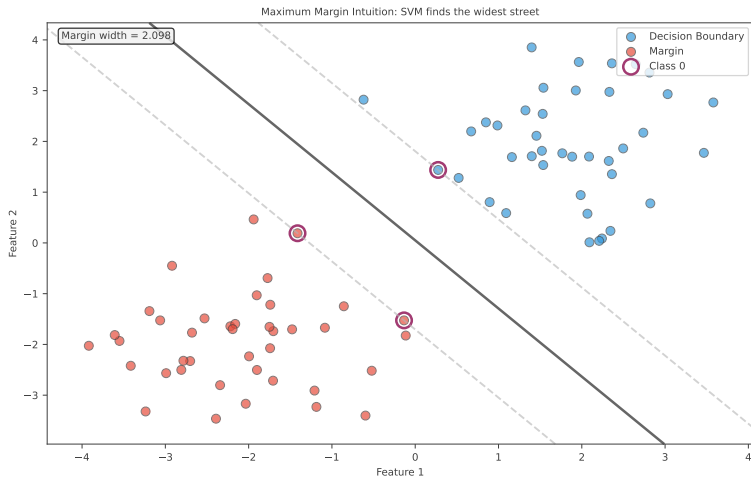
By the end of this week, you will be able to:

1. **Understand** the maximum margin principle and its geometric interpretation
2. **Derive** the soft margin SVM optimization problem
3. **Apply** the kernel trick to create non-linear decision boundaries
4. **Compare** different kernel functions (linear, polynomial, RBF, sigmoid)
5. **Tune** hyperparameters C and gamma for optimal performance
6. **Identify** support vectors and understand their role
7. **Extend** SVM to regression (SVR) and multi-class problems

**Support Vector Machines: Elegant theory meets practical power**
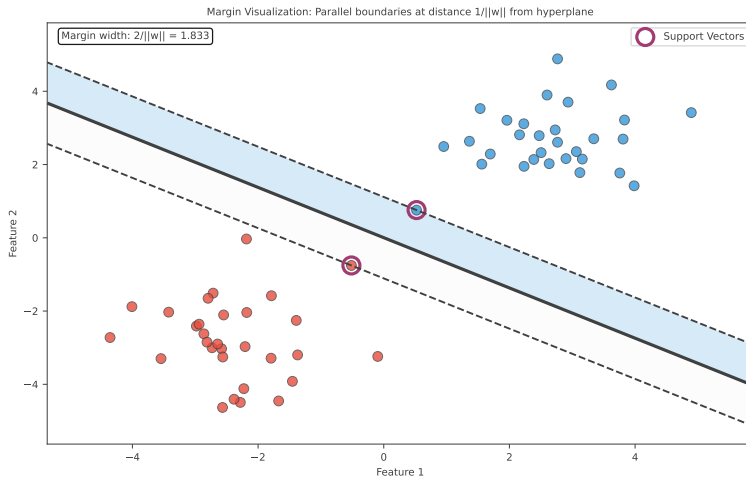
# Part I: Maximum Margin Principle

