

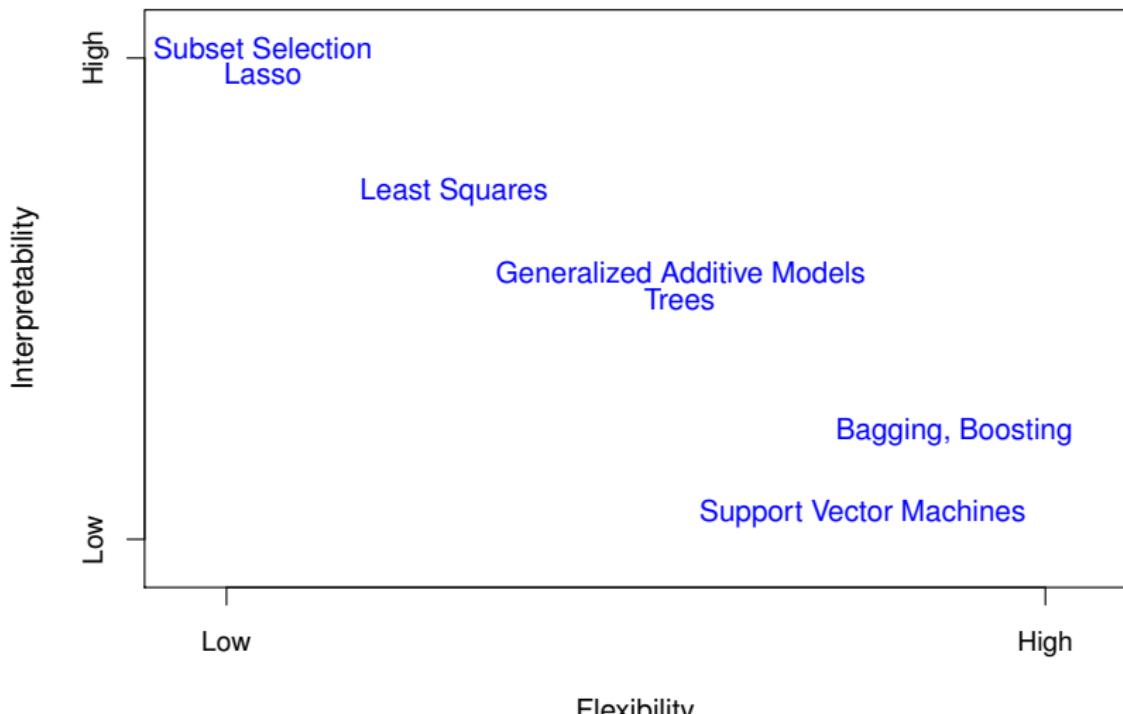
CPSC 8420 Advanced Machine Learning

Kernel Support Vector Machines

Support Vector Machines



Comparison of Different Machine Learning Models



Overview

- *Support Vector Machine* (SVM) is an approach for **supervised classification** that was developed in the computer science community in the 1990s.
- SVM is intended for **binary classification**.
- SVM can be considered an extension of the *support vector classifier*, which is an extension of the *maximum margin classifier*.

Vladimir N. Vapnik



Vladimir N. Vapnik

His main contributions are:

- ① Support Vector Machines
- ② Statistical Learning Theory
- ③ Vapnik–Chervonenkis theory
- ④ Kernel Method

He won numerous awards including:

- ① IEEE John von Neumann Medal (2017)
- ② Benjamin Franklin Medal (2012)
- ③ IEEE Neural Networks Pioneer Award (2010)
- ④ Fellow of the U.S. National Academy of Engineering (2006)

“Nothing is more practical than a good theory.”

Hyperplane

- The maximum margin classifier (and also its extension, the support vector classifier) is based on the idea of a **separating hyperplane**

Hyperplane

- The maximum margin classifier (and also its extension, the support vector classifier) is based on the idea of a **separating hyperplane**
- In two dimensions, a hyperplane is a flat one-dimensional subspace (i.e., a line), defined by:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

Hyperplane

- The maximum margin classifier (and also its extension, the support vector classifier) is based on the idea of a **separating hyperplane**
- In two dimensions, a hyperplane is a flat one-dimensional subspace (i.e., a line), defined by:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

- In three dimensions, a hyperplane is a flat two-dimensional subspace (i.e., a plane), defined by:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 = 0$$

Hyperplane

- The maximum margin classifier (and also its extension, the support vector classifier) is based on the idea of a **separating hyperplane**
- In two dimensions, a hyperplane is a flat one-dimensional subspace (i.e., a line), defined by:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

- In three dimensions, a hyperplane is a flat two-dimensional subspace (i.e., a plane), defined by:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 = 0$$

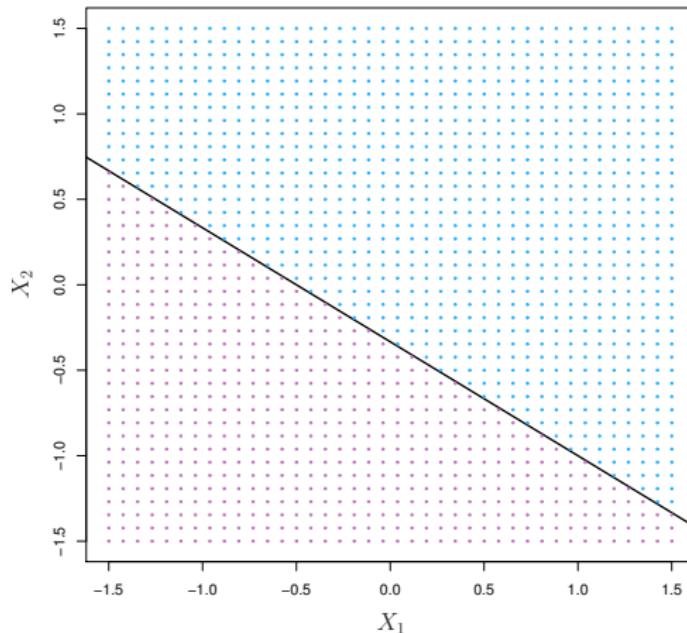
- In p dimensions, a hyperplane is difficult to visualize, but can easily be defined by extending the above equations to:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p = 0$$

Hyperplane

The hyperplane $1 + 2X_1 + 3X_2 = 0$.

