



INT104 ARTIFICIAL INTELLIGENCE

Support Vector Machine Review

Fang Kang

Fang.kang@xjtlu.edu.cn



Xi'an Jiaotong-Liverpool University

西交利物浦大學

- Linear SVM
 - ◆ Hard-margin Classification
 - ◆ Soft-margin Classification

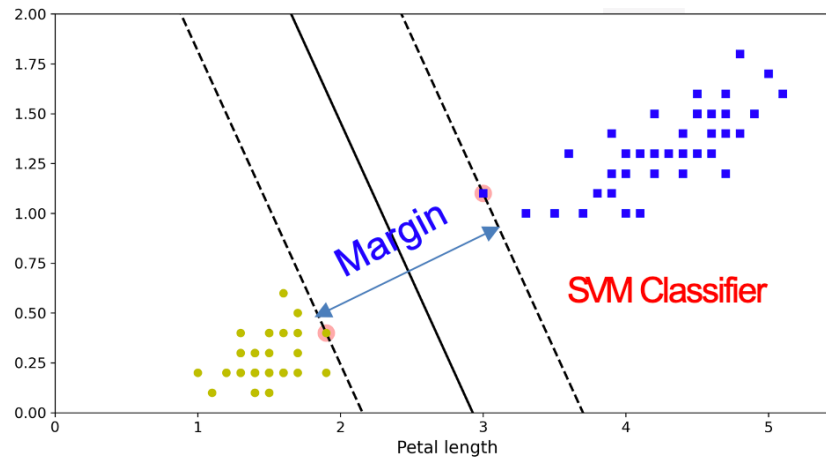
- Non-Linear SVM
 - ◆ Nonlinear SVM Classification
 - ◆ Kernel method
 - ◆ Polynomial Features

- SVM Regression



Hard Margin Classification

- All instances being off the street and on the right side is named “*hard margin classification*”



$$\text{Minimize: } \frac{1}{2} ||\mathbf{w}'||^2$$

$$\text{Subject to: } y_i(\mathbf{w}'^T \mathbf{x}_i + b) \geq 1 \quad (i = 1 \sim N)$$

- The main limitation of hard margin classification is
 - The data must be linearly separable
 - Sensitive to outliers

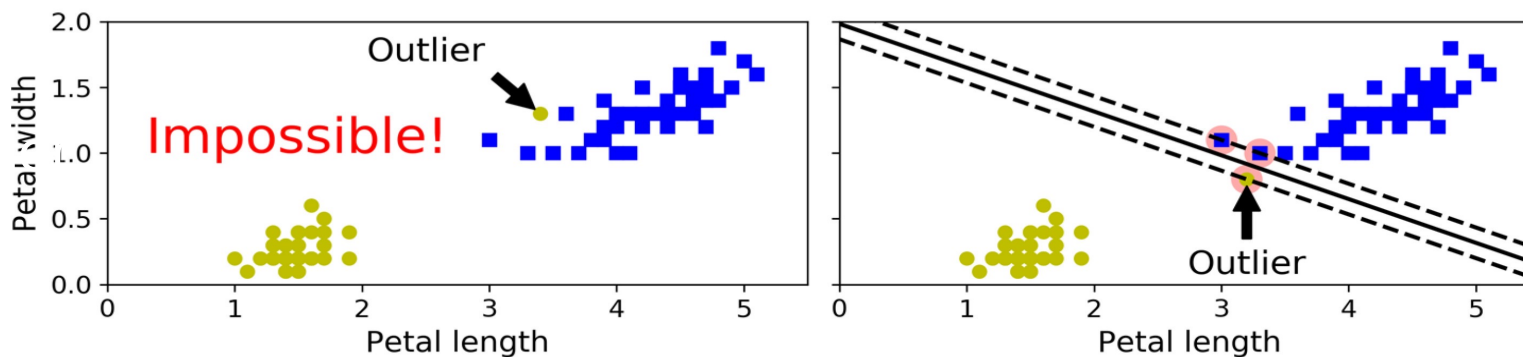


Figure 5-3. Hard margin sensitivity to outliers



Soft Margin Classification

- Allow margin violations
- The algorithm balances
 - The width of street
 - The amount of margin violations
- A hyper-parameter C is defined
 - A low value of C leads to more margin violations
 - A high value of C limits the flexibility