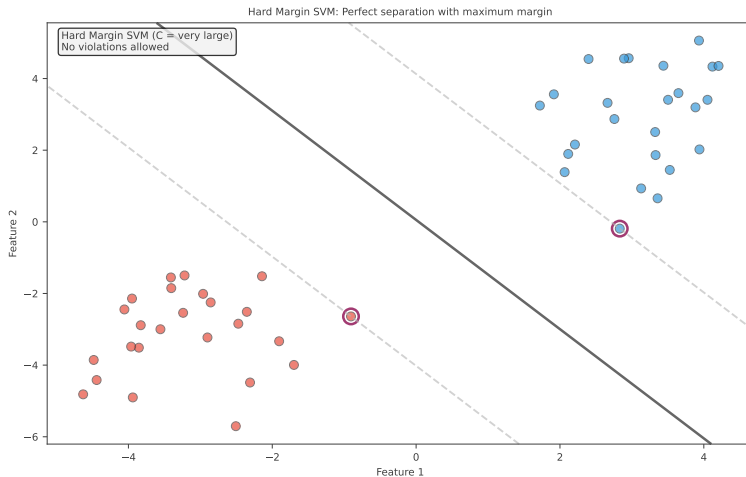*Maximize the gap between classes*

# Maximum Margin Intuition



Maximum Margin Intuition: SVM finds the widest street

**Among all separating hyperplanes, choose the one with maximum margin**

# Margin Visualization



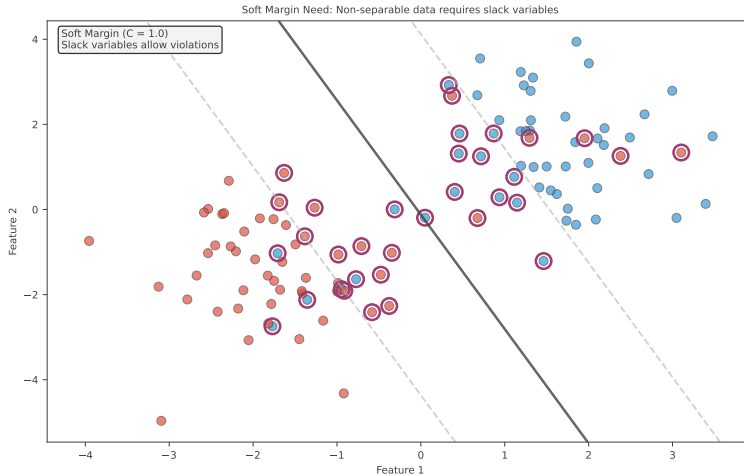Margin Visualization: Parallel boundaries at distance $1/\|w\|$ from hyperplane

**Margin:** Perpendicular distance from hyperplane to nearest points

# Hard Margin SVM



Hard Margin SVM: Perfect separation with maximum margin

Hard Margin SVM (C = very large)
No violations allowed

**Hard margin: Perfect separation required, all points outside margin**

# Why We Need Soft Margin



Soft Margin Need: Non-separable data requires slack variables

**Real data rarely perfectly separable: Allow some violations with penalty**

Support Vectors on Margin: Only these points matter (Dataset 1)

**Support vectors: Points on or inside margin that define the hyperplane**

# Decision Function Values



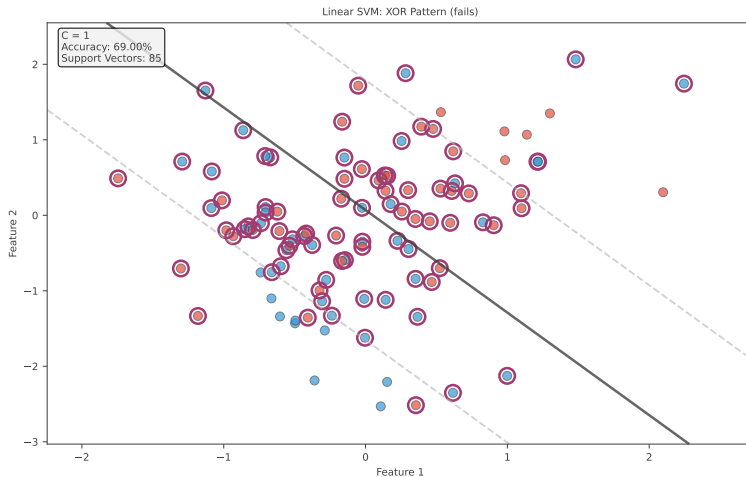Decision Function: Signed distance from hyperplane

**Signed distance from hyperplane indicates confidence**

# Linear SVM: When It Works



Linear SVM excels on linearly separable or near-separable data

Linear SVM: XOR Pattern (fails)

C = 1
Accuracy: 69.00%
Support Vectors: 85

**Linear boundary cannot solve XOR problem: Need non-linear solution**

Soft Margin SVM: C = 0.01

C = 0.01
Margin: 3.480
SVs: 61
Acc: 95.38%

**Very soft margin: Many support vectors, high tolerance for violations**

Soft Margin SVM: C = 0.1

C = 0.1
Margin: 1.932
SVs: 27
Acc: 96.15%

**Soft margin: Flexible boundary, tolerates some misclassifications**

# C Parameter Effect

C = 1 (Balanced Margin)



Soft Margin SVM: C = 1

C = 1
Margin: 1.373
SVs: 17
Acc: 96.15%

**Balanced margin: Default setting, moderate tolerance**

# C Parameter Effect

C = 10 (Stricter Margin)



Soft Margin SVM: C = 10

C = 10
Margin: 1.164
SVs: 15
Acc: 95.38%

**Stricter margin: Fewer violations allowed, more support vectors on margin**

# C Parameter Effect

C = 100 (Hard Margin)