- Blue: the set of points for which $1 + 2X_1 + 3X_2 > 0$
- Purple: the set of points for which $1 + 2X_1 + 3X_2 < 0$



Separating Hyperplane

- Suppose that it is possible to construct a hyperplane that separates the training observations perfectly according to their class labels.

Separating Hyperplane

- Suppose that it is possible to construct a hyperplane that separates the training observations perfectly according to their class labels.
- Label the observations above the hyperplane as $y_i = 1$ and those below the hyperplane as $y_i = -1$.

Separating Hyperplane

- Suppose that it is possible to construct a hyperplane that separates the training observations perfectly according to their class labels.
- Label the observations above the hyperplane as $y_i = 1$ and those below the hyperplane as $y_i = -1$.
- Then a separating hyperplane has the property that:

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} > 0 \quad \text{if } y_i = 1 \quad (1)$$

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} < 0 \quad \text{if } y_i = -1 \quad (2)$$

Separating Hyperplane

- Suppose that it is possible to construct a hyperplane that separates the training observations perfectly according to their class labels.
- Label the observations above the hyperplane as $y_i = 1$ and those below the hyperplane as $y_i = -1$.
- Then a separating hyperplane has the property that:

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} > 0 \quad \text{if } y_i = 1 \quad (1)$$

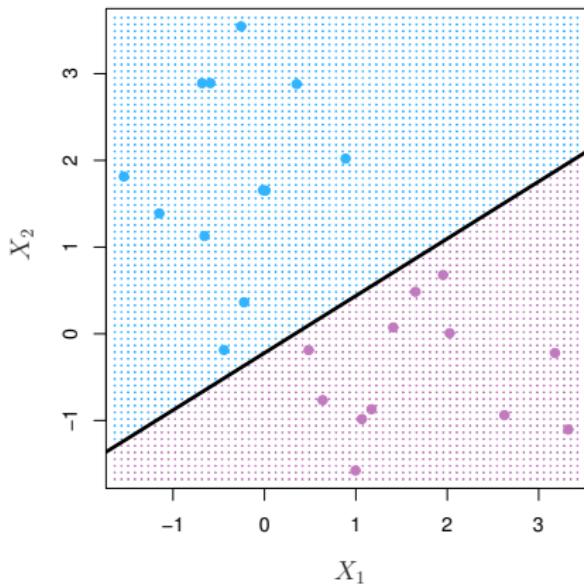
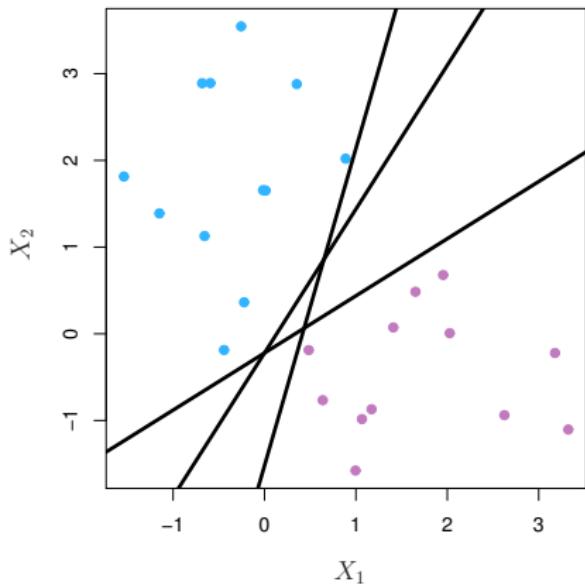
$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} < 0 \quad \text{if } y_i = -1 \quad (2)$$

- This can be written as:

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip}) > 0 \quad (3)$$

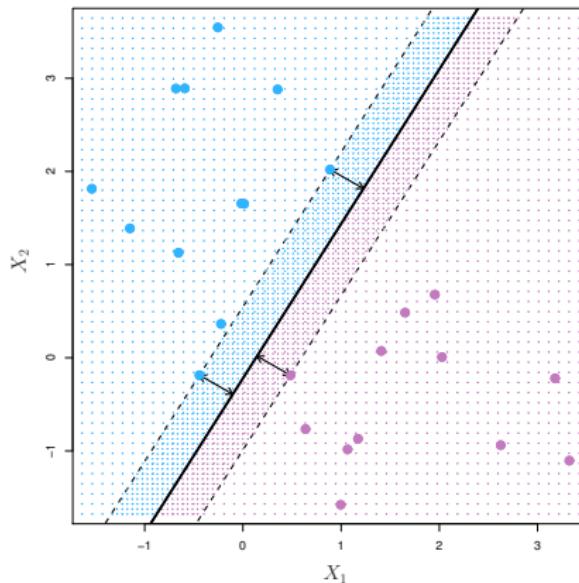
for all $i = 1, \dots, n$.

