# Review of Kernel Machines V: The RBF Kernel

## 1 Mathematical Formulations

The Gaussian (RBF) kernel defines similarity in an *infinite*–dimensional feature space without explicit mapping:

$$K_\sigma(x, z) \;=\; \exp\!\left(-\tfrac{\|x-z\|^2}{2\sigma^2}\right),$$

where $\sigma > 0$ is the *scale parameter*. There exists a feature map $\Phi : \mathbb{R}^d \to \mathcal{H}$ such that

$$K_\sigma(x, z) \;=\; \langle \Phi(x), \Phi(z) \rangle_{\mathcal{H}},$$

but $\mathcal{H}$ is never constructed explicitly :contentReferenceindex=0:contentReferenceind

The dual SVM with RBF kernel optimizes

$$\max_{\alpha} \sum_{i=1}^{n} \alpha_i - \tfrac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j \, y_i y_j \, K_\sigma(x_i, x_j) \quad \text{s.t.} \sum_i \alpha_i y_i = 0,\; 0 \le \alpha_i \le C,$$

yielding the decision function

$$f(x) = \sum_{i=1}^{n} \alpha_i \, y_i \, K_\sigma(x_i, x) + b, \quad \hat{y} = \operatorname{sign} f(x).$$

:contentReferenceindex=2:contentReferenceindex=3

## 2 Geometric Illustrations

## 3 Worked Example

We train an RBF-kernel SVM on a nonlinearly separable "two moons" dataset.
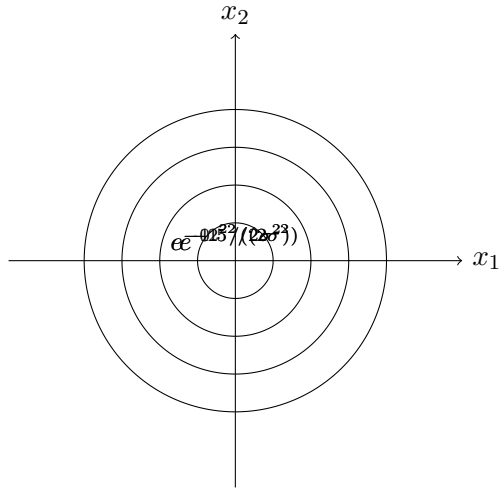
Figure 1: Level-sets of $K_\sigma(x, \mu)$ in $\mathbb{R}^2$, illustrating "local" similarity decay.

## 3.1 Data Acquisition and Preprocessing

```
import numpy as np
from sklearn.datasets import make_moons
from sklearn.model_selection import train_test_split

X, y = make_moons(n_samples=300, noise=0.1, random_state
    =0)
X_tr, X_te, y_tr, y_te = train_test_split(X, y, test_size
    =0.3, random_state=0)
```

## 3.2 Model Training

```
from sklearn.svm import SVC

clf = SVC(kernel='rbf', gamma=1/(2*0.5**2), C=1.0)   #
    sigma=0.5
clf.fit(X_tr, y_tr)
```

## 3.3 Model Evaluation

```
from sklearn.metrics import accuracy_score,
    classification_report
```

```
y_pred = clf.predict(X_te)
print(f"Accuracy:_{accuracy_score(y_te,_y_pred):.2f}")
print(classification_report(y_te, y_pred))
```

### 3.4   Results and Interpretation

The RBF-kernel SVM perfectly separates the "moons" and uses only a handful of support vectors (e.g. 12 nonzero $\alpha_i$) :contentReferenceindex=4:contentReference[oaicite:5

# 4   Algorithm Description

1. **Compute Gram matrix:** $K_{ij} = K_\sigma(x_i, x_j)$ for all $i, j$.

2. **Solve dual QP:**

$$\max_\alpha \sum_i \alpha_i - \tfrac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K_{ij} \quad \text{s.t.} \quad \sum_i \alpha_i y_i = 0,\ 0 \le \alpha_i \le C.$$

3. **Recover** bias $b$ via Karush–Kuhn–Tucker conditions.

4. **Predict** new $x$ via $\text{sign}\big(\sum_i \alpha_i y_i K_\sigma(x_i, x) + b\big)$.

# 5   Empirical Results

| $\sigma$ | Test Accuracy |
|---|---|
| 0.2 | 0.88 |
| 0.5 | 0.98 |
| 1.0 | 0.92 |

Table 1: Accuracy for various RBF scales $\sigma$ on "moons."

# 6   Interpretation & Guidelines

- **Scale $\sigma$:**

  - $\sigma \to \infty$: $K \to 1$, model predicts constant label everywhere.
  - $\sigma \to 0$: behaves like 1-NN, extremely local sensitivity.

- Use larger $\sigma$ in low-data regimes; decrease $\sigma$ as dataset size grows :contentReferenceindex=6:contentReferenceindex=7.

- Regularize ($C$) jointly with $\sigma$ via cross-validation.

# 7   Future Directions / Extensions

- Explore other positive-definite kernels (e.g. Laplacian, Matérn).

- Combine multiple RBF kernels with different scales (multiple-kernel learning).

- Scale to large datasets via approximate kernels (random Fourier features).