$$\begin{aligned} \text{Minimize: } & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \text{ (Slack variable)} \\ \text{Subject to: } & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad (i = 1 \sim N) \\ & \xi_i \geq 0 \end{aligned}$$

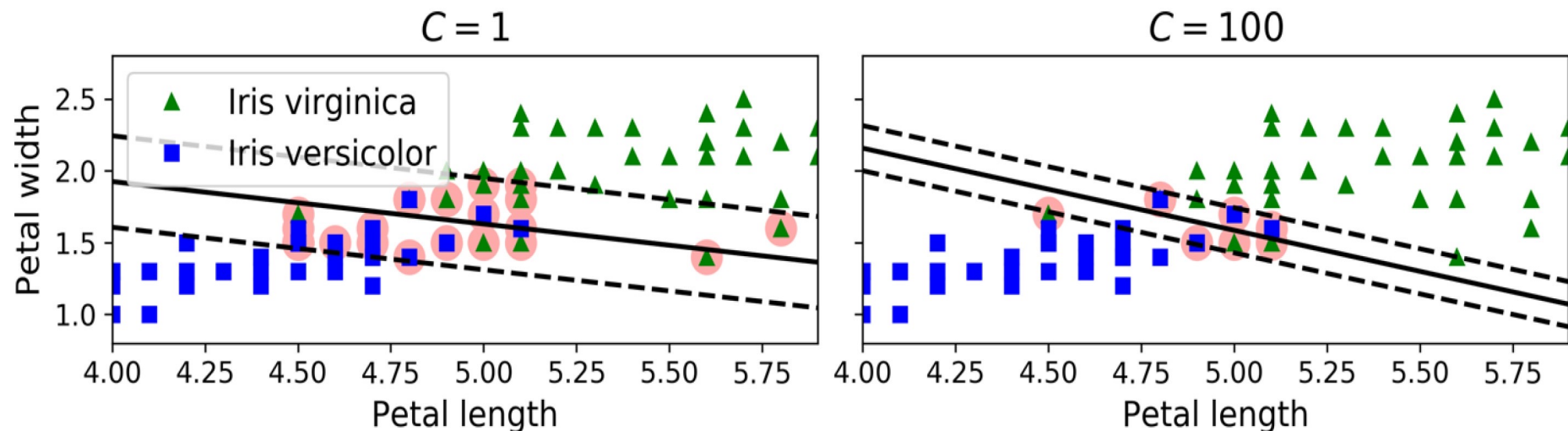


Figure 5-4. Large margin (left) versus fewer margin violations (right)

Tip:

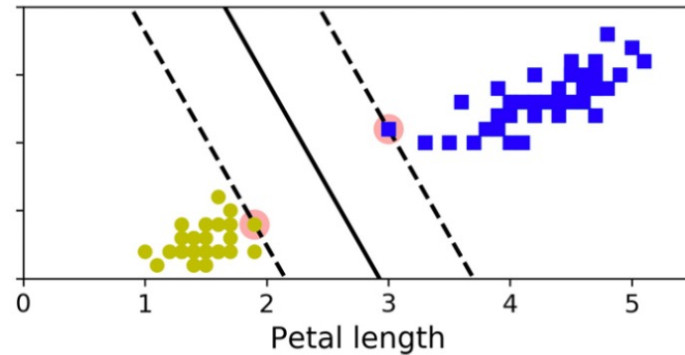
- To regularize SVM, reduce C
- C is inverse of regularizing hyperparameter α

$$J(\boldsymbol{\theta}) = \text{MSE}(\boldsymbol{\theta}) + \alpha \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