Soft Margin SVM: C = 100

C = 100
Margin: 1.121
SVs: 14
Acc: 95.38%

**Hard margin: Very strict, approaches perfect separation**

Soft Margin SVM: C = 1000

C = 1000
Margin: 1.120
SVs: 14
Acc: 95.38%

**Very hard margin: Minimal violations, risk of overfitting**

# Support Vectors vs C



Support Vector Count vs C: Decreasing trend

Larger C: fewer violations allowed
Fewer support vectors

**Fewer support vectors as C increases: Fewer violations allowed**

# Bias-Variance Trade-Off



Bias-Variance Trade-off: Error vs C

**C controls bias-variance: Small C = high bias, Large C = high variance**

## Why We Need Kernels

**Problem**

- Linear boundaries cannot solve non-linear problems
- XOR, circles, spirals are not linearly separable
- Need curved, complex decision boundaries

**Approach**

- Map data to higher-dimensional space
- Linear separation may exist in transformed space

**Solution**

- Kernel functions enable non-linear boundaries
- Implicit feature mapping via kernel trick

**Kernels solve the non-linearity problem without explicit high-dimensional computation**

# Kernel Comparison



**Different kernels create dramatically different decision boundaries**

# Polynomial Kernel: Degree 1 (Linear Boundary)



Degree 1 is equivalent to linear kernel: straight decision boundary

# Polynomial Kernel: Degree 3 (Curved Boundary)



Polynomial Kernel: degree = 3

Degree: 3
Accuracy: 92.00%
SVs: 45

**Degree 3 creates moderately complex curved boundaries**

Polynomial Kernel: degree = 5

Degree: 5
Accuracy: 97.33%
SVs: 30

**Higher degrees create very complex boundaries: flexible but risk overfitting**

# RBF Kernel: Gamma Comparison



gamma=0.01, acc=83.33%  gamma=0.1, acc=86.00%  gamma=0.5, acc=92.00%  gamma=1, acc=94.67%

gamma=5, acc=98.67%  gamma=10, acc=98.67%  gamma=50, acc=100.00%  gamma=100, acc=100.00%

**Gamma controls RBF complexity: Small = smooth, Large = complex**

## The Kernel Trick: Feature Space Mapping

**Problem**

- Some patterns cannot be separated by any linear boundary
- Example: Circular patterns in 2D

**Approach**

- Map to higher dimensions where linear separation exists
- Example: $(x_1, x_2) \rightarrow (x_1, x_2, x_1^2, x_1 x_2, x_2^2)$
- 2D circles become linearly separable in 5D

**Solution: Kernel Trick**

- Kernel function: $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$
- Compute inner product directly
- Never explicitly compute $\phi(\mathbf{x})$
- Efficient: $O(d)$ instead of $O(D)$ where $D \gg d$

**Kernel trick: Implicit high-dimensional mapping without computational cost**

# Feature Space Mapping



Original 2D Space (not linearly separable)

Transformed 3D Space (linearly separable)

**Kernel trick: Implicitly map to higher dimensions without computing explicitly**

## Part II: Mathematical Foundations

*Formalize the maximum margin concept*

## SVM Optimization Problem

**Hard Margin (Primal)**
Minimize:

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

Subject to:

$$y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1, \quad \forall i$$

**Margin Width**

$$\text{margin} = \frac{2}{\|\mathbf{w}\|}$$

Maximizing margin = Minimizing $\|\mathbf{w}\|^2$

**Soft Margin (Practical)**
Minimize:

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

Subject to:

$$y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$

$C$ controls trade-off:

- Large $C$: Hard margin
- Small $C$: Soft margin

**Soft margin allows misclassifications with penalty**

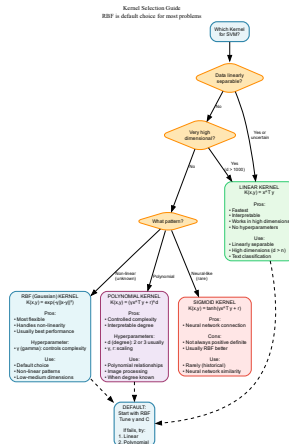**Hard Margin vs Soft Margin SVM**

| HARD MARGIN SVM | |
|---|---|
| Requirement | Data must be linearly separable<br>No points in margin allowed |
| Objective | Maximize margin<br>min 1/2 w squared |
| Constraints | y_i times w^T x_i + b greater or equal 1<br>for ALL points i |
| Support Vectors | Points exactly on margin<br>y_i times w^T x_i + b = 1 |
| C Parameter | Not applicable<br>(infinite C implicitly) |
| Advantages | Unique solution<br>Maximum separation<br>Simple formulation |
| Disadvantages | FAILS if not separable<br>Sensitive to outliers<br>Not robust<br>Rarely applicable |
| Use When | Perfect separation<br>Clean data<br>No noise |