Maximum Margin Classifier



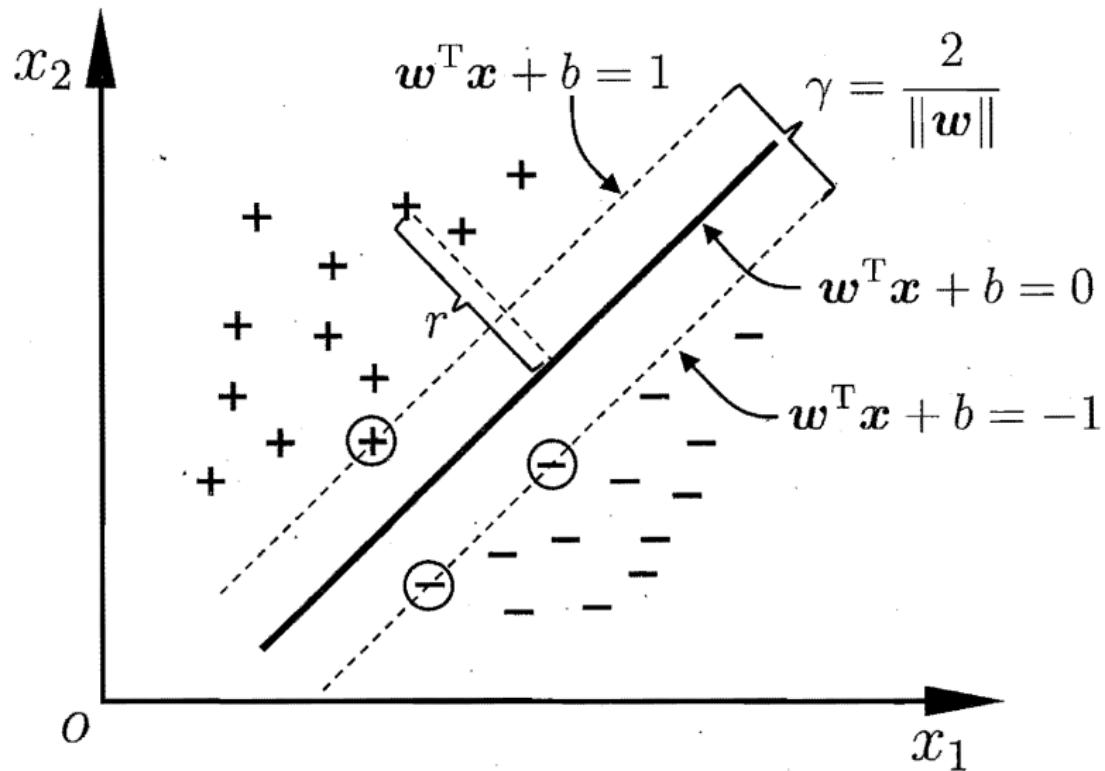
Example of data that can be perfectly separated using a hyperplane.
Problem: there exists an infinite number of such hyperplane.

Maximum Margin Classifier



- The **maximum margin classifier** uses the separating hyperplane that is farthest from the training observations (that is, has the largest margin).
- The three observations that lie along the dashed lines indicating the width of the margin are called **support vectors**.

Maximum Margin Classifier



Objective

The maximum margin hyperplane is the solution to the following optimization problem:

$$\max \gamma = \frac{y(w^T x + b)}{\|w\|_2} \quad s.t \quad y_i(w^T x_i + b) = \gamma'(i) \geq \gamma' \quad (i = 1, 2, \dots, m)$$

usually we set $\gamma' = 1$, then it is equivalent to:

$$\max \frac{1}{\|w\|_2} \quad s.t \quad y_i(w^T x_i + b) \geq 1 \quad (i = 1, 2, \dots, m)$$

$$\min \frac{1}{2} \|w\|_2^2 \quad s.t \quad y_i(w^T x_i + b) \geq 1 \quad (i = 1, 2, \dots, m)$$

Optimization

$$L(w, b, \alpha) = \frac{1}{2}||w||_2^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1] \quad s.t. \quad \alpha_i \geq 0$$
$$\underbrace{\min_{w,b}}_{\alpha_i \geq 0} \underbrace{\max_{\alpha_i \geq 0}}_{L(w,b,\alpha)} \quad (6)$$
$$\underbrace{\max_{\alpha_i \geq 0}}_{L(w,b,\alpha)} \underbrace{\min_{w,b}}_{\alpha_i \geq 0}$$

we can first optimize L w.r.t. w and b , namely $\underbrace{\min_{w,b}}_{L(w,b,\alpha)}$ by taking the derivative.

Optimization

$$L(w, b, \alpha) = \frac{1}{2}||w||_2^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1] \quad s.t. \quad \alpha_i \geq 0$$
$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \tag{7}$$
$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

Optimization

Now define $\psi(\alpha) = \min_{w,b} L(w, b, \alpha)$, we have:

$$\begin{aligned}\psi(\alpha) &= \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^m \alpha_i [y_i (w^T x_i + b) - 1] \\ &= \frac{1}{2} w^T \sum_{i=1}^m \alpha_i y_i x_i - w^T \sum_{i=1}^m \alpha_i y_i x_i - \sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i \\ &= -\frac{1}{2} w^T \sum_{i=1}^m \alpha_i y_i x_i - \sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i \\ &= -\frac{1}{2} \left(\sum_{i=1}^m \alpha_i y_i x_i \right)^T \left(\sum_{i=1}^m \alpha_i y_i x_i \right) - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^m \alpha_i y_i x_i^T \sum_{i=1}^m \alpha_i y_i x_i + \sum_{i=1}^m \alpha_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j\end{aligned}\tag{8}$$

Optimization

$$\underbrace{\max_{\alpha}}_{\alpha} - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^m \alpha_i \quad (9)$$

which is equivalent to:

$$\begin{aligned} & \underbrace{\min_{\alpha}}_{\alpha} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle - \sum_{i=1}^m \alpha_i \\ & s.t. \quad \sum_{i=1}^m \alpha_i y_i = 0 \\ & \alpha_i \geq 0 \quad i = 1, 2, \dots, m \end{aligned} \quad (10)$$

Optimization

We can make use of **Sequential Minimal Optimization (SMO)** to obtain α^* , which optimizes α pair-wise. By the definition

$w = \sum_{i=1}^m \alpha_i y_i x_i$, we have: $w^* = \sum_{i=1}^m \alpha_i^* y_i x_i$. To determine b^* , for

arbitrary (x_s, y_s) : $y_s(w^T x_s + b) = y_s(\sum_{i=1}^m \alpha_i y_i x_i^T x_s + b) = 1$, that is

$b_s^* = y_s - \sum_{i=1}^m \alpha_i y_i x_i^T x_s$. To determine the support vector, according to

KKT complementary condition: $\alpha_i^*(y_i(w^T x_i + b) - 1) = 0$, if $\alpha_i > 0$,
then $y_i(w^T x_i + b) = 1$. Therefore the hyperplane is: $w^* \bullet x + b^* = 0$,
and a new given data will be classified as $f(x) = sign(\langle w^*, x \rangle + b^*)$.

