

Support Vector Machines III



COMPCSI 361

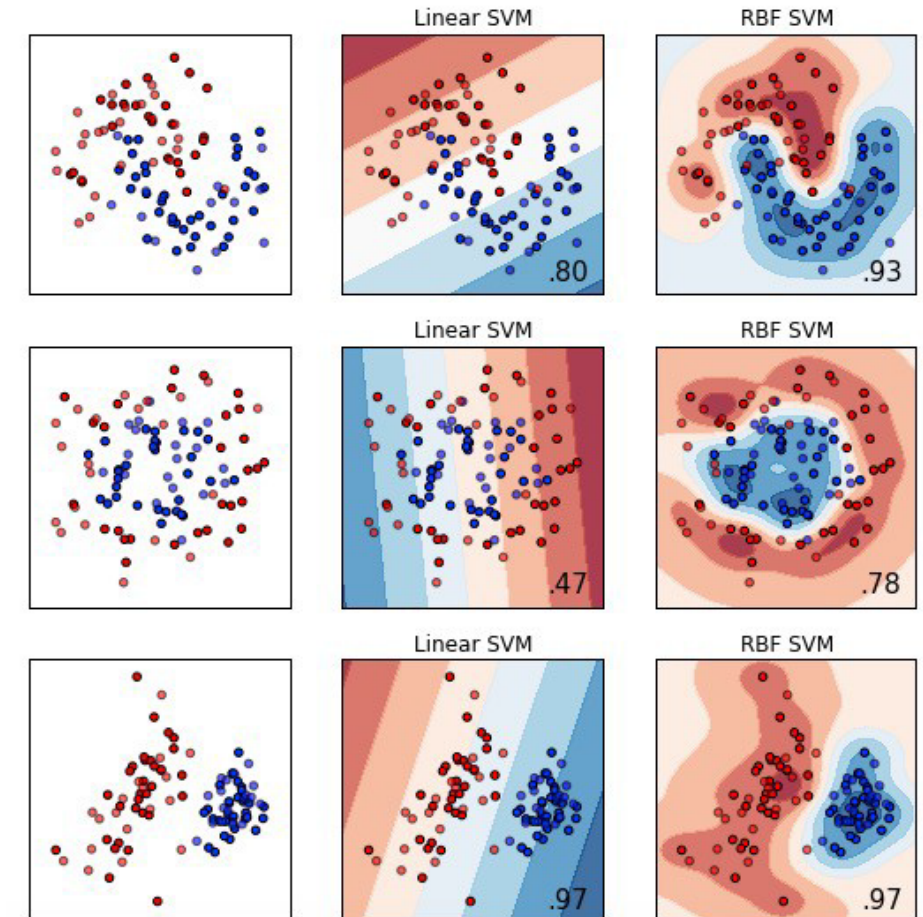
Instructor: Thomas Lacombe

Based on slides from Meng-Fen Chiang

WEEK 10

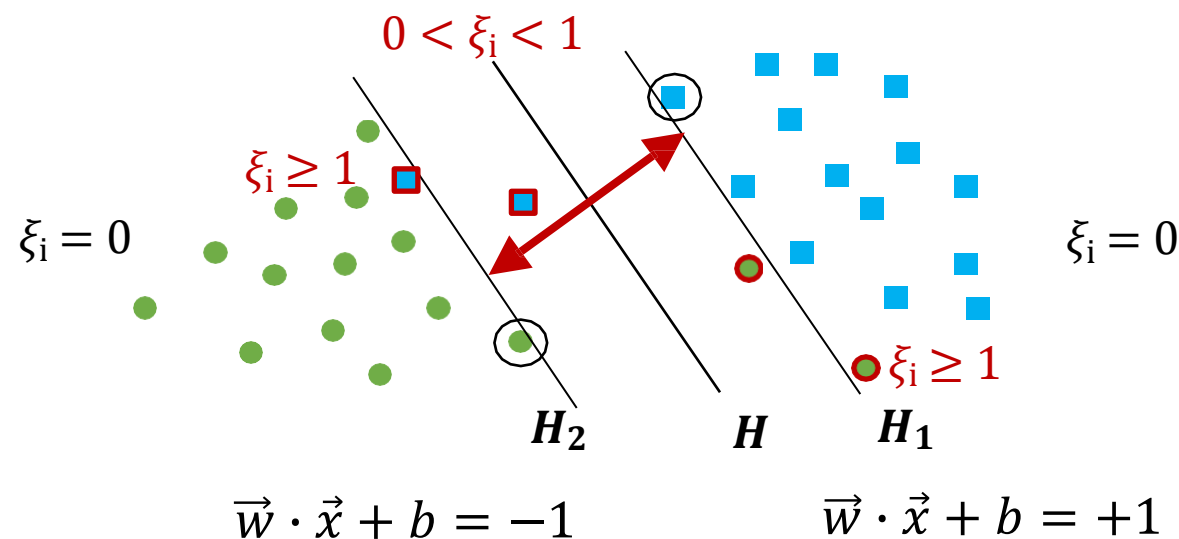
OUTLINE

- Data Characteristics
 - Linearly Separable Data
 - Non-Linearly separable Data
- SVM
 - Linearly Separable Data: Hard-margin SVMs (9.1)
 - Non-Linearly Separable Data: Soft-margin SVMs (9.2)
 - Non-Linearly Separable Data: Kernelized SVMs (9.3)
- Summary



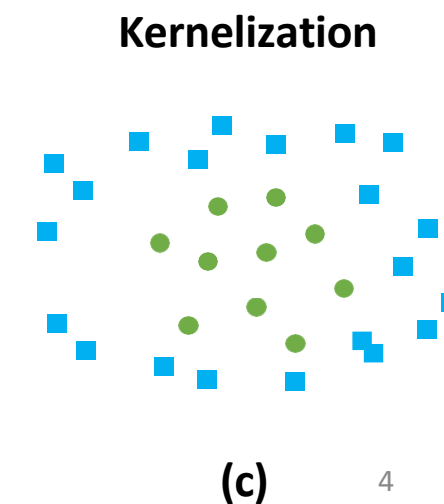
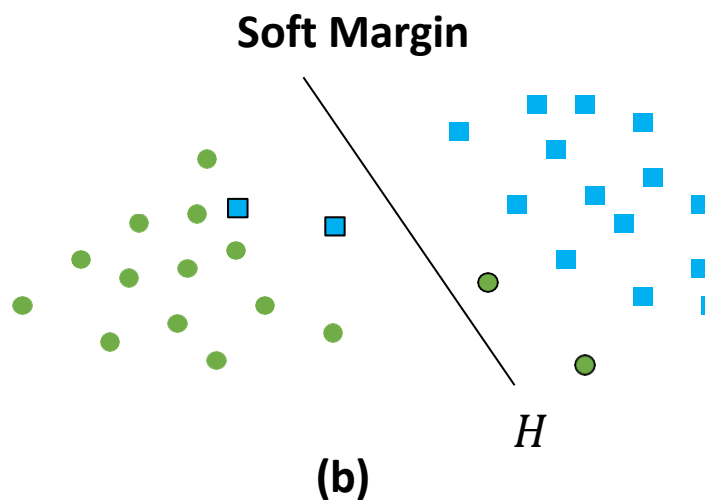
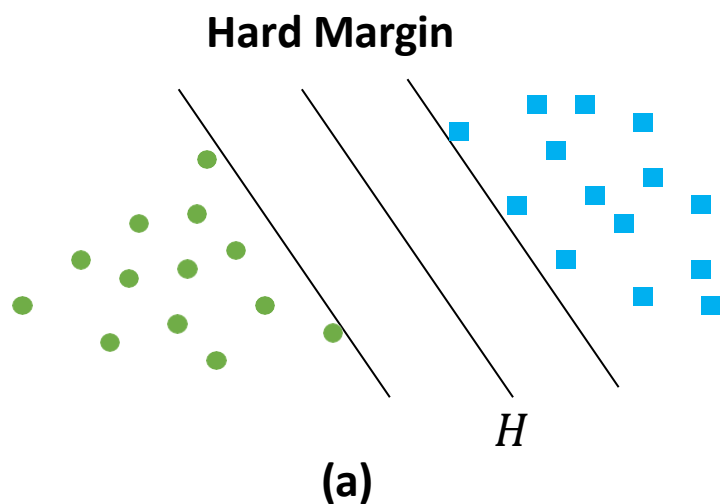
RECAP: Soft-margin Maximization

