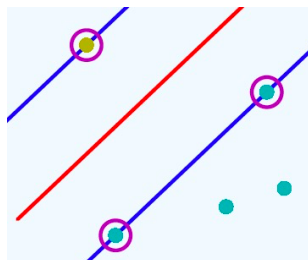




MACHINE LEARNING

COMP8220

04 – Support Vector Machine



Acknowledgement



MACQUARIE
University

- ❖ Some slides are from the S1 version designed by Rolf Schwitter (Rolf.Schwitter@mq.edu.au)



Some slides are also based on book of Christopher Bishop "Pattern Recognition and Machine Learning" Springer-Verlag New York (2006), and the slides made by Torsten Möller from this book.



❖ Linear Support Vector Machine

- How SVM comes?
- SVM learning problem formulation
- Solving SVM optimization problem

❖ SVM for Non-linearity

- SVM with soft margin
- Kernel trick

- ❖ Linear Support Vector Machine
 - How SVM comes?
 - SVM learning problem formulation
 - Solving SVM optimization problem

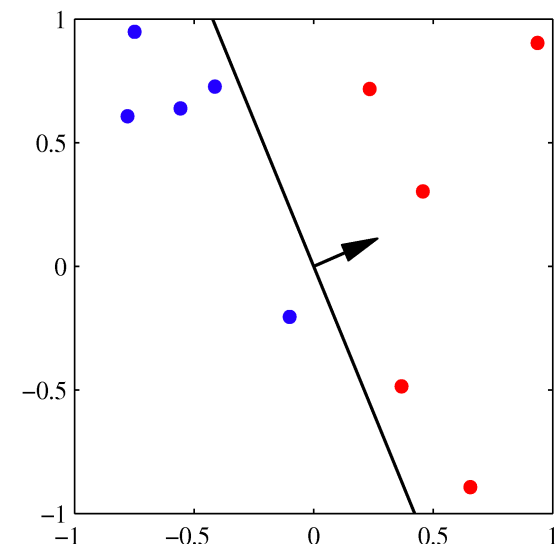
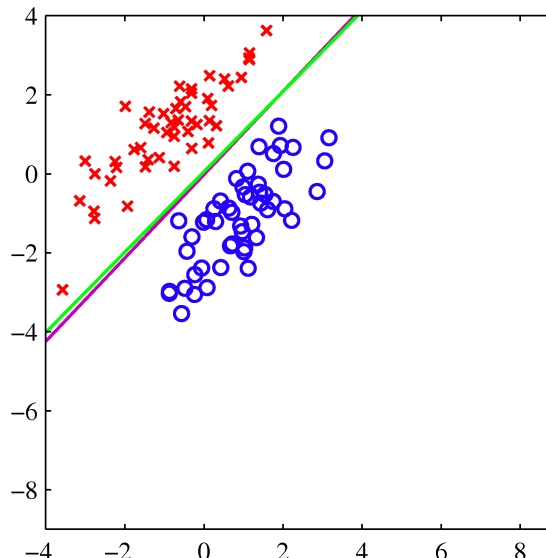
- ❖ Two methods for the discriminant function