| SOFT MARGIN SVM | |
|---|---|
| Requirement | Works with any data<br>Allows margin violations |
| Objective | Maximize margin + penalize violations<br>min 1/2 w squared + C sum of slack |
| Constraints | y_i times w^T x_i + b greater or equal 1 - slack_i<br>slack_i greater or equal 0 |
| Support Vectors | Points on margin OR inside<br>alpha_i greater than 0 |
| C Parameter | Controls trade-off<br>Large C to hard margin<br>Small C to soft margin |
| Advantages | Always has solution<br>Robust to outliers<br>Handles noise<br>Practical |
| Disadvantages | Requires tuning C<br>More hyperparameters<br>Less unique depends on C |
| Use When | Real-world data always<br>Noise present<br>Overlapping classes<br>Default choice |

In practice: ALWAYS use soft margin (set finite C)
Hard margin is theoretical concept, soft margin is practical necessity

---

**Soft margin is practical necessity for real-world data**

Kernel Selection Guide
RBF is default choice for most problems

Decision tree for choosing appropriate kernel based on data characteristics

## Kernel Functions

**Linear Kernel**

$$K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$$

**Polynomial Kernel**

$$K(\mathbf{x}, \mathbf{y}) = (\gamma \mathbf{x}^T \mathbf{y} + r)^d$$

**RBF (Gaussian) Kernel**

$$K(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2)$$

**Sigmoid Kernel**

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\gamma \mathbf{x}^T \mathbf{y} + r)$$

**Kernel Trick**
Decision function:

$$f(\mathbf{x}) = \sum_{i=1}^{n} \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

Only sum over **support vectors**!

**Properties**

- Linear: Fast, interpretable
- Polynomial: Controlled complexity
- RBF: Most flexible, default choice
- Sigmoid: Neural network-like

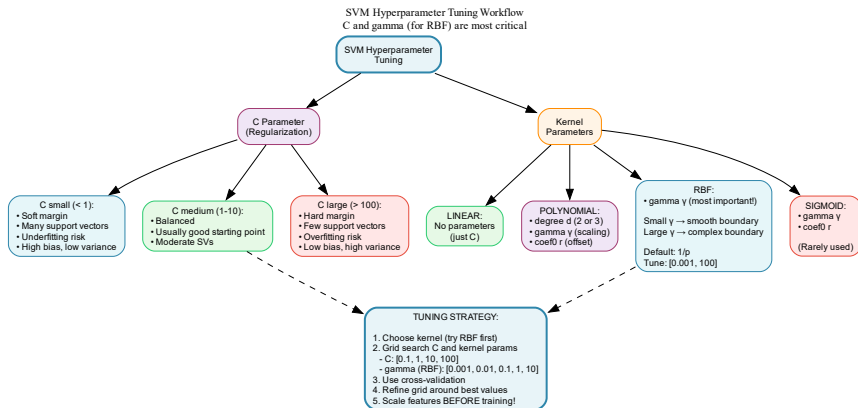**Kernels enable non-linear boundaries without explicit feature mapping**