- Given a set of training data $S = ((x_1, y_1), \dots, (x_n, y_n))$, $y_i \in \{+1, -1\}$
- Goal: The **soft-margin** SVM algorithm aims to find a linear classifier that
 - Maximizes (γ) the margin on S and
 - Minimize the misclassification error $C \sum_{i=1}^n \xi_i$

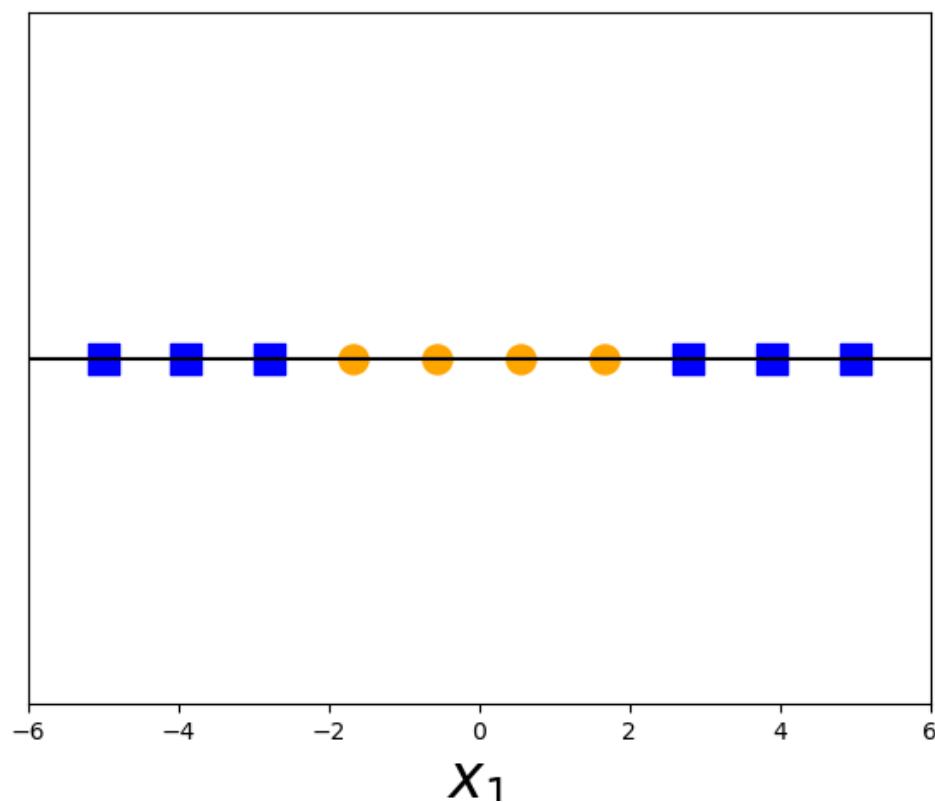


RECAP: Non-Linearly Separable Data

- Input data is not linearly separable in a two-dimensional space
 - Soft-margin SVMs: a linear separating hyperplane that allows misclassifications
 - Kernelized SVMs: non-linear classifier that is a linear separating hyperplane(s) in higher feature space