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

Or more generally with feature map function,

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + w_0$$

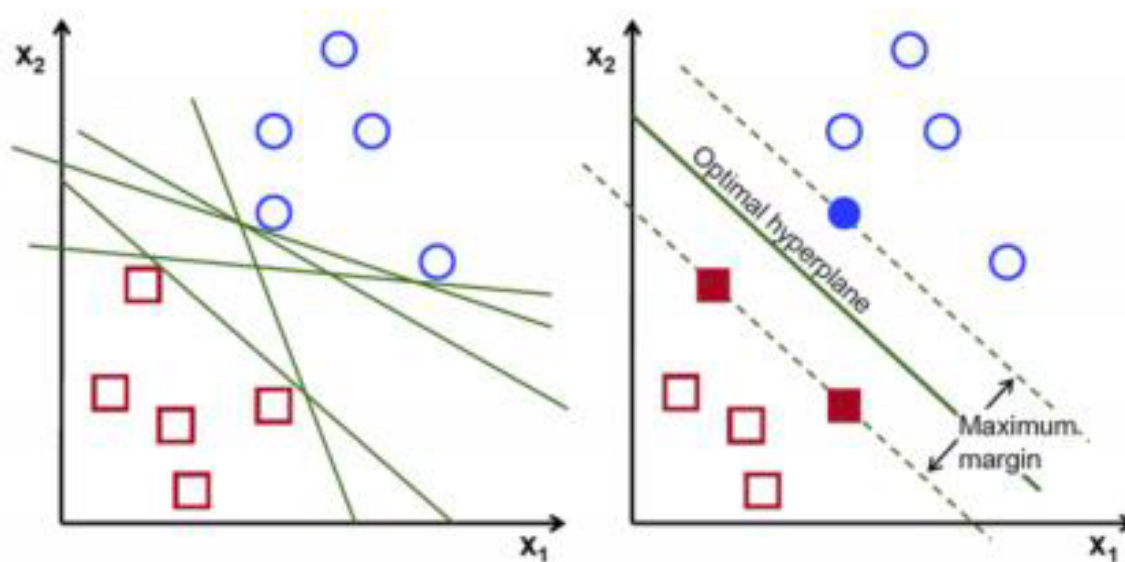
- Least squares for classification
- Perceptron
 - There are many possible optimal decision boundaries for perceptron
 - Which is the best?



Max Margin Criterion

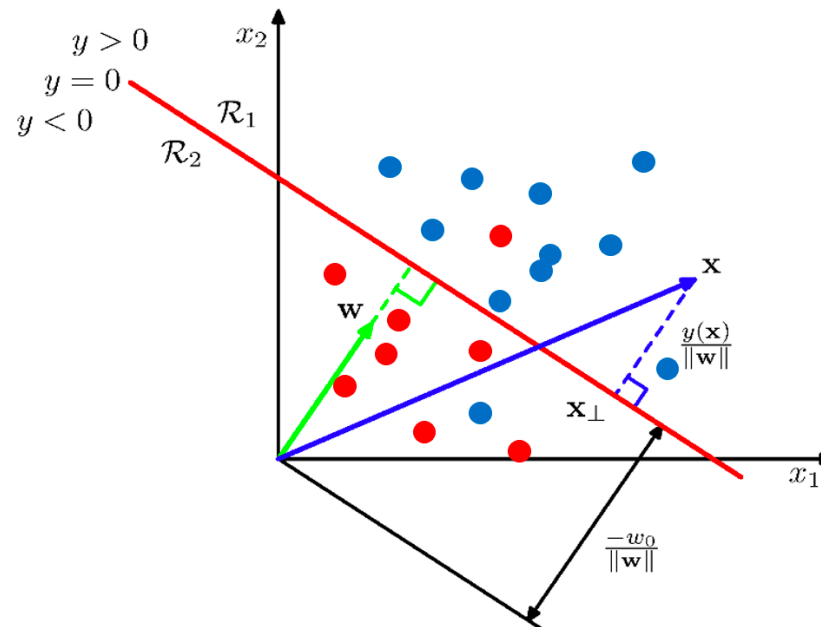


- ❖ **Margin** of a classification boundary (hyperplane): the minimum distance to any example



- ❖ Idea of Support Vector Machines (SVM)
 - Choose the decision boundary which maximizes the margin

- ❖ How to calculate the margin? (two-class case)
 - How to calculate distance between a point and a hyperplane?



- **Signed distance** of x to the decision surface (denoted as r)

$$r = \frac{y(x)}{\|w\|}$$

- ❖ Then, we can get (unsigned) distance

$$t \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

- t takes values in $\{1, -1\}$

- ❖ The margin can be calculated by

$$\begin{aligned} & \min_n \left\{ t_n \frac{y(\mathbf{x}_n)}{\|\mathbf{w}\|} \right\} \\ &= \frac{1}{\|\mathbf{w}\|} \min_n \{ t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \} \end{aligned}$$