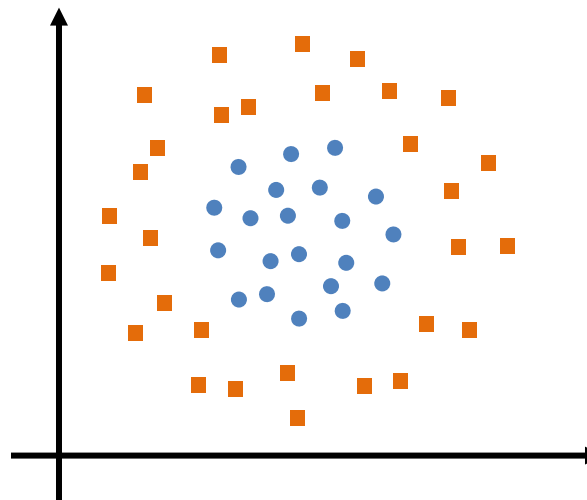
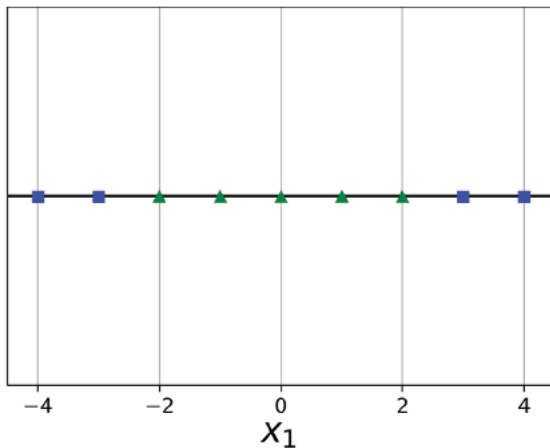


Nonlinear SVM Classification

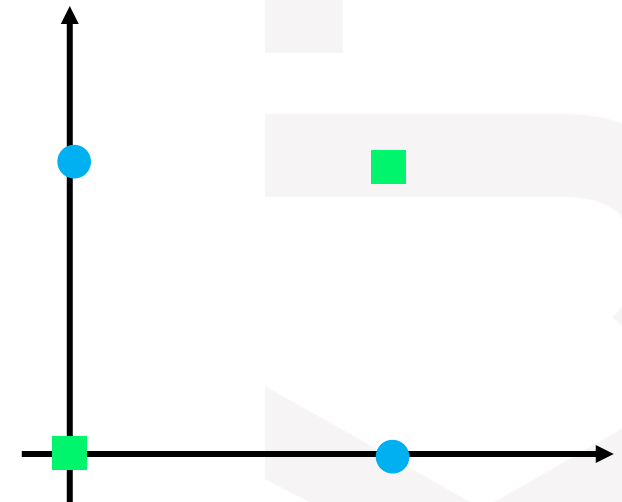
Linearly separable



How do we deal with these cases?



Not linearly separable





Nonlinear SVM Classification

How do we deal with nonlinearity?

Solution

1. Directly tackle nonlinear problems

Multilayer Perceptron, Deep Learning, etc.