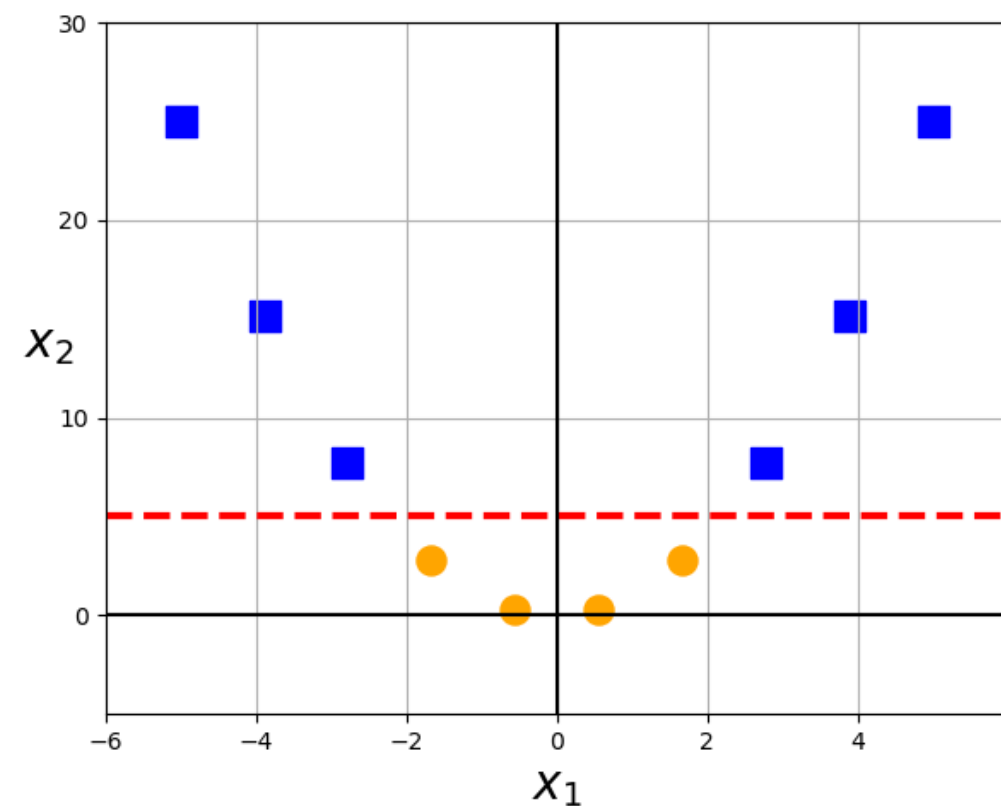


Intuitive example 1



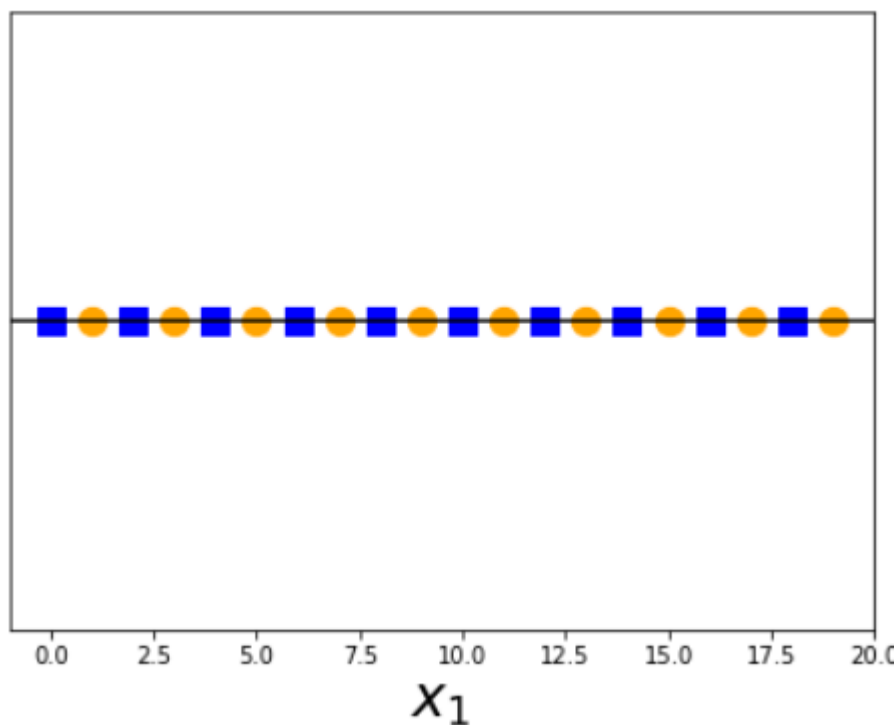
Non-linearly separable data

$$\Phi(x) = x^2$$



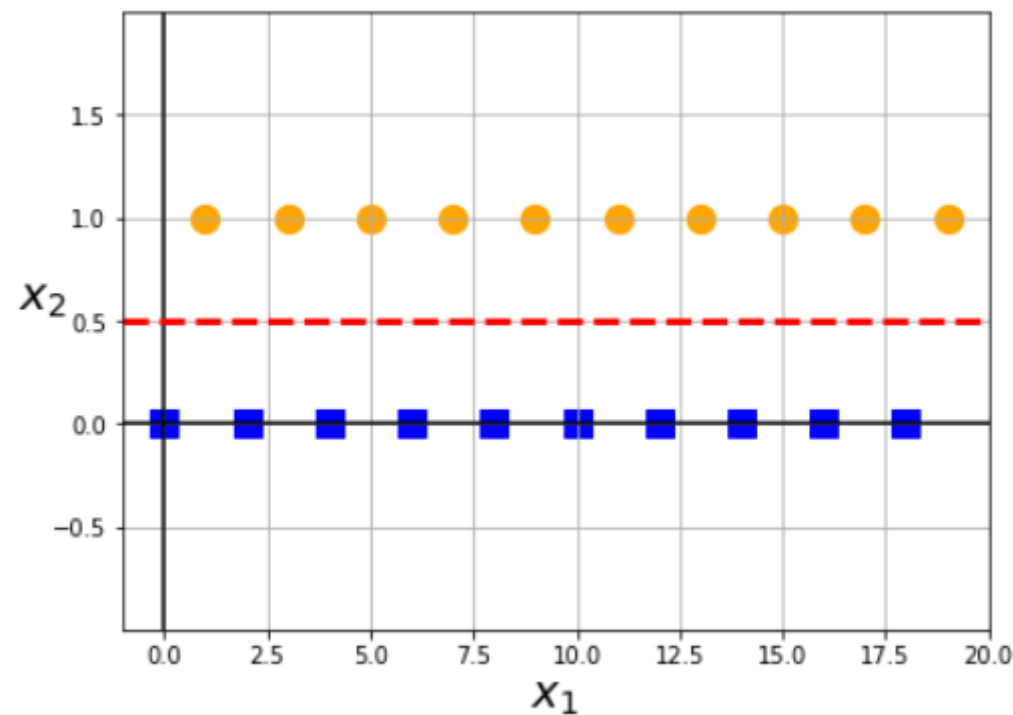
Linearly separable data!

Intuitive example 2



Non-linearly separable data

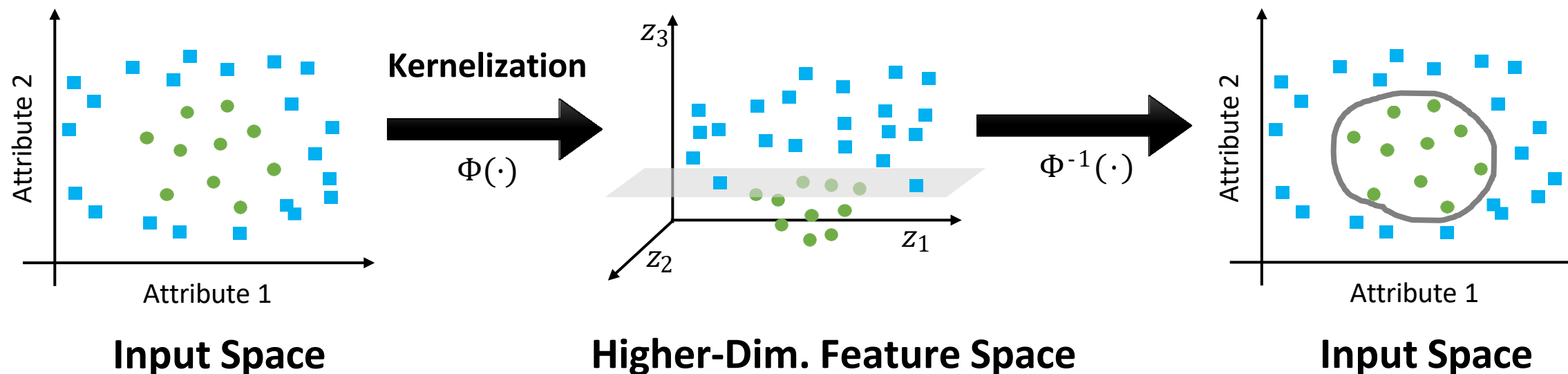
$$\Phi(x) = x \bmod(2)$$



Linearly separable data!