- $\|\mathbf{w}\|$ is independent of n
- The minimum distance to any example

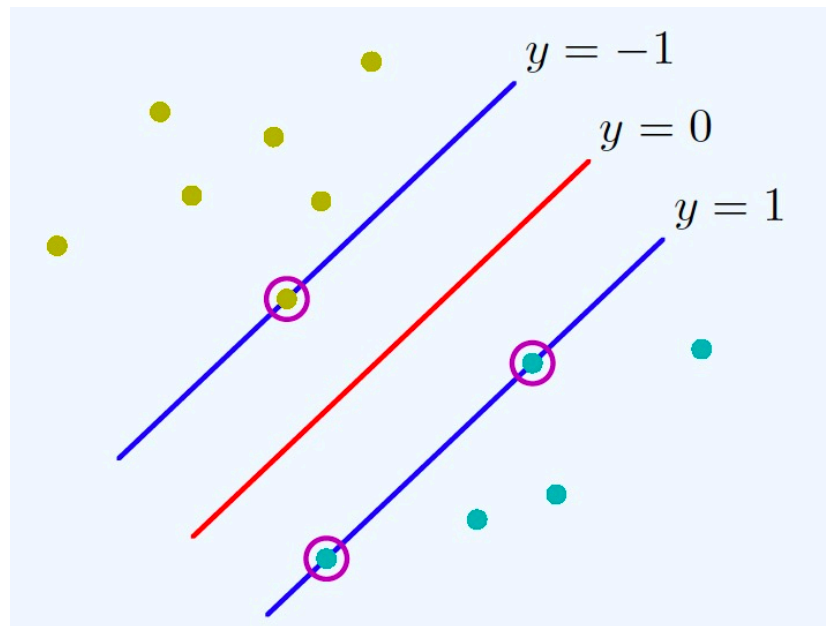
SVM Learning Problem



- ❖ The learning problem is to choose \mathbf{w} and w_0 to maximise the margin

$$\arg \max_{\mathbf{w}, w_0} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n \{t_n (\mathbf{w}^T \mathbf{x}_n + w_0)\} \right\}$$

- Points with the min value are known as **support vectors**



- ❖ The **unconstrained optimization problem** is complex

$$\arg \max_{\mathbf{w}, w_0} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n \{t_n(\mathbf{w}^T \mathbf{x}_n + w_0)\} \right\}$$

- ❖ Let $C = \min_n \{t_n(\mathbf{w}^T \mathbf{x}_n + w_0)\}$, we can have

$$t_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq C$$

- ❖ So, the optimization problem can be re-formulated as

$$\arg \max_{\mathbf{w}, w_0} \left\{ C \frac{1}{\|\mathbf{w}\|} \right\}$$

s.t.

$$t_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq C$$

❖ An use property of a hyperplane

- If $\mathbf{w} = \kappa \mathbf{v}$ and $w_0 = \kappa v_0$, the two hyperplanes are the same, where κ is a non-zero constant

$$\mathbf{w}^T \mathbf{x} + w_0 = 0 \text{ and } \mathbf{v}^T \mathbf{x} + v_0 = 0$$

- So, we can always find a κ to make

$$\min_n \{t_n (\mathbf{v}^T \mathbf{x}_n + v_0)\} = 1$$

❖ Then, we can have another optimization problem

$$\arg \max_{\mathbf{v}, v_0} \left\{ \frac{1}{\|\mathbf{v}\|} \right\}$$

s.t.

$$t_n (\mathbf{v}^T \mathbf{x}_n + v_0) \geq 1$$

❖ Observation

- The two optimization problem will produce the same hyperplane (classification decision boundary)
- Just denotation difference: \mathbf{w} vs \mathbf{v} , and w_0 vs v_0
- But the second form is easier to solve!

❖ Then, we just need to solve the following constrained optimization problem

$$\arg \max_{\mathbf{w}, w_0} \left\{ \frac{1}{\|\mathbf{w}\|} \right\}$$

s.t.

$$t_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$$

❖ Further observations

- Maximising $\frac{1}{\|\mathbf{w}\|}$ is equivalent to minimizing $\|\mathbf{w}\|$
- Minimizing $\|\mathbf{w}\|$ is equivalent to minimizing $\frac{1}{2} \|\mathbf{w}\|^2$

❖ So, we can have the **canonical representation** of the SVM learning problem

$$\mathbf{w}^*, w_0^* = \arg \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2$$

s.t.