More Flexible Case

For non-separable case, we introduce slack variable $\xi_i \geq 0$ for each data (x_i, y_i) , such that: $y_i(w \bullet x_i + b) \geq 1 - \xi_i$.

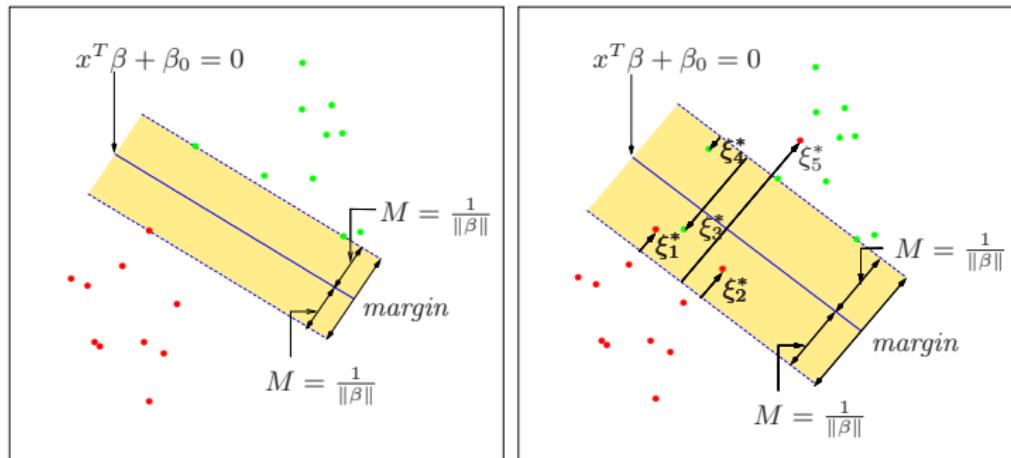


FIGURE 12.1. Support vector classifiers. The left panel shows the separable case. The decision boundary is the solid line, while broken lines bound the shaded maximal margin of width $2M = 2/\|\beta\|$. The right panel shows the nonseparable (overlap) case. The points labeled ξ_j^* are on the wrong side of their margin by an amount $\xi_j^* = M\xi_j$; points on the correct side have $\xi_j^* = 0$. The margin is maximized subject to a total budget $\sum \xi_i \leq \text{constant}$. Hence $\sum \xi_j^*$ is the total distance of points on the wrong side of their margin.

New Objective

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t. } & y_i(w^T x_i + b) \geq 1 - \xi_i \quad (i = 1, 2, \dots, m) \\ & \xi_i \geq 0 \quad (i = 1, 2, \dots, m) \end{aligned} \tag{11}$$

Similar to separable case above following Lagrangian Multipliers, we have:

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i \tag{12}$$

where $\mu_i \geq 0, \alpha_i \geq 0$ are Lagrangian variables.

Optimization

If we denote $J(w, b, \xi) = \frac{1}{2}||w||_2^2 + C \sum_{i=1}^m \xi_i$, then we have:

$$J = \underbrace{\max}_{\alpha_i \geq 0, \mu_i \geq 0} L(w, b, \alpha, \xi, \mu), \quad (13)$$

therefore, $\min J(w, b, \xi) = \underbrace{\min}_{w, b, \xi} \underbrace{\max}_{\alpha_i \geq 0, \mu_i \geq 0} L(w, b, \alpha, \xi, \mu)$. According to KKT condition, it is equivalent to $\underbrace{\max}_{\alpha_i \geq 0, \mu_i \geq 0} \underbrace{\min}_{w, b, \xi} L(w, b, \alpha, \xi, \mu)$.

Karush-Kuhn-Tucker (KKT) conditions

KKT is an extension for Lagrangian Multipliers, where there are inequality as well as equality constraints. assume we are given:

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k \\ & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

To solve it, we start by defining the generalized Lagrangian:

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Karush-Kuhn-Tucker (KKT) conditions

KKT will tell us the necessary condition of the minimizers:

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, d$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

Karush-Kuhn-Tucker (KKT) Practice

$$\min \frac{1}{2}(x^2 + y^2), \quad s.t. \quad x + y \geq 4 \quad (14)$$

Karush-Kuhn-Tucker (KKT) Practice

$$\min \frac{1}{2}(x^2 + y^2), \quad s.t. \quad x + 2y \leq 4 \quad (15)$$

Optimization

For $\max_{\alpha_i \geq 0, \mu_i \geq 0} \min_{w, b, \xi} L(w, b, \alpha, \xi, \mu)$, now we can first optimize

$\min_{w, b, \xi} L(w, b, \alpha, \xi, \mu)$ by taking the derivative and set it to be zero:

$$\begin{aligned}\frac{\partial L}{\partial w} = 0 &\Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \\ \frac{\partial L}{\partial b} = 0 &\Rightarrow \sum_{i=1}^m \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \xi} = 0 &\Rightarrow C - \alpha_i - \mu_i = 0\end{aligned}\tag{16}$$

Optimization

Now plugin the above to L :

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i \quad (17)$$

we have:

$$\begin{aligned} L(w, b, \xi, \alpha, \mu) &= \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i \\ &= \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1 + \xi_i] + \sum_{i=1}^m \alpha_i \xi_i \\ &= \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1] \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \end{aligned}$$

Optimization

$$\begin{aligned} & \underbrace{\max_{\alpha}}_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \\ & \text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0 \\ & C - \alpha_i - \mu_i = 0 \\ & \alpha_i \geq 0 \ (i = 1, 2, \dots, m) \\ & \mu_i \geq 0 \ (i = 1, 2, \dots, m) \end{aligned} \tag{19}$$

which is equivalent to:

$$\begin{aligned} & \underbrace{\min_{\alpha}}_{\alpha} \frac{1}{2} \sum_{i=1, j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^m \alpha_i \\ & \text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned} \tag{20}$$