Notion of Feature Transformation

- Non-linearly separable in input data space
- Linearly separable in a ***higher dimensional*** feature space



Problem Definition: Kernelized SVMs

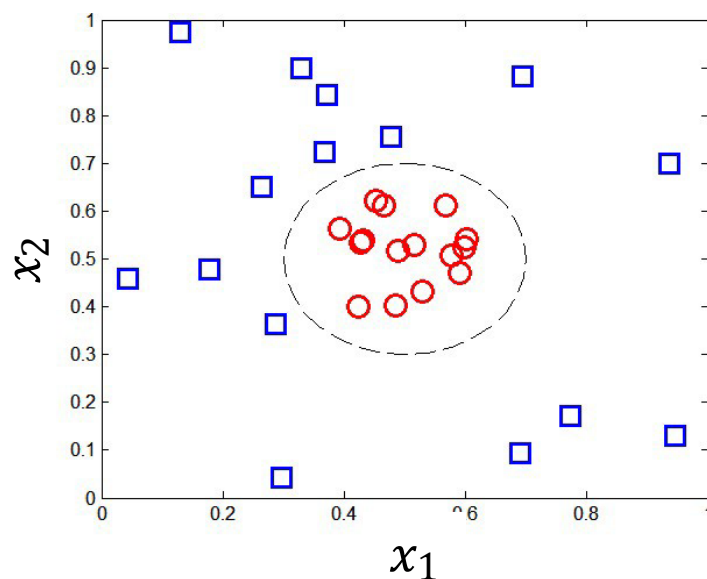
- Given a set of training data $S = ((x_1, y_1), \dots, (x_n, y_n))$, $y_i \in \{+1, -1\}$
- Goal: The **kernelized** SVM algorithm aims to find a **non-linear classifier**, which is a separating hyperplane(s) in a higher dimensional space
- Approach
 1. Transform the original input data into a higher dimensional space
 2. Search for a linear separating hyperplane in the new space

Step1: Space Transformation

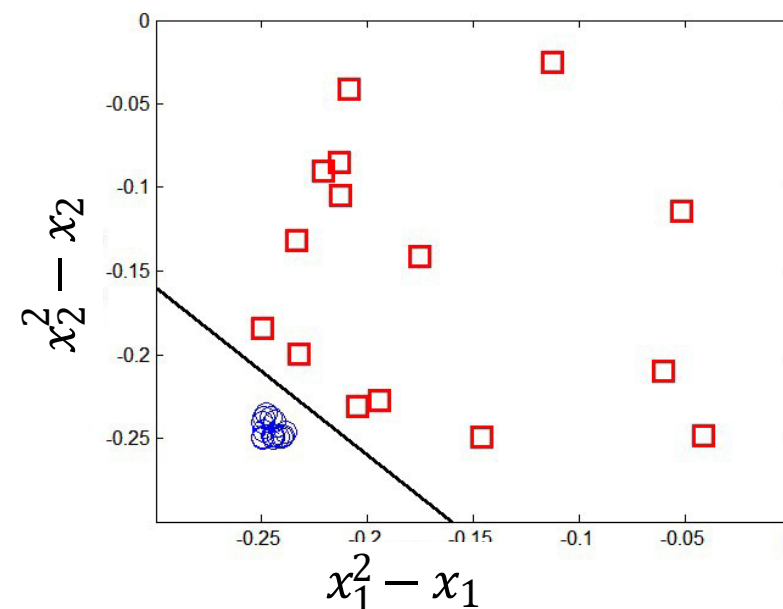
- $\Phi(\cdot)$: Transform the original input data into a higher dimensional space

$$\Phi:(x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2)$$

- Hyperplane: $\vec{w} \cdot \Phi(\vec{x}) + b = 0 \rightarrow w_4x_1^2 + w_3x_2^2 + w_2\sqrt{2}x_1 + w_1\sqrt{2}x_2 + b = 0$



$\Phi(\cdot)$



Step2: Search for a Linear Separating Hyperplane

- Dual Optimization Problem:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \overrightarrow{x_i} \cdot \overrightarrow{x_j} \quad \text{s.t. } \alpha_i \geq 0 \text{ and } \sum_{i=1}^n \alpha_i y_i = 0, \quad i = 1, 2, \dots, n$$

- Kernelized Optimization Problem:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \Phi(\overrightarrow{x_i}) \cdot \Phi(\overrightarrow{x_j}) \quad \text{s.t. } \alpha_i \geq 0 \text{ and } \sum_{i=1}^n \alpha_i y_i = 0, \quad i = 1, 2, \dots, n$$

- Same set of equations as dual problem optimization except that involve $\Phi(\vec{x})$ in feature space, instead of \vec{x} in input space

Step2: Search for a Linear Separating Hyperplane

- Training: Solve the Kernelized Optimization Problem by Quadratic Programming (QP)

$$\vec{w} = \sum_{x_i \in SV} \alpha_i y_i \Phi(\vec{x}_i) \qquad \vec{b} = \frac{1}{|SV|} \sum_{x_i \in SV} y_i - (\vec{w} \cdot \Phi(\vec{x}_i))$$

- Testing: Determine the class label for a test point \vec{z} by using the learned kernelized SVM (\vec{w} and \vec{b}) with support vectors (SV)

$$\begin{aligned} f(\vec{z}) &= \text{sign}(\vec{w} \cdot \Phi(\vec{z}) + b) \\ &= \text{sign}(\sum_{x_i \in SV} \alpha_i y_i (\Phi(\vec{x}_i) \cdot \Phi(\vec{z})) + b) \end{aligned}$$

Kernel Trick

- Kernelized Optimization Problem:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j)$$

- In practice, applying the transformation Φ to the original data is often impractical (high number of features and Φ can involve polynomial combinations).
- Kernel Trick:** Instead of transforming the input data, and then applying the inner product between pairs of transformed data points, a **Kernel function** is applied:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j)$$

Kernel Functions

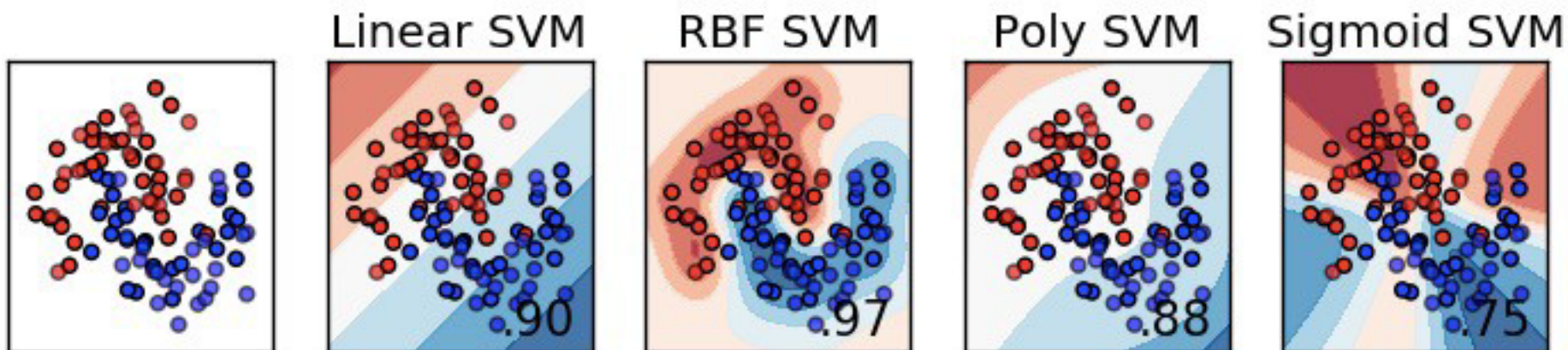
- **Kernel function:** A kernel function K takes vectors in the original space as inputs, and returns the inner product of the vectors in the transformed space:

$$K(\vec{x}_i, \vec{x}_j) = \Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j)$$

- A kernel function calculates the similarity between pairs of data points.
- Applying the function is equivalent to calculating the inner product in the transformed space, but Φ is never directly calculated (no need to know Φ explicitly).
- A valid kernel function is a function that can be used to compute the inner product between two feature vectors in a high-dimensional space without explicitly mapping them.