$$t_n(\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1$$

- ❖ To solve the constrained optimization problem, we need to use the Lagrangian Multipliers method to convert it to an unconstrained one

$$L(\mathbf{w}, w_0, \mathbf{a}) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{n=1}^N a_n [t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1]$$

- ❖ Set the derivatives of L w.r.t. \mathbf{w} and w_0 to 0, we get

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

$$0 = \sum_{n=1}^N a_n t_n$$

- ❖ Plugging those equations into L removes \mathbf{w} and w_0 , and results in a version of L where $\nabla_{\mathbf{w}, w_0} L = 0$:

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \mathbf{x}_n^T \mathbf{x}_m$$

- ❖ This is the dual representation of the problem
$$\arg \max_{\mathbf{a}} \tilde{L}(\mathbf{a})$$

s.t.

$$\sum_{n=1}^N a_n t_n = 0$$

❖ The dual representation is easier to solve

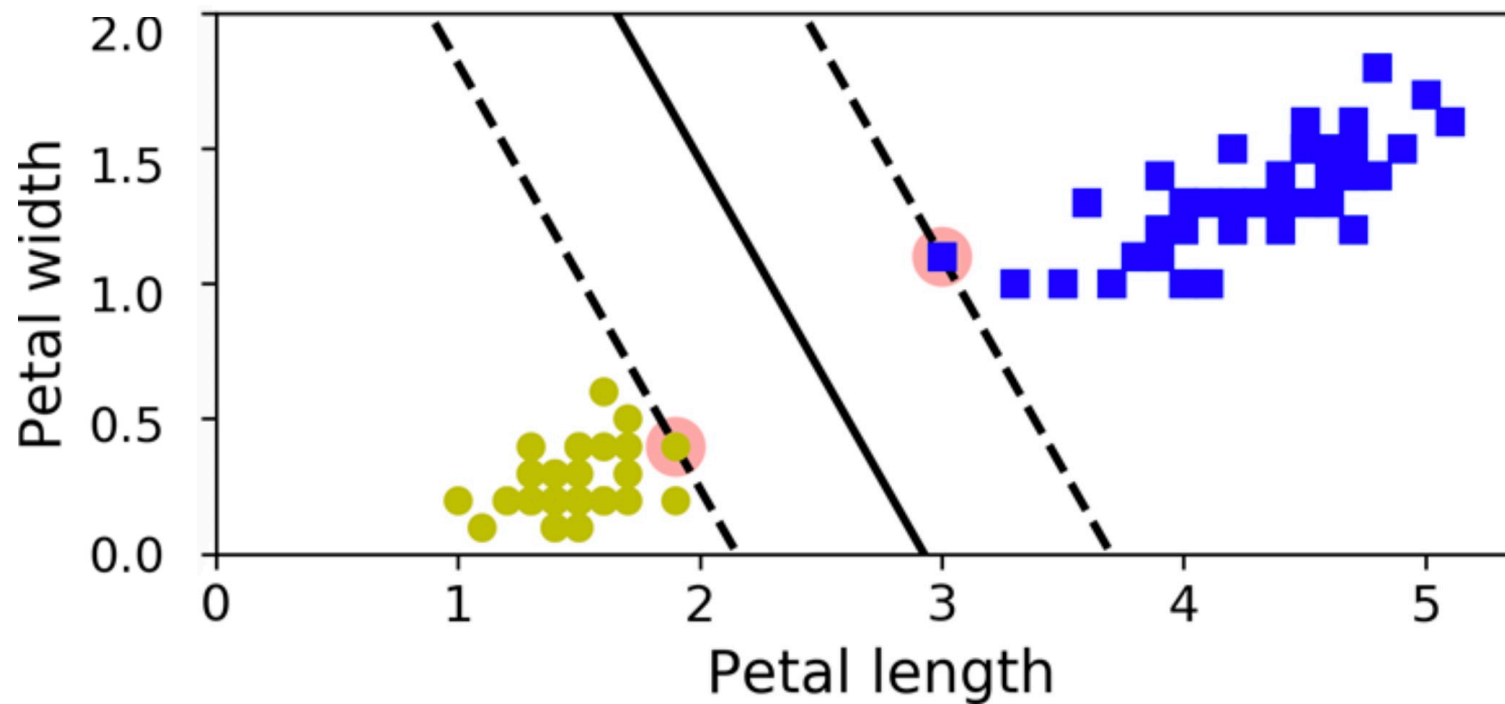
- Quadratic and convex in α
- Kernelized with K positive semi-definite

