

Support Vector Machine (SVM) Complete Cheatsheet

1. Problem Setup

Goal: Find a hyperplane that separates two classes with maximum margin

Training Data:

- Dataset: $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- $x_i \in \mathbb{R}^d$ (feature vectors)
- $y_i \in \{-1, +1\}$ (class labels)

Hyperplane Equation:

$$w \cdot x + b = 0$$

- w : weight vector (normal to hyperplane)
- b : bias term
- x : input vector

Decision Function:

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

2. Linear SVM - Hard Margin

Geometric Margin

Distance from point x to hyperplane:

$$\text{distance} = |w \cdot x + b| / ||w||$$

Margin (γ): Distance from hyperplane to nearest point

$$\gamma = \min(|w \cdot x_i + b| / ||w||) \text{ for all } i$$

Optimization Problem (Primal Form)

Maximize margin:

$$\max (1/||w||) \iff \min (1/2)||w||^2$$

Subject to constraints:

$$y_i(w \cdot x_i + b) \geq 1 \quad \text{for all } i = 1, \dots, n$$

This ensures all points are correctly classified with margin $\geq 1/||w||$

Complete Primal Problem:

$$\begin{aligned} \min_{w,b} \quad & (1/2)||w||^2 \\ \text{subject to: } & y_i(w \cdot x_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

3. Linear SVM - Soft Margin

Why Soft Margin?

- Data may not be perfectly separable
- Allow some misclassifications
- Introduces slack variables ξ_i

Slack Variables

$$\xi_i \geq 0 \quad (\text{amount of constraint violation for point } i)$$

Interpretation:

- $\xi_i = 0$: point is correctly classified
- $0 < \xi_i < 1$: point is inside margin but correctly classified
- $\xi_i \geq 1$: point is misclassified

Optimization Problem

Primal Form:

$$\begin{aligned} \min_{w,b,\xi} \quad & (1/2)||w||^2 + C \cdot \sum \xi_i \\ \text{subject to: } & \\ & y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad \text{for all } i \\ & \xi_i \geq 0 \quad \text{for all } i \end{aligned}$$

C (Regularization Parameter):

- Large C: prioritize correct classification (small margin, low bias, high variance)
- Small C: prioritize large margin (allow misclassifications, high bias, low variance)

4. Dual Formulation (Lagrangian)

Lagrangian Function

Primal variables: w, b, ξ

Dual variables (Lagrange multipliers): α, μ

$$L(w, b, \xi, \alpha, \mu) = (1/2) \|w\|^2 + C \cdot \sum \xi_i - \sum \alpha_i [y_i (w \cdot x_i + b) - 1 + \xi_i] - \sum \mu_i \xi_i$$

KKT Conditions

Stationarity:

$$\begin{aligned} \partial L / \partial w &= 0 \implies w = \sum \alpha_i y_i x_i \\ \partial L / \partial b &= 0 \implies \sum \alpha_i y_i = 0 \\ \partial L / \partial \xi_i &= 0 \implies \alpha_i = C - \mu_i \end{aligned}$$

Primal Feasibility:

$$\begin{aligned} y_i (w \cdot x_i + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned}$$

Dual Feasibility:

$$\begin{aligned} \alpha_i &\geq 0 \\ \mu_i &\geq 0 \end{aligned}$$

Complementary Slackness:

$$\begin{aligned} \alpha_i [y_i (w \cdot x_i + b) - 1 + \xi_i] &= 0 \\ \mu_i \xi_i &= 0 \end{aligned}$$

Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & \sum \alpha_i - (1/2) \sum \sum \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \\ \text{subject to:} \quad & \sum \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \quad \text{for all } i \end{aligned}$$

Why use dual?

- Easier to solve for high-dimensional data
- Enables kernel trick
- Only depends on dot products $x_i \cdot x_j$

5. Support Vectors

Classification of points based on α :

1. **Non-support vectors:** $\alpha_i = 0$
 - Correctly classified, outside margin
 - Don't affect decision boundary
2. **Support vectors on margin:** $0 < \alpha_i < C$
 - Exactly on margin boundary
 - $\xi_i = 0, y_i(w \cdot x_i + b) = 1$
3. **Support vectors inside margin/misclassified:** $\alpha_i = C$
 - Inside margin or misclassified
 - $\xi_i > 0$

Computing w and b :

$$\begin{aligned} w &= \sum (\alpha_i y_i x_i) \quad [\text{sum over support vectors}] \\ b &= y_s - w \cdot x_s \quad [\text{for any support vector } x_s \text{ with } 0 < \alpha_s < C] \end{aligned}$$

Better (averaged):

$$b = (1/N_{sv}) \sum [y_s - w \cdot x_s] \quad [\text{average over all support vectors with } 0 < \alpha < C]$$

6. Kernel Trick

Motivation

Transform data to higher-dimensional space where it's linearly separable

Feature Mapping:

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D \quad (D \gg d)$$

Kernel Function:

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

Computes dot product in high-dimensional space without explicit transformation!