Kernel Functions

- There exist a lot of possible kernel functions:



Kernel Function: Polynomial Kernel

- Polynomial Kernel of Degree d :

$$k(\vec{x}, \vec{z}) = (\vec{x} \cdot \vec{z})^d$$

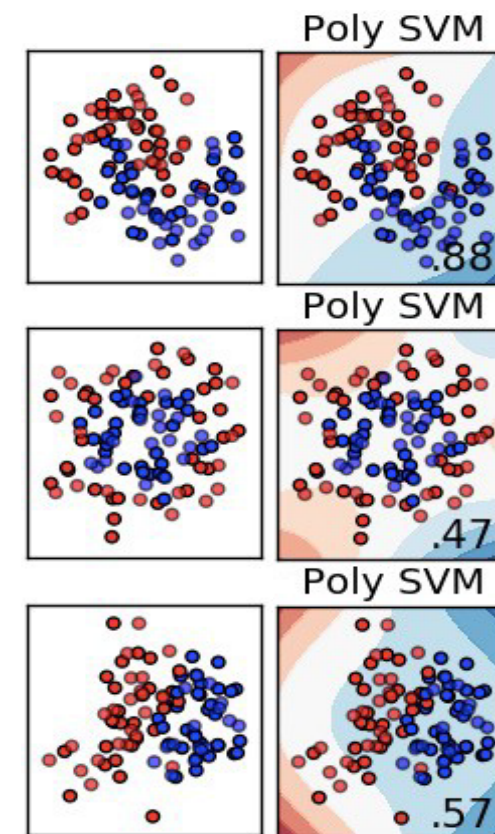
$$d=1: \Phi(\vec{x})\Phi(\vec{z}) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = x_1z_1 + x_2z_2 = (\vec{x} \cdot \vec{z})$$

Second degree polynomial mapping:

$$\Phi(\vec{x}) = \Phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$

$$\begin{aligned} d=2: \Phi(\vec{x})\Phi(\vec{z}) &= \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix} = x_1^2z_1^2 + 2x_1x_2z_1z_2 + x_2^2z_2^2 \\ &= (x_1z_1 + x_2z_2)^2 \\ &= (\vec{x} \cdot \vec{z})^2 \end{aligned}$$

Example Decision Boundaries

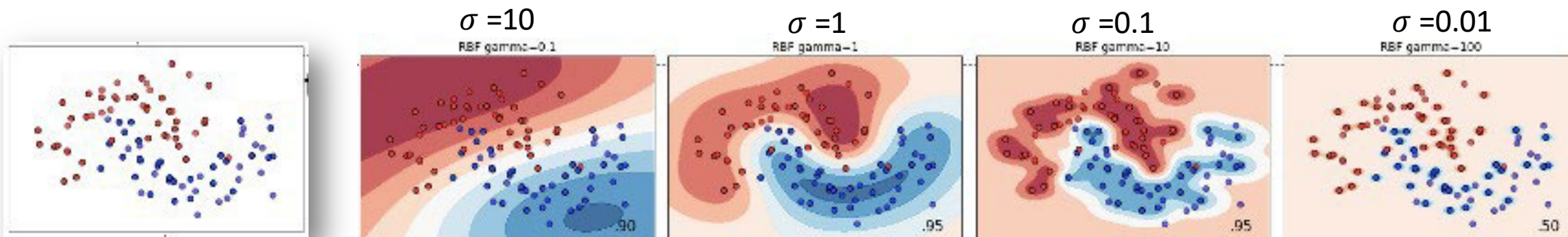
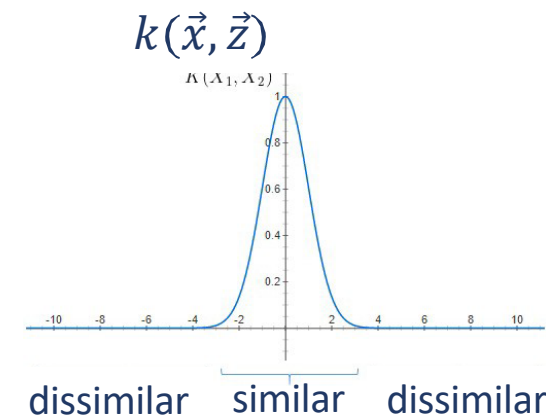


Kernel Function: RBF Kernel

- Radial Basis Function (RBF):

$$k(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|}{2\sigma^2}\right), \quad k(\cdot) \in [0,1]$$

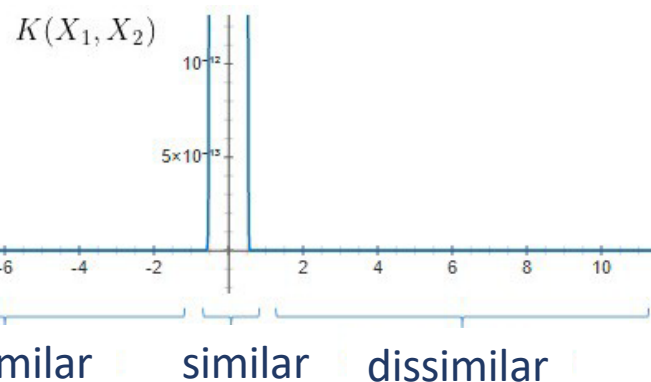
- sigma (σ) defines kernel width (radius) of a single training data
 - i.e., high σ values mean 'wider kernel' \rightarrow smoother decision boundary
 - i.e., small σ values mean 'narrow kernel' \rightarrow like the nearest neighbors



Example: Kernel Width (σ)

$$\sigma=0.1$$

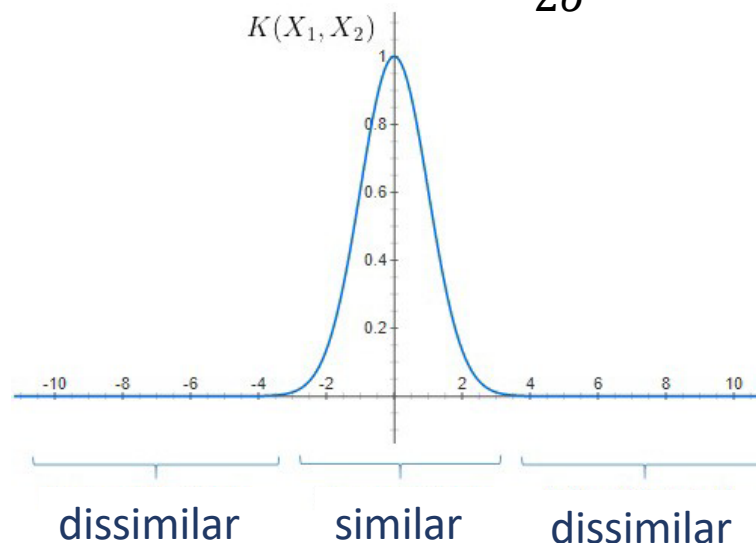
$$k(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|}{2\sigma^2}\right)$$