❖ Results

$$\mathbf{w}^* = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

- Recall $a_n(t_n y(\mathbf{x}_n) - 1) = 0$ condition Lagrange. This makes either $a_n = 0$ or \mathbf{x}_n is a support vector
- a_n will be sparse – many zeros
- Another formula for computing w_0^*
- Generally, \mathbf{x}_n can be replaced by $\phi(\mathbf{x}_n)$

Example





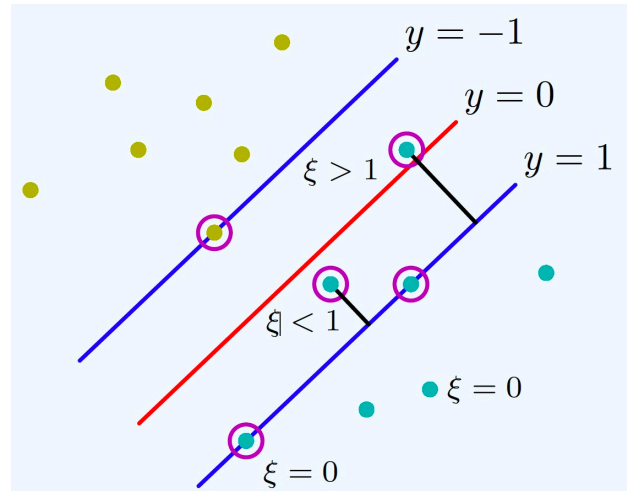
❖ Linear Support Vector Machine

- How SVM comes?
- SVM learning problem formulation
- Solving SVM optimization problem

❖ SVM for Non-linearity

- SVM with soft margin
- Kernel trick

- ❖ For most problems, data will not be linearly separable (even in feature space $\phi(\mathbf{x})$)
- ❖ We can relax the constraints from $t_n y(\mathbf{x}_n) \geq 1$ to
$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n$$
 - $\xi_n \geq 0$ are called **slack variables**
 - $\xi_n = 0$ satisfy original problem