Classification Discussion

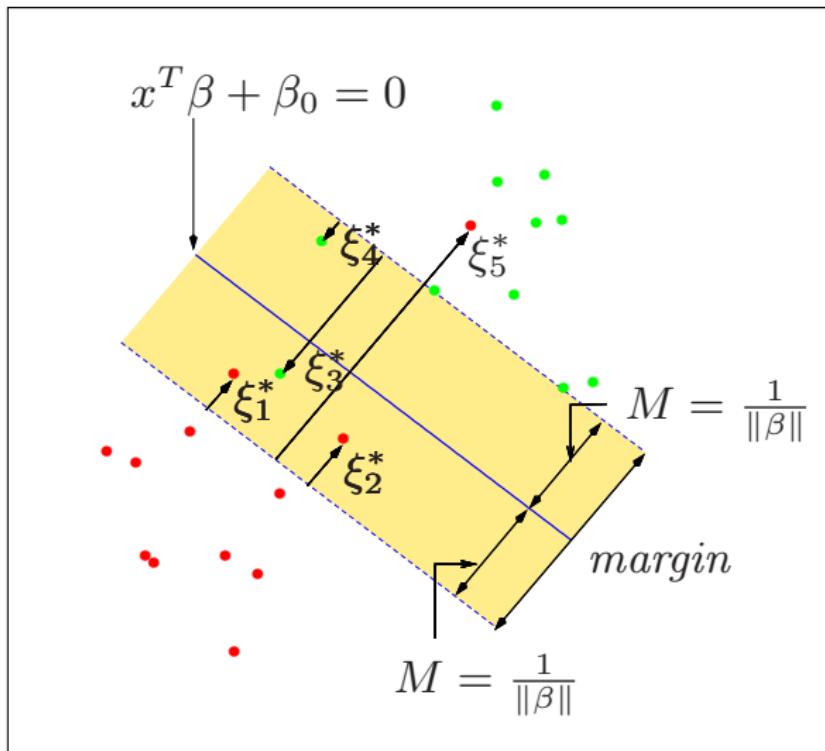
According to the complementary slackness condition, we have:

$$\begin{aligned}\alpha_i^*(y_i(w^T x_i + b) - 1 + \xi_i^*) &= 0 \\ \mu_i \xi_i &= 0\end{aligned}\tag{21}$$

also, since $C = \alpha_i + \mu_i$, we now discuss different cases:

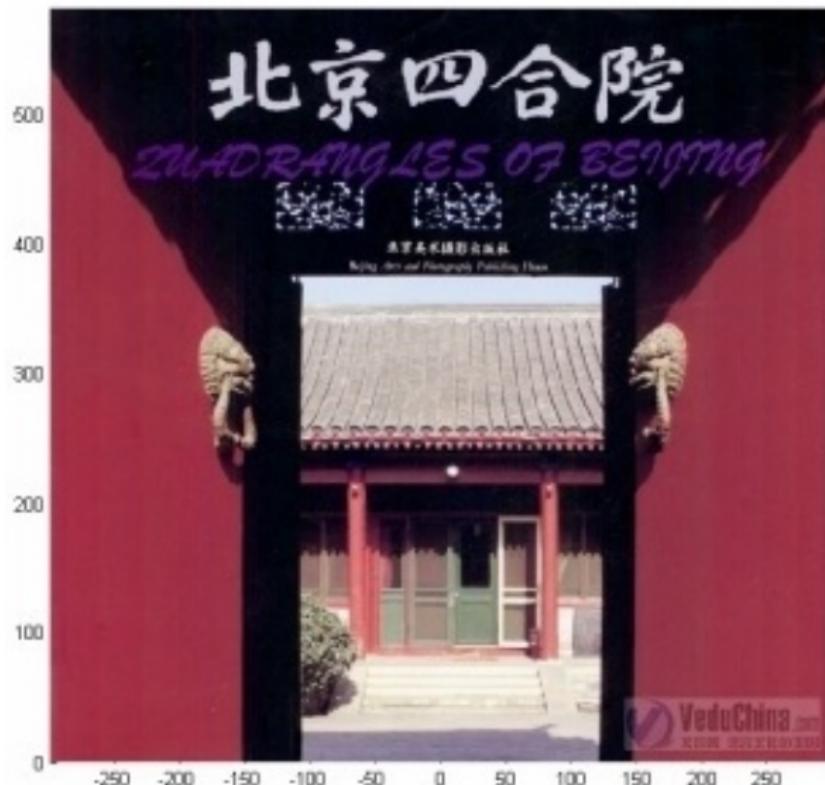
- ① $\alpha_i = 0 \implies \mu_i = C \implies \xi_i = 0$.
- ② $0 < \alpha_i < C \implies y_i(w^T x_i + b) - 1 + \xi_i^* = 0$, also
 $\mu_i > 0 \implies \xi_i = 0$, then $y_i(w^T x_i + b) - 1 = 0$, which is the support vector.
- ③ $\alpha_i = C$
 - ① $0 < \xi_i < 1$, it is correctly classified but on the wrong side of its margin.
 - ② $\xi_i = 1$, it lies exactly on the hyperplane.
 - ③ $\xi_i > 1$, mis-classified.