(a) Narrower Kernel Width

$$\sigma=1$$

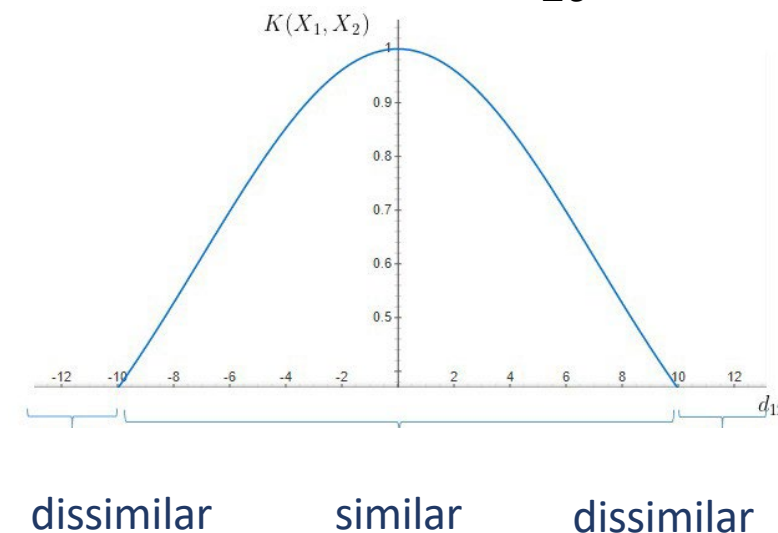
$$k(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|}{2\sigma^2}\right)$$