2. Transform nonlinear problems into linear issues

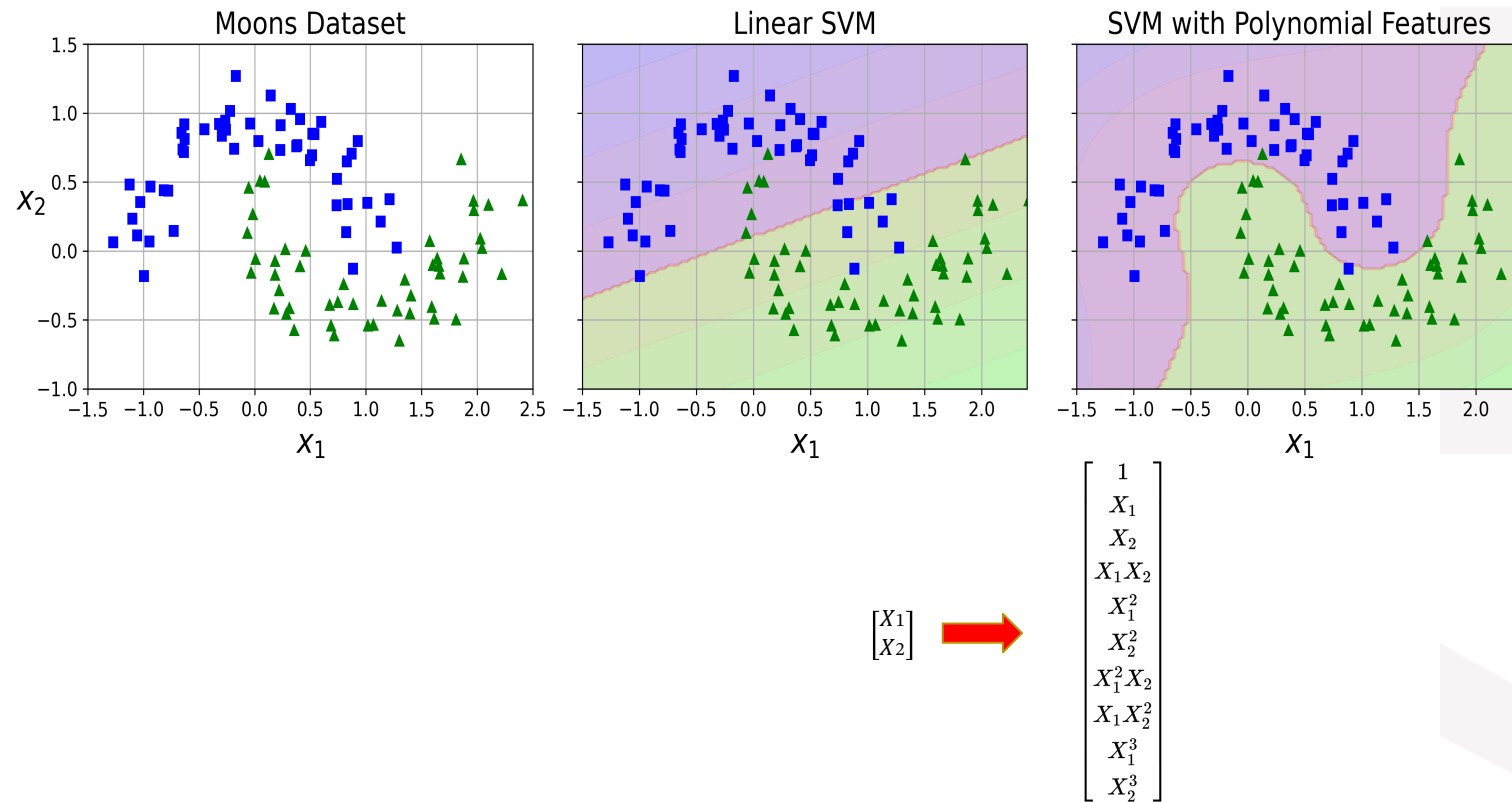
Kernel SVM



Nonlinear SVM Classification

- Polynomial Features

Polynomial features involve taking an existing feature and raising it to a power. This is useful for capturing non-linear relationships between the feature and the target variable. For example, if you have a feature X , polynomial features could include X^2 , X^3 , etc.