Dual Problem with Kernels

$$\begin{aligned} \max_{\alpha} \quad & \sum \alpha_i - (1/2) \sum \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{subject to:} \quad & \sum \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

Decision Function with Kernels

$$f(x) = \text{sign}(\sum \alpha_i y_i K(x_i, x) + b)$$

7. Common Kernel Functions

1. Linear Kernel

$$K(x_i, x_j) = x_i \cdot x_j$$

- Use when data is linearly separable
- Equivalent to no kernel
- Fast computation

2. Polynomial Kernel

$$K(x_i, x_j) = (\gamma \cdot x_i \cdot x_j + r)^d$$

- γ : kernel coefficient (default: $1/n_{\text{features}}$)
- r : independent term (default: 0)
- d : degree (default: 3)
- Use for non-linear boundaries with polynomial shape

Special case (Homogeneous):

$$K(x_i, x_j) = (x_i \cdot x_j)^d$$

3. Radial Basis Function (RBF/Gaussian) Kernel ★

$$K(x_i, x_j) = \exp(-\gamma ||x_i - x_j||^2)$$

- γ : kernel coefficient (default: $1/n_{\text{features}}$)
- Most popular kernel
- Can handle highly non-linear boundaries
- Maps to infinite-dimensional space

γ parameter effect:

- Large γ : narrow Gaussian, complex decision boundary (high variance)
- Small γ : wide Gaussian, smooth decision boundary (high bias)

4. Sigmoid Kernel

$$K(x_i, x_j) = \tanh(\gamma \cdot x_i \cdot x_j + r)$$

- Similar to neural network activation
- Not always positive semi-definite (can cause issues)

5. Custom Kernels

Mercer's Condition: Kernel must satisfy

$$\iint K(x, x') g(x) g(x') dx dx' \geq 0 \quad \text{for all } g$$

8. Multi-Class Classification

SVMs are binary classifiers. For multi-class:

One-vs-Rest (OvR)

- Train K binary classifiers (K = number of classes)
- Classifier k: class k vs all others
- Predict: class with highest decision function value

Decision:

$$\hat{y} = \underset{k}{\operatorname{argmax}}(w_k \cdot x + b_k)$$

One-vs-One (OvO)

- Train $K(K-1)/2$ binary classifiers
- One for each pair of classes
- Predict: class that wins most pairwise comparisons (voting)

9. SVM Regression (SVR)

 ϵ -insensitive Loss

Goal: Find function $f(x) = w \cdot x + b$ such that $|y_i - f(x_i)| \leq \epsilon$

ϵ -tube: No penalty for errors within $\pm\epsilon$

Optimization Problem

$$\min_{w, b, \xi, \xi^*} (1/2) \|w\|^2 + C \cdot \sum (\xi_i + \xi_i^*)$$

subject to:

$$\begin{aligned} y_i - (w \cdot x_i + b) &\leq \epsilon + \xi_i \\ (w \cdot x_i + b) - y_i &\leq \epsilon + \xi_i^* \\ \xi_i, \xi_i^* &\geq 0 \end{aligned}$$

Dual Form:

$$\max_{\alpha, \alpha^*} -\epsilon \cdot \sum (\alpha_i + \alpha_i^*) + \sum y_i (\alpha_i - \alpha_i^*) - (1/2) \sum \sum (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) K(x_i, x_j)$$

subject to:

$$\begin{aligned} \sum (\alpha_i - \alpha_i^*) &= 0 \\ 0 &\leq \alpha_i, \alpha_i^* \leq C \end{aligned}$$

Prediction:

$$f(x) = \sum (\alpha_i - \alpha_i^*) K(x_i, x) + b$$

10. Key Formulas Summary

Training (Dual Optimization)

$$\begin{aligned} \max_{\alpha} \quad & \sum \alpha_i - (1/2) \sum \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad & \sum \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C \end{aligned}$$

Weight Vector

$$w = \sum \alpha_i y_i \phi(x_i) \quad [\text{in feature space}]$$

Bias Term

$$b = (1/N_{sv}) \sum [y_s - \sum \alpha_i y_i K(x_i, x_s)] \quad [\text{average over support vectors}]$$

Prediction

$$f(x) = \text{sign}(\sum \alpha_i y_i K(x_i, x) + b)$$

Margin

$$\text{margin} = 2 / ||w||$$

11. Hyperparameter Tuning

C (Regularization)

Effect:

- $\uparrow C$: Lower bias, higher variance (harder margin, fit training data closely)
- $\downarrow C$: Higher bias, lower variance (softer margin, better generalization)

Typical values: 0.1, 1, 10, 100

γ (RBF Kernel)

Effect:

- $\uparrow \gamma$: Lower bias, higher variance (more complex boundary)

- $\downarrow \gamma$: Higher bias, lower variance (smoother boundary)