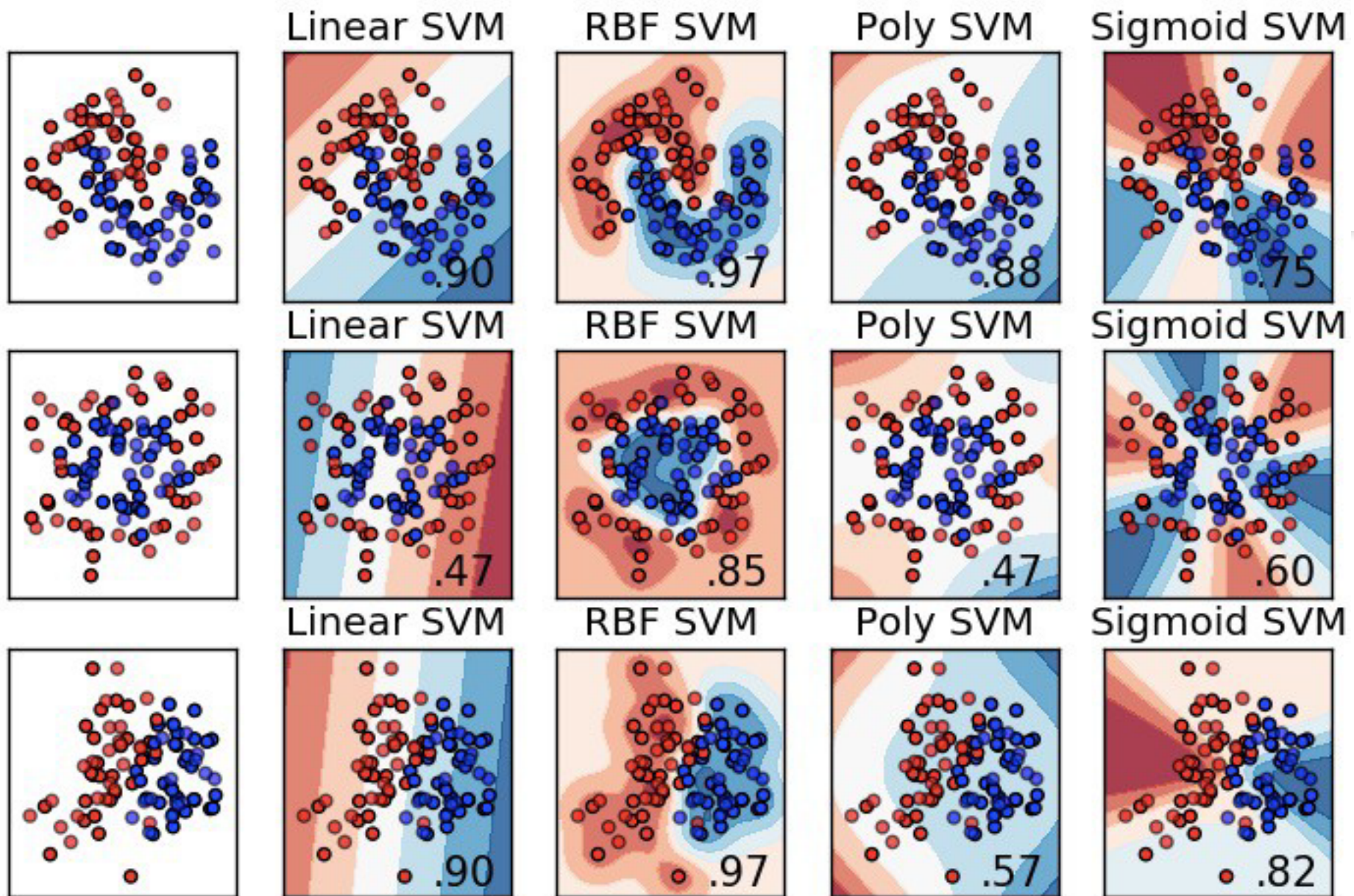
(b) Kernel Width

$$\sigma=10$$

$$k(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|}{2\sigma^2}\right)$$

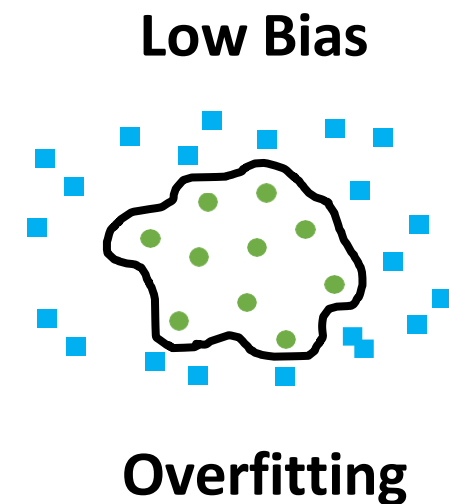
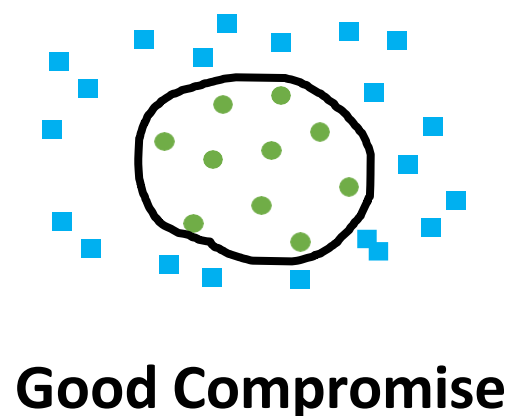
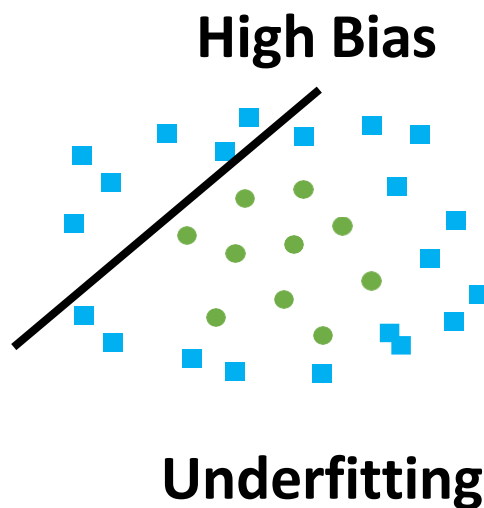


(c) Wider Kernel Width



Bias and Variance

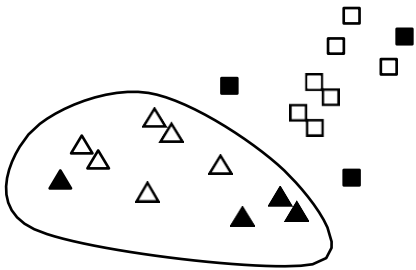
- A non-linear separation that trade-off between the bias and variance



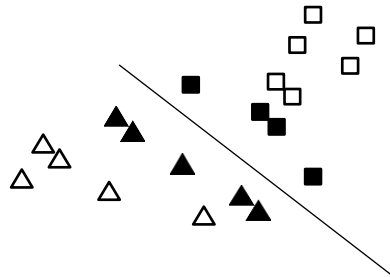
Quiz: Kernelized SVM with Soft/Hard Margin

- Which decision boundaries refer to kernelize SVMs? (Note: support vectors are represented by solid square/triangle)

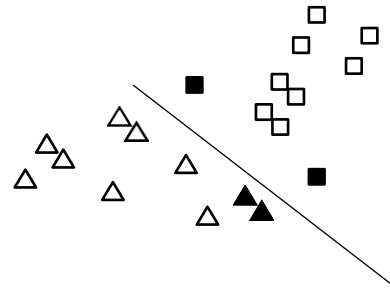
Kernelized SVM



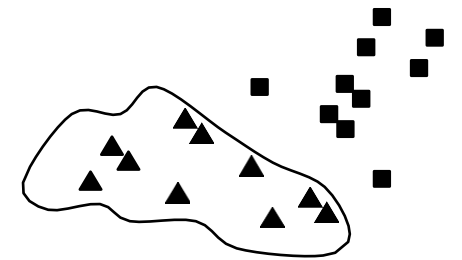
Soft-margin SVM



Hard-margin SVM



Kernelized SVM



Advantages v.s. Disadvantages

Advantages

- Prediction accuracy is generally high
- As compared to Bayesian methods generally
- Robust when training examples contain errors
- Fast evaluation of the learned target function
- Bayesian networks are normally slow

Disadvantages

- Long training time
- Difficult to understand the learned function (weights)
- Bayesian networks can be used easily for pattern discovery
- Not easy to incorporate domain knowledge
- Easy in the form of priors on the data or distributions

Why Kernelized SVMs Work?

- It uses a **nonlinear mapping** to transform the original training data into a higher dimension
- With the new dimension, it searches for the linear optimal separating **hyperplane** (i.e., “decision boundary”)
- With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane
- The hyperplane is discovered based on **support vectors** (“essential” training tuples) and **margins** (defined by the support vectors)

Jupyter Notebook

Kernelized SVMs Coding Example

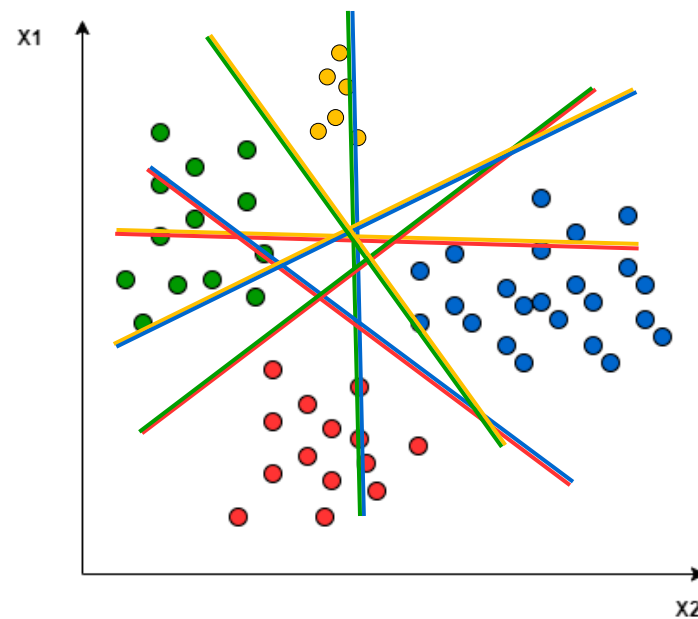
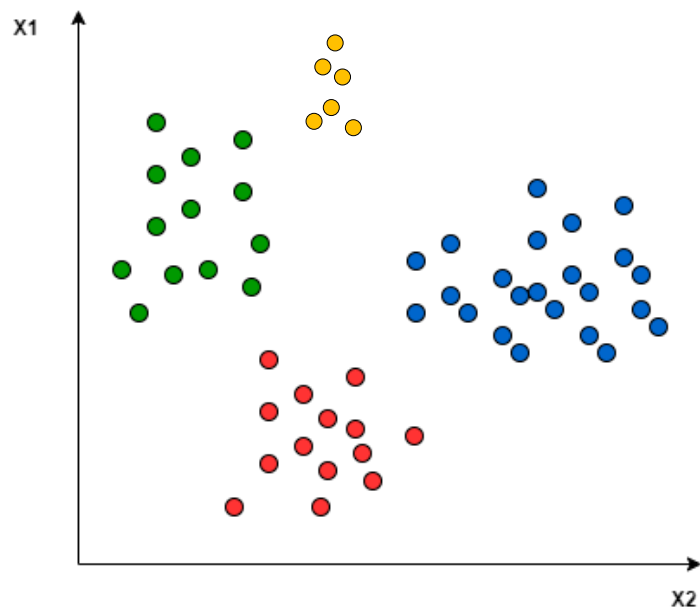
What about if we have more than 2 classes?