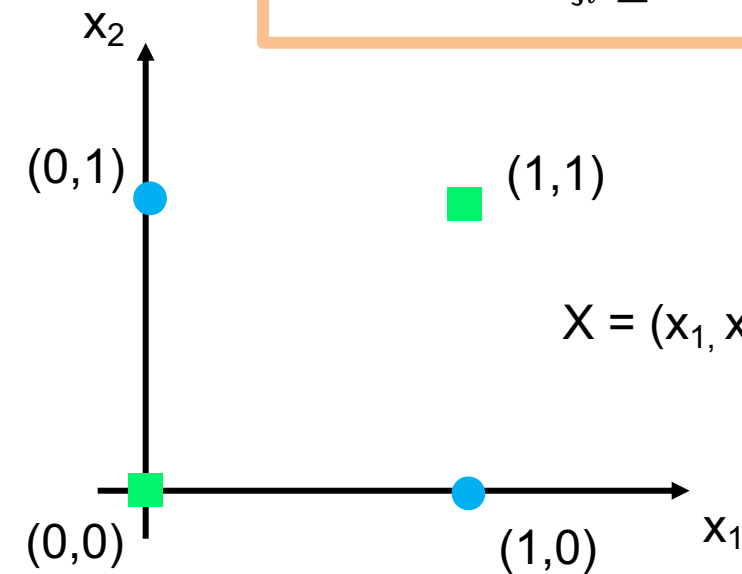
Polynomial features of degree 3
2D to 10D



Nonlinear SVM Classification

$$\text{Minimize: } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \text{ (Slack variable)}$$

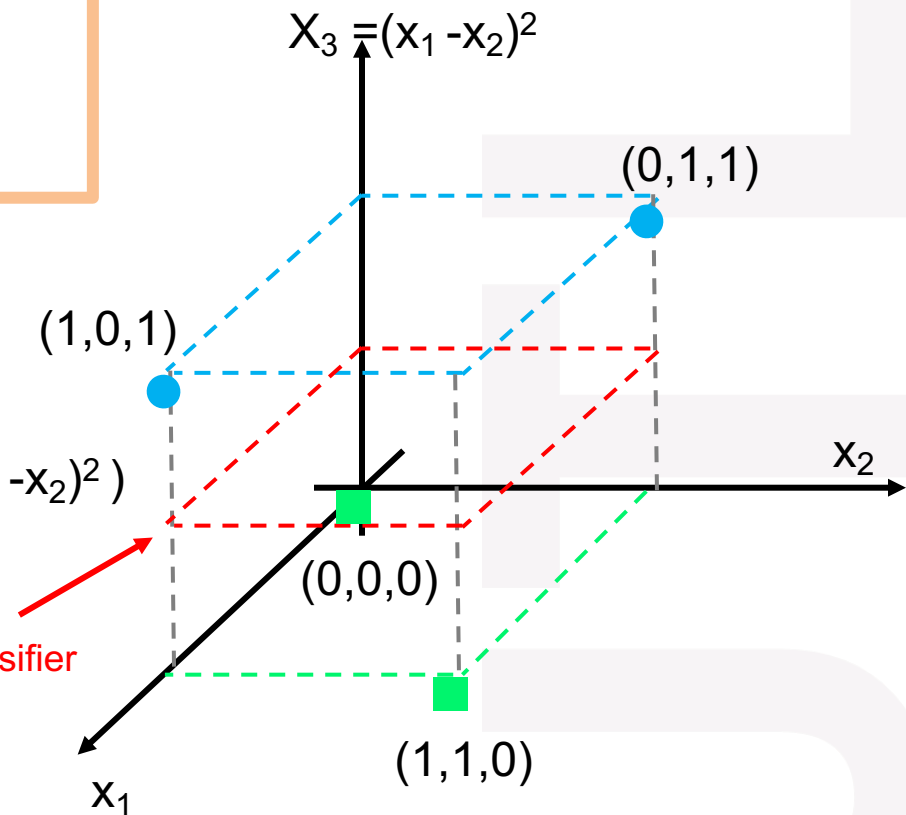
$$\text{Subject to: } \xi_i \geq 0$$



$$X = (x_1, x_2) \xrightarrow{\phi(x)} Z = (x_1, x_2, (x_1 - x_2)^2)$$

Add polynomial features

linear SVM classifier



Nonlinear transformation $\phi(x)$: not only one form

Cover's theorem: High-dimensional space is more likely to be linearly separable than in a low-dimensional space.

