

Minor in AI

SVM: Classifying the Unclassifiable

1 The Magic of Separation: Spam or Ham?

Imagine your phone buzzing with a new message: "Congrats! You won \$1,000,000. Click here to claim!" Is this legitimate or spam? This is exactly what Support Vector Machines (SVM) can solve! SVMs are powerful algorithms that find optimal boundaries between different categories of data.

Consider email classification:

- **Problem:** Separate genuine emails (ham) from spam
- **Challenge:** Messages don't come with clear labels
- **SVM Solution:** Finds the best dividing line using key examples

But what about more complex patterns? Like distinguishing between galaxy spirals in astronomy? Traditional straight-line classifiers fail here. This is where SVM's true power emerges - it can find curved boundaries in complex data landscapes.

2 The Geometry of Separation

2.1 Core Idea: Maximum Margin

SVM creates the widest possible "safety buffer" between classes:

$$\text{Margin Width} = \frac{2}{\|\mathbf{w}\|}$$

where \mathbf{w} is the weight vector. Maximizing this margin improves generalization.

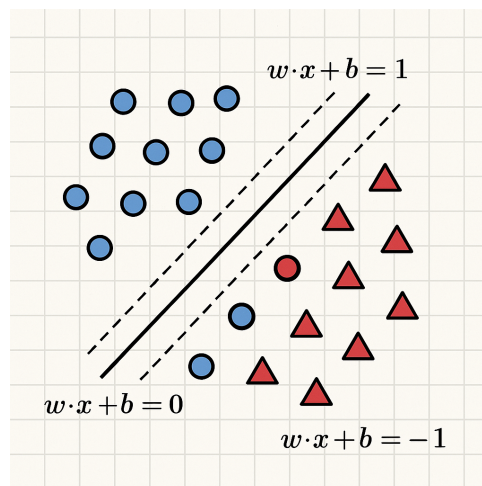


Figure 1: SVM creates the widest possible separation between classes

2.2 Support Vectors: The Heroes

These are the critical data points that "support" the decision boundary:

- Lie closest to the separating hyperplane
- Determine the margin's position and orientation
- If removed, the decision boundary changes

2.3 The Kernel Trick: Seeing in Higher Dimensions

When data isn't linearly separable (like our spiral galaxy data), SVM uses kernel functions to project data into higher dimensions:

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

Common kernels:

Kernel	Use Case
Linear	Simple separations
Polynomial	Curved boundaries
RBF (Radial Basis Function)	Complex patterns like spirals

3 Galaxy Classification: A Cosmic Case Study

3.1 Creating Cosmic Data

We generate synthetic spiral data mimicking galaxy patterns:

Listing 1: Generating Spiral Data

```

1 import numpy as np
2
3 def generate_spiral(n_points, noise=0.5):
4     n = np.sqrt(np.random.rand(n_points)) * 780 * (2*np.pi)/360
5     dx = -np.cos(n)*n + np.random.rand(n_points) * noise
6     dy = np.sin(n)*n + np.random.rand(n_points) * noise
7     # ... (full code from lecture_code.py)
8     return X, y
9
10 X, y = generate_spiral(200, noise=0.5)

```

This creates two interleaved spirals (Class 0 = blue, Class 1 = red) with controlled noise.

3.2 Why Linear Models Fail

Traditional logistic regression creates a straight boundary that fails spectacularly:

```

1 from sklearn.linear_model import LogisticRegression
2 clf = LogisticRegression()
3 clf.fit(X, y)
4 # Visualize poor performance

```

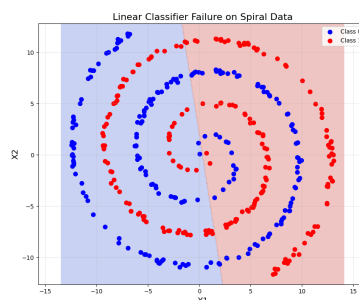


Figure 2: Linear classifier fails on spiral data

3.3 RBF Kernel to the Rescue

We project 2D data into 3D using Radial Basis Function:

```
1 def rbf_lift(X, gamma=0.1, center=np.array([0,0])):
2     diff = X - center
3     squared_dist = np.sum(diff**2, axis=1)
4     z = np.exp(-gamma * squared_dist)
5     return z
```

The transformation creates a 3D landscape where classes become separable:

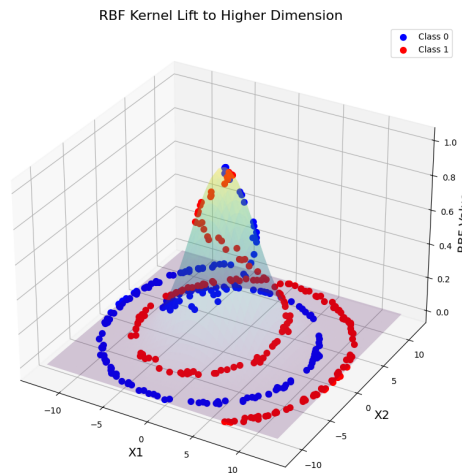


Figure 3: Spiral data lifted to 3D via RBF kernel

3.4 SVM in Action

Implement SVM with RBF kernel and visualize results:

```
1 from sklearn.svm import SVC
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5 X_scaled = scaler.fit_transform(X)
6 clf = SVC(kernel='rbf', gamma='auto')
7 clf.fit(X_scaled, y)
8
9 # Visualize decision boundaries and support vectors
```

Key parameters:

- **gamma**: Controls influence range of single points
- **C**: Regularization parameter (balance accuracy vs. simplicity)

4 Real-World Applications

4.1 SMS Spam Detection

Practical implementation with TF-IDF vectorization:

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.svm import SVC
3
4 # Load SMS dataset
5 vectorizer = TfidfVectorizer(stop_words='english')
6 X_train_tfidf = vectorizer.fit_transform(X_train)
7
8 clf = SVC(kernel='linear')
9 clf.fit(X_train_tfidf, y_train)
10 # Achieves >98% accuracy!
```

4.2 Diabetes Prediction

Using the Pima Indians Diabetes Database:

```
1 from sklearn.svm import SVC
2
3 # Load diabetes data
4 clf = SVC(kernel='rbf')
5 clf.fit(X_train, y_train)
6
7 # Compare kernel performance:
8 # Linear: 70% accuracy
9 # RBF: 77% accuracy
```

Pro Tip: Kernel Selection

Start simple:

1. Try linear kernel first
2. If accuracy low, switch to RBF
3. Tune gamma and C parameters
4. Visualize decision boundaries

5 Key Takeaways: The SVM Advantage

1. **Margin Maximization:** Creates robust decision boundaries resistant to noise
2. **Kernel Magic:** Handles non-linear data through dimension lifting
3. **Support Vectors:** Only critical points affect the model (memory efficient)
4. **Versatility:** Works with images, text, biological data, astronomy
5. **Parameter Tuning:** Gamma and C dramatically affect performance

Remember

"Support Vector Machines find the widest possible road between classes, using the most critical data points as guardrails."

6 SVM Exploration Kit

Continue experimenting with:

- **Datasets:** - SMS Spam Collection: <https://archive.ics.uci.edu/dataset/228/sms+spam+collection> - Pima Indians Diabetes: <https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>
- **Interactive Tools:** - SVM Playground: <https://www.cs.cmu.edu/~svmplay/> - Kernel Visualization: https://scikit-learn.org/stable/auto_examples/svm/plot_iris_svc.html