- Original SVMs can only perform binary classification (2 classes).
- SVM can be extended to multiclass problems (more than 2 classes).
- 2 approaches:
 1. One-vs-One (OVO):

Train one SVMs for each binary problem, i.e., each 2 classes, ignoring the other $\frac{m(m-1)}{2}$ SVMs in total).
 2. One-vs-Rest (OVR):

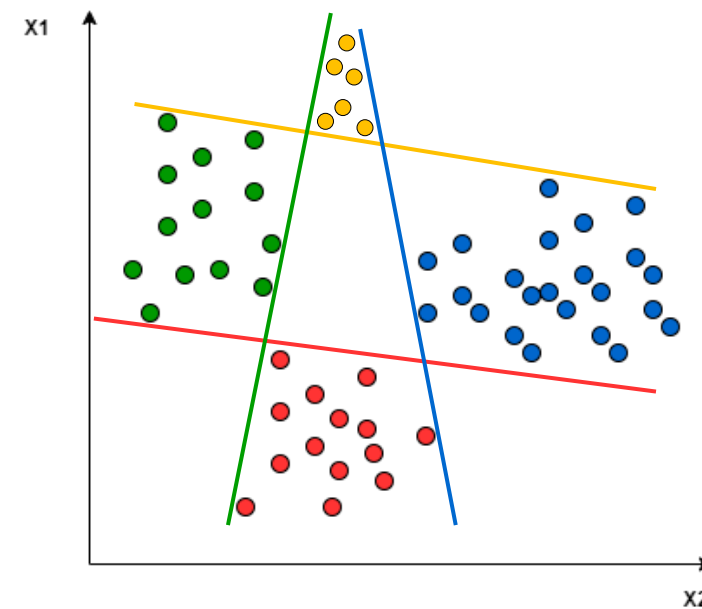
Train m SVMs. Each SVM learns to separate 1 class from all the other ones.

OVO vs OVR



OVO

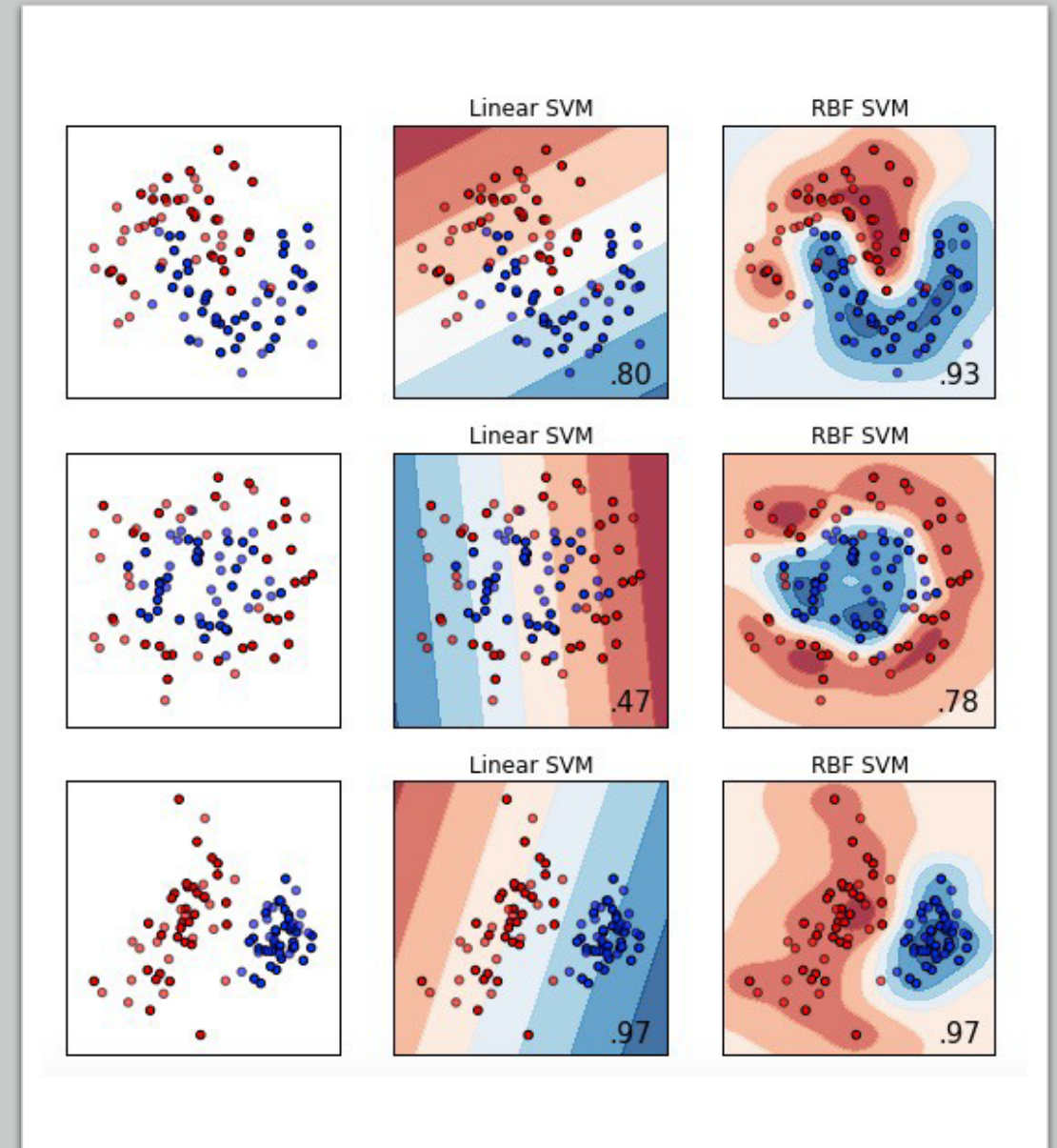
$$\frac{m(m-1)}{2} = 6 \text{ SVMs}$$



OVR
 4 SVMs

SUMMARY

- Kernelized SVMs
 - Non-linearly Separable Data
 - Dual Problem Optimization
 - Kernel Problem Optimization
 - Testing Stage
 - Kernel Functions
- SVMs: Advantages and Disadvantages



Resources

- SVM Website: <http://www.kernel-machines.org/>
- Representative Implementation
 - **LIBSVM**: an efficient implementation of SVM, multi-class classifications, nu-SVM, one-class SVM, including also various interfaces with java, python, etc.
 - **SVM-light**: simpler but performance is not better than LIBSVM, support only binary classification and only in C
 - **SVM-torch**: another recent implementation also written in C
 - **Scikit-Learn**: a set of supervised learning methods used for classification, regression and outliers detection. [\[link\]](#)

Resources (Contd.)

- Book Chapters: Christopher Bishop, “Pattern Recognition and Machine Learning” ([PDF](#))
 - Sec 7.1.1-7.1.3
 - Sec 4.1.1, 4.1.2
 - Sec 6.1, 6.2
 - Appendix E
- Literatures
 - C.J.C. Burges, Chris J.C. Burges "A Tutorial on Support Vector Machines for Pattern Recognition." Data Mining and Knowledge Discovery, 1998 (PDF)