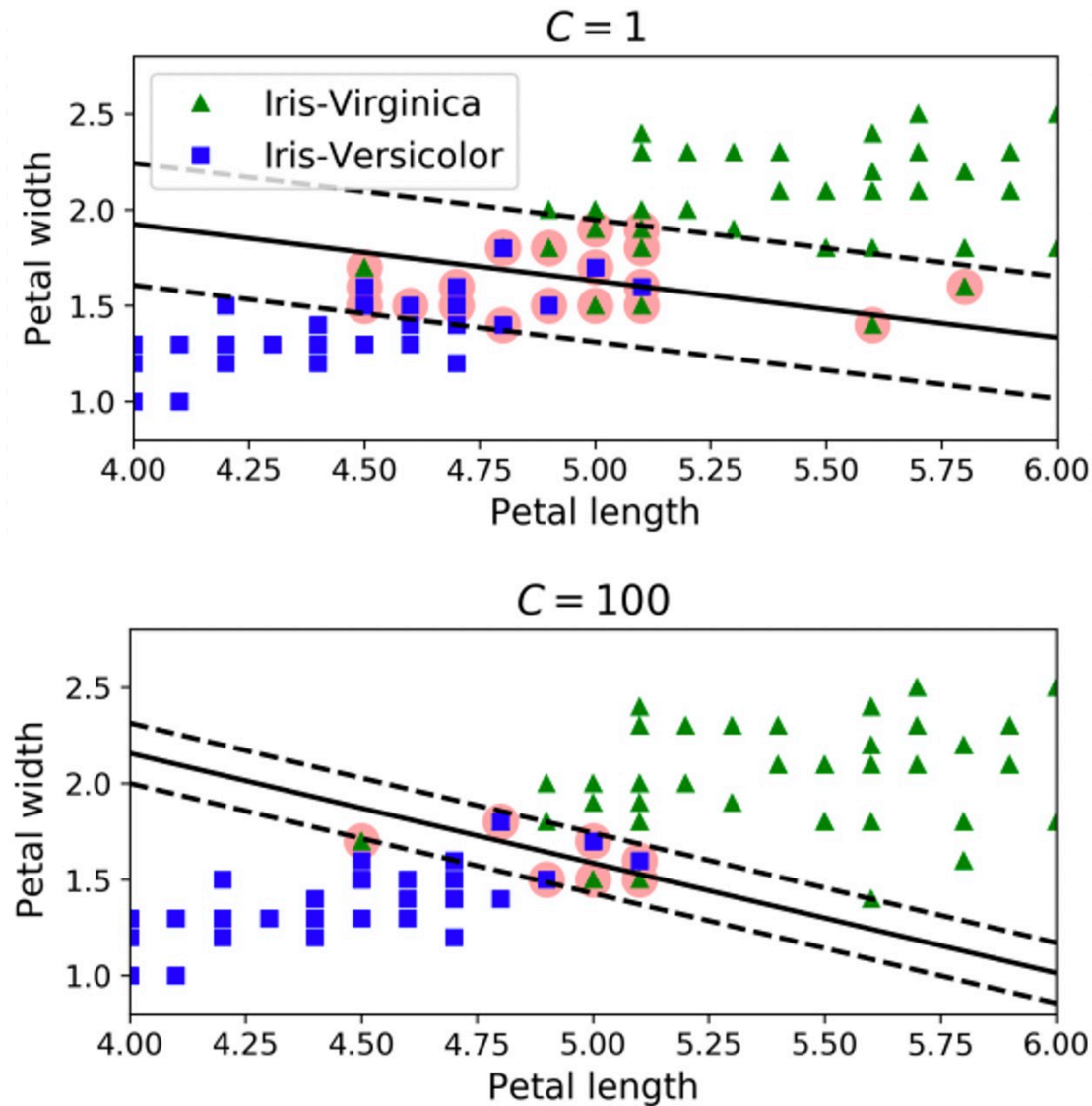


- ❖ Basic idea: non-zero slack variables are bad, penalize while maximizing the margin

$$\arg \min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \right\}$$

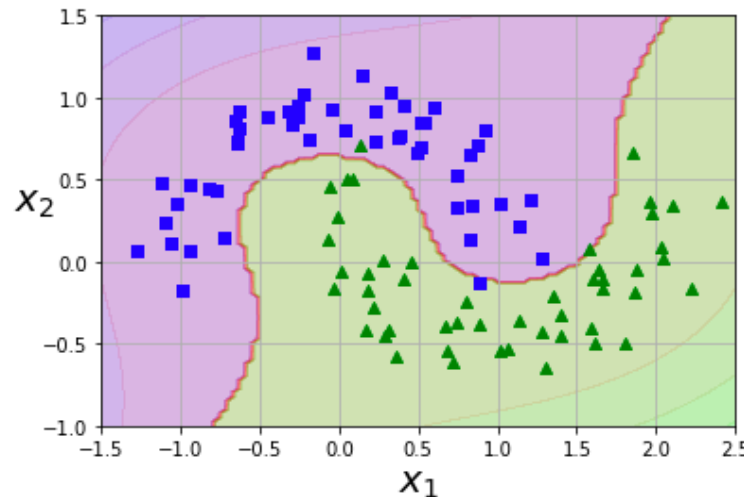
- Optimization is same quadratic, different constraints, convex
- Constant $C > 0$ controls the importance of large margin versus violation (non-zero slack)
 - Set using cross-validation (hyperparameter)
- Also, C controls the complexity of the model
 - If the SVM model is overfitting, then reduce C