**Training: Solve QP — Prediction: Kernel evaluation with support vectors only**

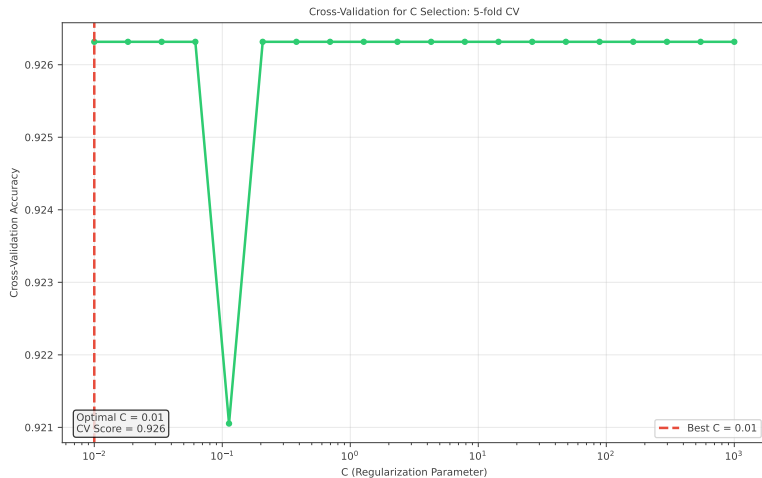## Part III: Hyperparameters & Practical Use
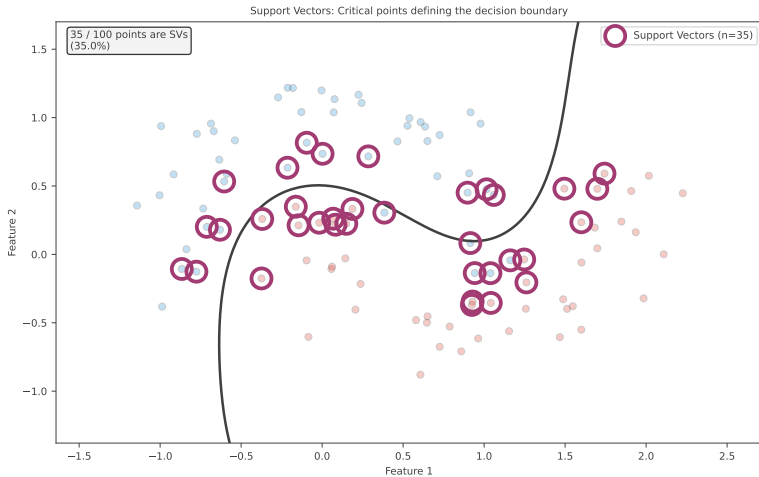
*Tuning and application*

# Hyperparameter Tuning



SVM Hyperparameter Tuning Workflow
C and gamma (for RBF) are most critical

SVM Hyperparameter Tuning

C Parameter (Regularization)

Kernel Parameters

C small (< 1):
• Soft margin
• Many support vectors
• Underfitting risk
• High bias, low variance

C medium (1-10):
• Balanced
• Usually good starting point
• Moderate SVs

C large (> 100):
• Hard margin
• Few support vectors
• Overfitting risk
• Low bias, high variance

LINEAR:
No parameters
(just C)

POLYNOMIAL:
• degree d (2 or 3)
• gamma γ (scaling)
• coef0 r (offset)

RBF:
• gamma γ (most important!)

Small γ → smooth boundary
Large γ → complex boundary

Default: 1/p
Tune: [0.001, 100]

SIGMOID:
• gamma γ
• coef0 r

(Rarely used)

TUNING STRATEGY:
1. Choose kernel (try RBF first)
2. Grid search C and kernel params
   - C: [0.1, 1, 10, 100]
   - gamma (RBF): [0.001, 0.01, 0.1, 1, 10]
3. Use cross-validation
4. Refine grid around best values
5. Scale features BEFORE training!

**C and gamma (for RBF) are most critical hyperparameters**

# Cross-Validation for C



Cross-Validation for C Selection: 5-fold CV

**Use CV to find optimal C value**

Support Vectors: Critical points defining the decision boundary

**Support vectors are the critical points that define the margin**

# Support Vector Regression (SVR)



Support Vector Regression: Epsilon-tube defines margin

**SVR: Epsilon-insensitive loss creates margin for regression**

Multi-Class SVM Strategies
Conceptual overview - sklearn handles automatically

**SVM naturally binary, extended to multi-class via One-vs-Rest or One-vs-One**

## Summary: Key Concepts

**Core Ideas**

- Maximum margin principle
- Support vectors define boundary
- Soft margin with C parameter
- Kernel trick for non-linearity

**Optimization**

- Quadratic programming
- Convex optimization
- Global optimum guaranteed

**Kernels**

- Linear: Default, interpretable
- Polynomial: Controlled complexity
- RBF: Most flexible, default choice
- Sigmoid: Neural network-like

**Hyperparameters**

- C: Margin softness
- gamma: RBF complexity
- degree: Polynomial complexity

**SVM: Powerful, elegant, kernel-based classification and regression**

**Next Week: Naive Bayes**
Shift from discriminative to generative models:

- Probabilistic classification
- Bayes theorem application
- Gaussian, Multinomial, Bernoulli variants
- Independence assumption
- Text classification
- Fast training and prediction

**Key Difference**
SVM models boundaries. Naive Bayes models class distributions.

**From maximum margin to maximum likelihood**

Thank you

Questions?