Formula: $\gamma = 1/(2\sigma^2)$

Typical values: 0.001, 0.01, 0.1, 1

Selection Strategy

Use **Grid Search** or **Random Search** with **Cross-Validation**

Common ranges:

```
param_grid = {  
    'C': [0.1, 1, 10, 100, 1000],  
    'gamma': [1, 0.1, 0.01, 0.001, 0.0001],  
    'kernel': ['rbf', 'poly', 'sigmoid']  
}
```

12. Computational Complexity

Training

- **Primal:** $O(nd)$ per iteration
- **Dual:** $O(n^2 \text{ to } n^3)$ depending on solver
 - n : number of samples
 - d : number of features

For large n : Use linear SVM or approximations

Prediction

$O(n_{sv} \cdot d)$

- n_{sv} : number of support vectors
- Typically $n_{sv} \ll n$

13. Advantages vs Disadvantages

☒ Advantages

1. **Effective in high dimensions** ($d > n$)
2. **Memory efficient** (only stores support vectors)
3. **Versatile** (different kernels)
4. **Robust** to outliers (when C is small)
5. **Global optimum** (convex optimization)
6. **Good generalization** with proper tuning

✗ Disadvantages

1. **Slow for large datasets** ($O(n^2)$ to $O(n^3)$)
 2. **No probability estimates** (requires Platt scaling)
 3. **Sensitive to feature scaling**
 4. **Kernel/parameter selection** can be difficult
 5. **Black box** with non-linear kernels (interpretability)
 6. **Not suitable for large n** (use linear SVM or logistic regression)
-

14. Feature Scaling

Critical: SVMs are sensitive to feature scales

Standardization (Recommended)

$$x' = (x - \mu) / \sigma$$

- Mean = 0, Std = 1
- Preserves outliers

Normalization

$$x' = (x - \min) / (\max - \min)$$

- Range [0, 1]
- Affected by outliers

Rule: Always scale features before training!

15. Implementation Tips

Scikit-learn Example

```
from sklearn.svm import SVC
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import GridSearchCV

# Scale features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Train SVM
svm = SVC(kernel='rbf', C=1.0, gamma='scale')
svm.fit(X_train_scaled, y_train)
```

```
# Hyperparameter tuning
param_grid = {'C': [0.1, 1, 10], 'gamma': [0.001, 0.01, 0.1]}
grid = GridSearchCV(SVC(), param_grid, cv=5)
grid.fit(X_train_scaled, y_train)

# Access support vectors
n_support = svm.n_support_ # number per class
support_vectors = svm.support_vectors_
support_indices = svm.support_

# Prediction
y_pred = svm.predict(X_test_scaled)
```

16. When to Use SVM

☒ Use SVM when:

- High-dimensional data (text classification, genomics)
- Clear margin of separation exists
- More features than samples
- Need robust model against outliers
- Non-linear decision boundaries (with kernels)

☒ Avoid SVM when:

- Very large datasets ($n > 100,000$)
- Need probability estimates (use logistic regression)
- Need interpretability (use decision trees)
- Data has significant noise/overlapping classes
- Computational resources are limited

17. Common Applications

1. **Text Classification** (spam detection, sentiment analysis)
2. **Image Recognition** (face detection, handwriting recognition)
3. **Bioinformatics** (protein classification, gene expression)
4. **Medical Diagnosis** (cancer classification)
5. **Time Series Prediction** (with SVR)
6. **Financial Forecasting**
7. **Remote Sensing** (land cover classification)

18. Important Notes

Probability Estimates

SVM doesn't naturally output probabilities. Use **Platt Scaling**:

$$P(y=1|x) = 1 / (1 + \exp(A \cdot f(x) + B))$$

- Requires additional cross-validation
- Enable with `probability=True` in sklearn

Class Imbalance

- Use `class_weight='balanced'` parameter
- Or manually set: `class_weight={0: w0, 1: w1}`
- Adjusts C: $C_i = C \times w_i$

Kernel Selection Guidelines

1. Start with **linear** kernel (baseline)
2. Try **RBF** if linear doesn't work well
3. Use **polynomial** for specific problem knowledge
4. Avoid **sigmoid** (rarely useful)

Quick Reference Card

Concept	Formula
Hyperplane	$w \cdot x + b = 0$
Primal Problem	$\min (1/2) \ w\ ^2 + C \cdot \sum \xi_i$
Dual Problem	$\max \sum \alpha_i - (1/2) \sum \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j)$
Decision Function	$f(x) = \text{sign}(\sum \alpha_i y_i K(x_i, x) + b)$
Margin	$2 / \ w\ $
RBF Kernel	$K(x, x') = \exp(-\gamma \ x - x'\ ^2)$
Support Vectors	Points where $0 < \alpha_i \leq C$

Remember:

- Always scale features
- Tune C and γ carefully
- Use cross-validation
- Consider linear SVM first for large datasets