Soft Margin Example



Nonlinear SVM classification

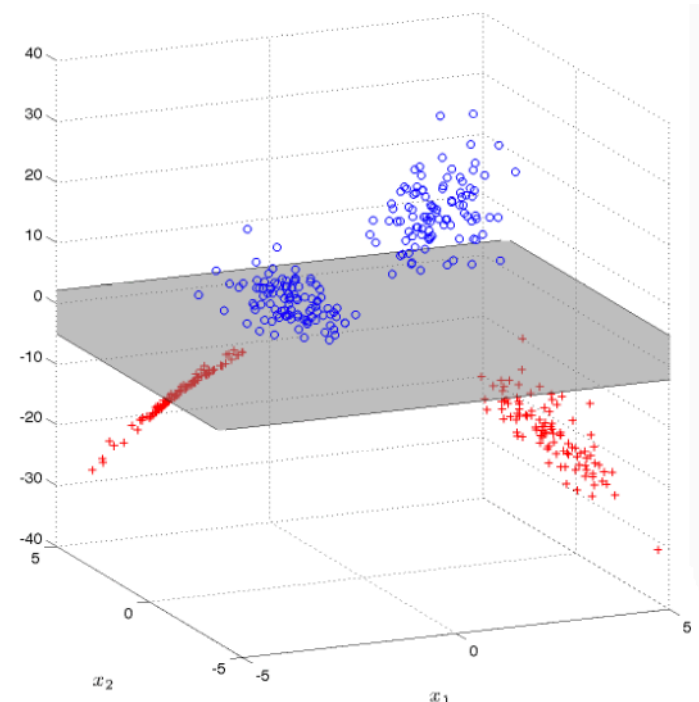
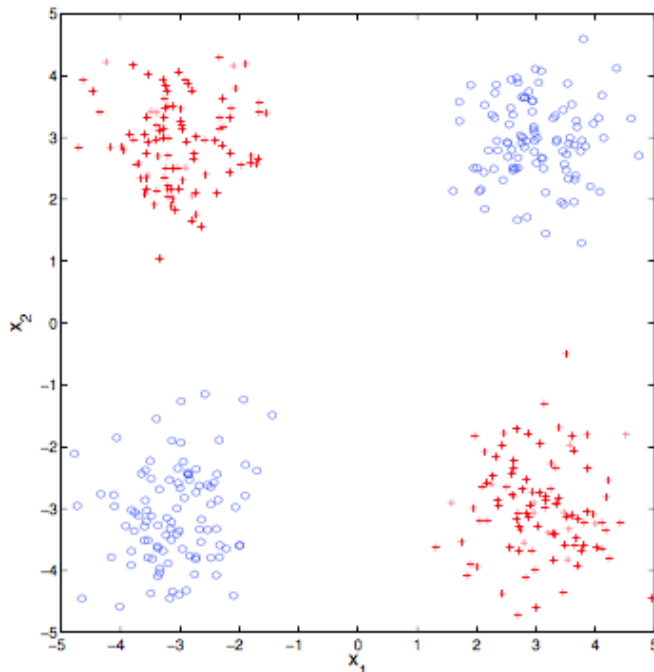
- ❖ Many datasets are not linearly separable



- ❖ One solution is to add polynomial features
 - May work with all sorts of ML algorithms, not only SVMs
 - This can result in linearly separable datasets
- ❖ We can then use a linear SVM classifier (SVC) together

Separation in High Dimension

- ❖ Separation may be easier in a higher dimension
 - As we investigate more features
 - XOR data example: Mapping the data from a 2-D space to a 3-D space allows us to find a separation hyperplane



- ❖ Mapping to high-dimensional feature space is promising
 - But the computation is usually expensive
- ❖ This is where the kernel trick comes in
 - The kernel trick allows us to operate in the original feature space without computing the coordinates of the data in a higher dimensional space
- ❖ A kernel is a function $K(\mathbf{x}_i, \mathbf{x}_j)$ that can compute the dot (inner) product $\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$ based only on the original vector \mathbf{x}_i and \mathbf{x}_j
 - No need to know the transformation function $\phi(\cdot)$
 - E.g., $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^2$ is a 2nd-degree polynomial kernel

Kernel Trick Example



Given two points: $x = (2, 3, 4)$ and $y = (3, 4, 5)$.

We have $K(x, y) = \langle f(x), f(y) \rangle$.

f is a map from n -dimensional to m -dimensional space.

$\langle x, y \rangle$ denotes the dot product.