https://en.wikipedia.org/wiki/Cover's_theorem



Nonlinear SVM Classification

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

Subject to: $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad (i = 1 \sim N)$



$$\max_{\lambda} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \boxed{x_i^\top x_j} + \sum_{i=1}^N \lambda_i, \text{ s.t. } \lambda_i \geq 0$$



$\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$

$\phi(\mathbf{x})_i$ High dimension (infinite)

Difficult to calculate!

How to solve $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$?

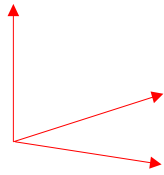


Nonlinear SVM: Kernel Trick

Input Space: dimension n

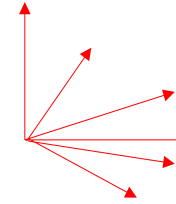
High-dimensional Feature Space: dimension $N \gg n$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$



$$\phi(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) \\ \phi_2(\mathbf{x}) \\ \phi_3(\mathbf{x}) \\ \vdots \\ \phi_N(\mathbf{x}) \end{bmatrix}$$

$N \gg n$



Expensive operation and requires large memory

$$K(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T \mathbf{b} = \begin{bmatrix} a_1 & a_2 & a_3 & \dots & a_n \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}$$

$$K(\phi(\mathbf{a}), \phi(\mathbf{b})) = \underbrace{\phi(\mathbf{a})^T}_{\phi(\mathbf{a})^T \phi(\mathbf{b}) = \text{function}(\mathbf{a}^T \mathbf{b})} \phi(\mathbf{b}) = \begin{bmatrix} \phi_1(\mathbf{a}) & \phi_2(\mathbf{a}) & \phi_3(\mathbf{a}) & \dots & \phi_N(\mathbf{a}) \end{bmatrix} \begin{bmatrix} \phi_1(\mathbf{b}) \\ \phi_2(\mathbf{b}) \\ \phi_3(\mathbf{b}) \\ \vdots \\ \phi_N(\mathbf{b}) \end{bmatrix}$$

Kernel Trick

Common kernels:

Linear: $K(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T \mathbf{b}$

Polynomial: $K(\mathbf{a}, \mathbf{b}) = (\gamma \mathbf{a}^T \mathbf{b} + r)^d$

Gaussian Radial Basis Function: $K(\mathbf{a}, \mathbf{b}) = \exp(-\gamma \|\mathbf{a} - \mathbf{b}\|^2)$

Sigmoid: $K(\mathbf{a}, \mathbf{b}) = \tanh(\gamma \mathbf{a}^T \mathbf{b} + r)$

Universal approximator.
Corresponding feature space
 $\phi(\mathbf{x})$ is infinite dimensional space

non-linearly separable data

infinite-dimensional space



Nonlinear SVM: Kernel Method

Linearly separable	A little violations	Strict nonlinearity
Hard-margin SVM	Soft-margin SVM	$\phi(x)$ +hard-margin -> kernel SVM

Kernel Method: Ideologically, transform low-dimensional non-linear space to high-dimensional linear space, using kernel function.

Kernel Function: Kernel Function = $\langle \phi(x), \phi(x') \rangle = \phi(x_i)^T \phi(x_j)$, $\langle \cdot \rangle$ means dot-product

It covers non-linear transformations and an inner product operation on nonlinear transformations.

Kernel Trick: Computationally, avoiding explicitly computing the transformation to another feature space.

- Project input space into a very high-dimensional feature space, may be even infinity
- Problem:
 - Projecting training data in to a high-dimensional space is expensive
 - large number of parameters
- Trick:
 - Compute dot-product between training samples in the projected high-dimensional space without ever projecting.