Classification Discussion



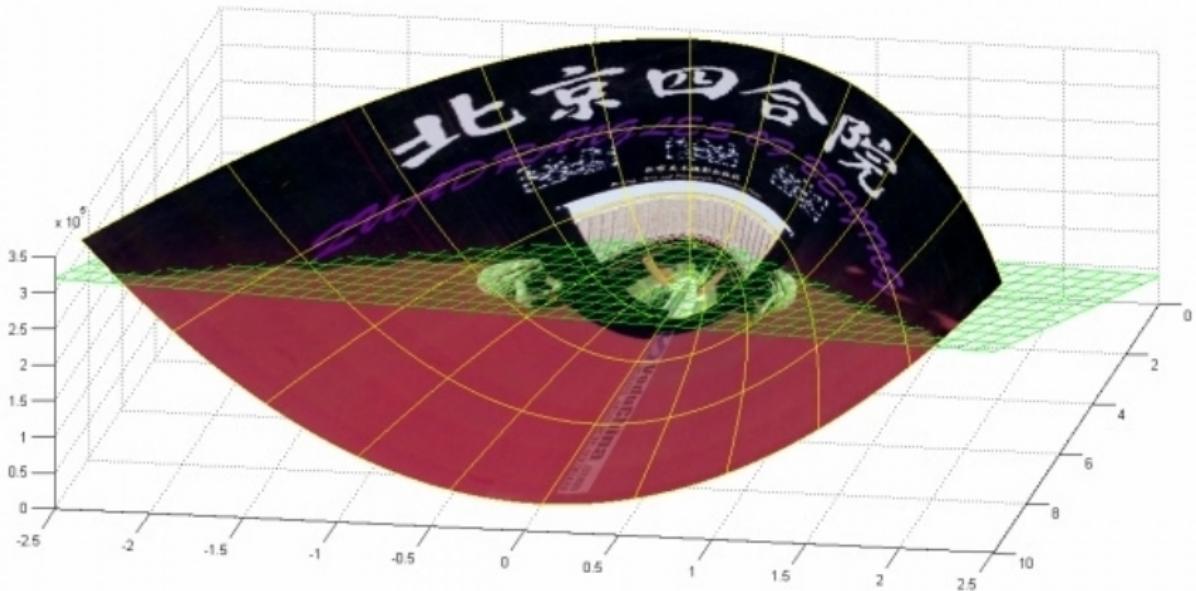
A Gentle Start

How can we separate the purple characters from the red part?



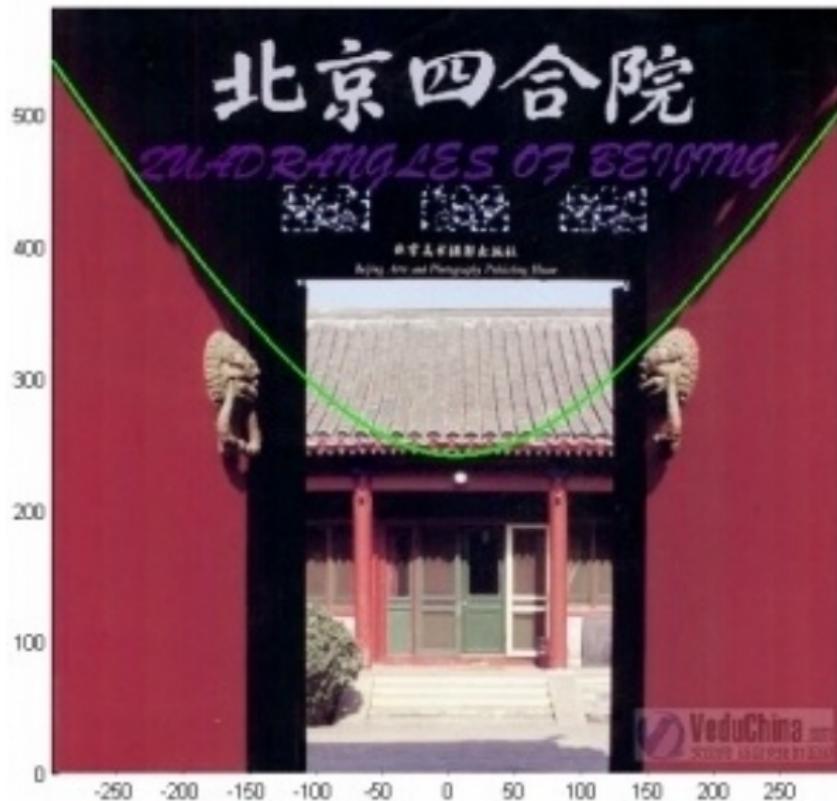
A Gentle Start

Consider the case $P(x, y) = (x^2, \sqrt{2}xy, y^2)$.



A Gentle Start

The green hyperplane corresponds to the green curve here:



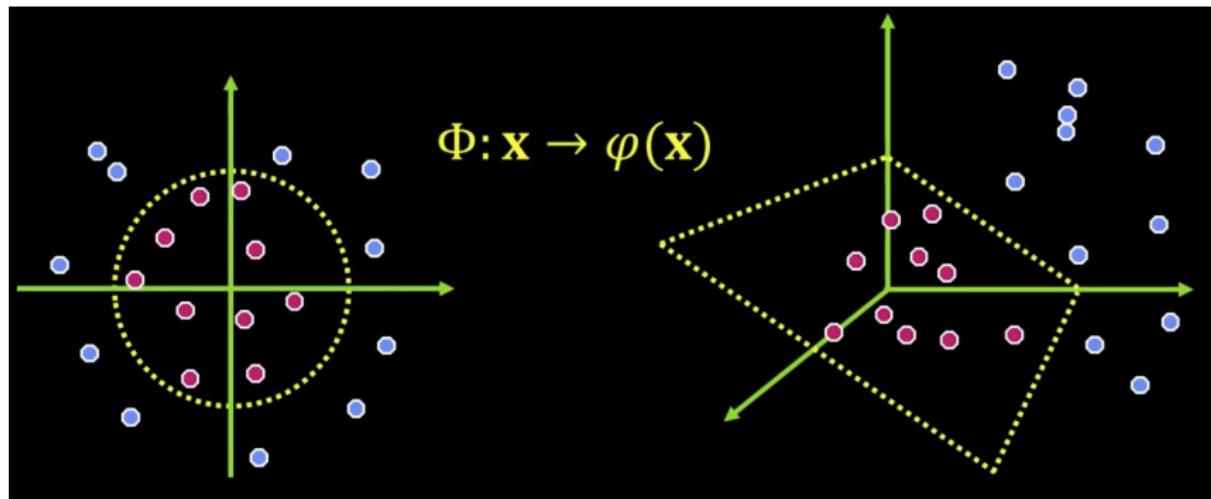
Kernel

In real-world, there are some datapoints they can't be separated by a hyperplane, even we leverage soft-margin in SVM. But as a theorem:
Almost all data points can be linearly separated in a (sufficiently) high dimension space.