Let us first calculate $\langle f(x), f(y) \rangle$.

$$f(x) = (x_1x_1, x_1x_2, x_1x_3, x_2x_1, x_2x_2, x_2x_3, x_3x_1, x_3x_2, x_3x_3)$$

$$f(y) = (y_1y_1, y_1y_2, y_1y_3, y_2y_1, y_2y_2, y_2y_3, y_3y_1, y_3y_2, y_3y_3)$$

In our case:

$$f(2, 3, 4) = (4, 6, 8, 6, 9, 12, 8, 12, 16) \text{ and}$$

$$f(3, 4, 5) = (9, 12, 15, 12, 16, 20, 15, 20, 25)$$

Kernel Trick Example

So the dot product is:

$$f(x) \cdot f(y) = f(2, 3, 4) \cdot f(3, 4, 5) = \\ (36 + 72 + 120 + 72 + 144 + 240 + 120 + 240 + 400) = 1444$$

Now let us use the kernel instead:

$$K(x, y) = (2*3 + 3*4 + 4*5)^2 = (6 + 12 + 20)^2 = \\ 38 * 38 = 1444$$

- ❖ The first method requires a lot of calculations because of projecting 3 dimensions into 9 dimensions, while using the kernel is much easier

- ❖ As we already have

$$\mathbf{w}^* = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

- ❖ Then, the discriminant function from SVM

$$\begin{aligned} y(\mathbf{x}) &= \mathbf{w}^{*T} \phi(\mathbf{x}) + w_0 \\ &= \sum_{n=1}^N a_n t_n \underbrace{\phi^T(\mathbf{x}_n) \phi(\mathbf{x})}_{K(\mathbf{x}, \mathbf{x}_n)} + w_0 \end{aligned}$$

Linear kernel $k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2$

- $\phi(\mathbf{x}) = \mathbf{x}$

Polynomial kernel $k(\mathbf{x}_1, \mathbf{x}_2) = (1 + \mathbf{x}_1^T \mathbf{x}_2)^d$

- Contains all polynomial terms up to degree d

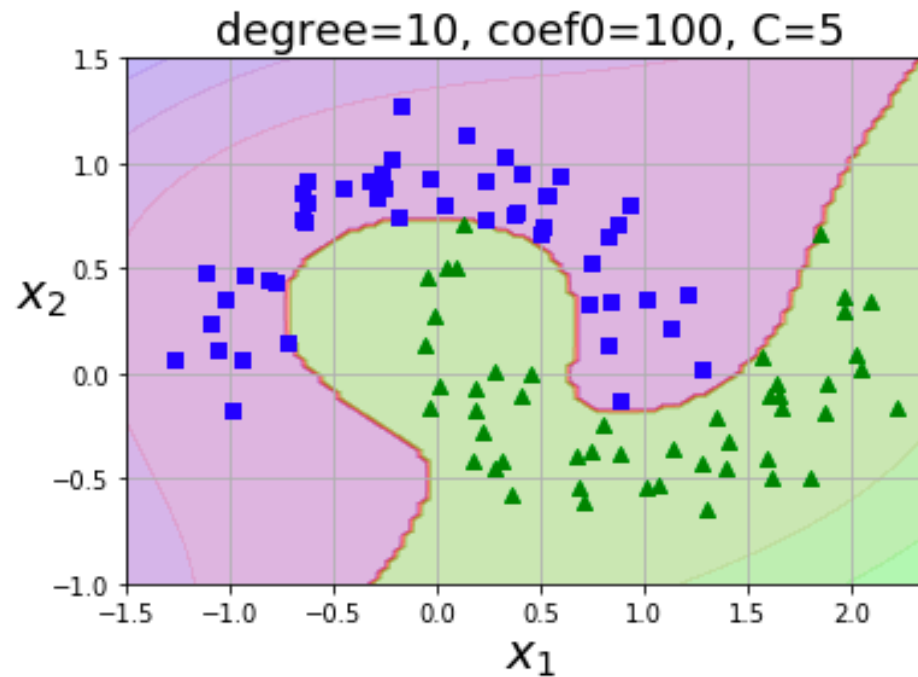
Gaussian kernel $k(\mathbf{x}_1, \mathbf{x}_2) = \exp(-||\mathbf{x}_1 - \mathbf{x}_2||^2 / 2\sigma^2)$

- Infinite dimension feature space

Example



❖ Output: SVC with Polynomial Kernel



Example

❖ Output: Output: SVC with RBF Kernel