Now let's define function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and recall soft margin SVM objective:

$$\begin{aligned} & \min \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t. } & y_i(w^T x_i + b) \geq 1 - \xi_i \quad (i = 1, 2, \dots, m) \\ & \xi_i \geq 0 \quad (i = 1, 2, \dots, m) \end{aligned} \tag{22}$$

Kernel



Kernel

In the new dimension space, we formulate the objective as:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t. } & y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i \quad (i = 1, 2, \dots, m) \\ & \xi_i \geq 0 \quad (i = 1, 2, \dots, m) \end{aligned} \tag{23}$$

the generalized Lagrangian is:

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(w^T \phi(x_i) + b) - 1 + \xi_i] - \sum_{i=1}^m \mu_i \xi_i \tag{24}$$

Optimization

For $\max_{\alpha_i \geq 0, \mu_i \geq 0} \min_{w, b, \xi} L(w, b, \alpha, \xi, \mu)$, now we can first optimize

$\min_{w, b, \xi} L(w, b, \alpha, \xi, \mu)$ by taking the derivative and set it to be zero:

$$\begin{aligned}\frac{\partial L}{\partial w} = 0 &\Rightarrow w = \sum_{i=1}^m \alpha_i y_i \phi(x_i) \\ \frac{\partial L}{\partial b} = 0 &\Rightarrow \sum_{i=1}^m \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \xi} = 0 &\Rightarrow C - \alpha_i - \mu_i = 0\end{aligned}\tag{25}$$

Optimization

$$\begin{aligned} & \underbrace{\max_{\alpha}}_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^m \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) \\ & \text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0 \\ & C - \alpha_i - \mu_i = 0 \\ & \alpha_i \geq 0, \mu_i \geq 0 \quad (i = 1, 2, \dots, m) \end{aligned} \tag{26}$$

if we define $K(x, z) = \phi(x)^T \phi(z)$, then the above equation is equivalent to:

$$\begin{aligned} & \underbrace{\min_{\alpha}}_{\alpha} \frac{1}{2} \sum_{i=1, j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \sum_{i=1}^m \alpha_i \\ & \text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned} \tag{27}$$

Polynomial Kernel

Now since $\phi(x) \in \mathbb{R}^p$, if p is very large, then it is computationally demanding. Then **Kernel Method** is proposed, such that instead of operating on high dimension $\phi(x) \in \mathbb{R}^p$, we can alternatively compute on original space \mathbb{R}^d . Recall $P(x, y) = (x^2, \sqrt{2}xy, y^2)$, we can verify:

$$\begin{aligned}\langle P(v_1), P(v_2) \rangle &= \langle (x_1^2, \sqrt{2}x_1y_1, y_1^2), (x_2^2, \sqrt{2}x_2y_2, y_2^2) \rangle \\&= x_1^2x_2^2 + 2x_1x_2y_1y_2 + y_1^2y_2^2 \\&= (x_1x_2 + y_1y_2)^2 \\&= \langle v_1, v_2 \rangle^2 \\&= K(v_1, v_2)\end{aligned}\tag{28}$$

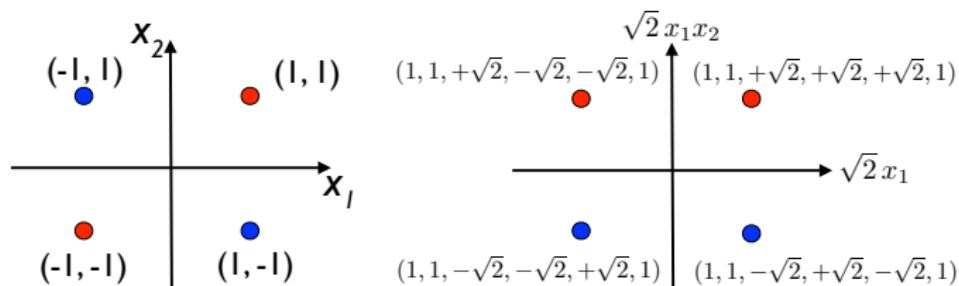
where we define *Polynomial Kernels* as:

$$\forall x, y \in \mathbb{R}^d, K(x, y) = (\langle x, y \rangle + c)^k, c > 0.$$

Polynomial Kernel ($d = 2, k = 2, c = 1$)

$$K(x, y) = (x_1 y_1 + x_2 y_2 + c)^2$$

$$= \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}cx_1 \\ \sqrt{2}cx_2 \\ c \end{bmatrix} \cdot \begin{bmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2}y_1y_2 \\ \sqrt{2}cy_1 \\ \sqrt{2}cy_2 \\ c \end{bmatrix} \quad (29)$$



Polynomial Kernel

$$\begin{aligned}\phi(\langle x_1, x_2, x_3 \rangle) \\ = & \langle x_1^3, x_2^3, x_3^3, \sqrt{3}x_1^2x_2, \sqrt{3}x_1^2x_3, \sqrt{3}x_2^2x_1, \sqrt{3}x_2^2x_3, \sqrt{3}x_3^2x_1, \sqrt{3}x_3^2x_2, \sqrt{6}x_1x_2x_3 \rangle \\ \phi(\langle y_1, y_2, y_3 \rangle) \\ = & \langle y_1^3, y_2^3, y_3^3, \sqrt{3}y_1^2y_2, \sqrt{3}y_1^2y_3, \sqrt{3}y_2^2y_1, \sqrt{3}y_2^2y_3, \sqrt{3}y_3^2y_1, \sqrt{3}y_3^2y_2, \sqrt{6}y_1y_2y_3 \rangle\end{aligned}\tag{30}$$

we can verify that $K(X, Y) = \phi(X)^T \phi(Y) = \langle X, Y \rangle^3$, which means we can reduce the computation from p to d .

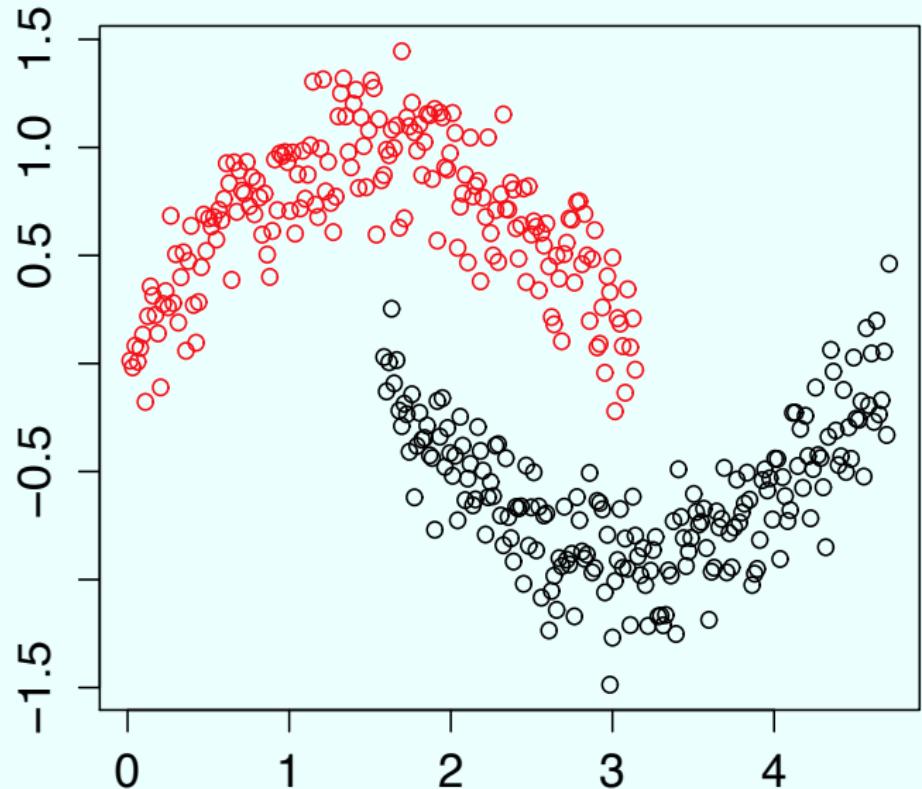
More broadly, the kernel $K(x, z) = (x^T z + c)^k$ (**Polynomial Kernel**) corresponds to a feature mapping to an $C(k + d, k)$ feature space. However, despite working in this $\mathcal{O}(d^k)$ -dimensional space, computing $K(x, z)$ still takes only $\mathcal{O}(d)$ time, and hence we never need to explicitly represent feature vectors in this very high dimensional feature space.

Gaussian Kernel

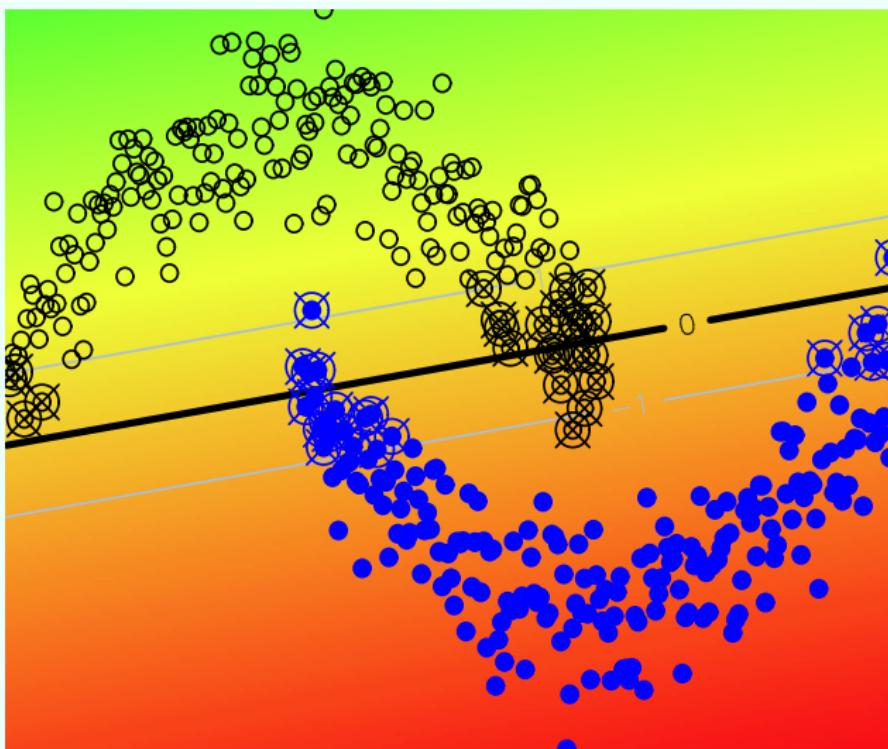
Radial Basis Function (RBF) is the Kernel used in libsvm by default:

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x} - \mathbf{y})^T(\mathbf{x} - \mathbf{y})\right) \\ &= \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x}^T \mathbf{x} + \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{x})\right) \\ &= \exp\left(-\frac{1}{2\sigma^2}(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)\right) \cdot \exp\left(\frac{1}{\sigma^2} \mathbf{y}^T \mathbf{x}\right) \quad (31) \\ &= C \cdot \exp\left(\frac{1}{\sigma^2} \mathbf{y}^T \mathbf{x}\right) \\ &= C \cdot \sum_{i=0}^{\infty} \frac{\left(\frac{1}{\sigma^2} \mathbf{y}^T \mathbf{x}\right)^i}{i!} \\ &= \sum_{i=0}^{\infty} \frac{C}{\sigma^{2i} i!} (\mathbf{y}^T \mathbf{x})^i \end{aligned}$$

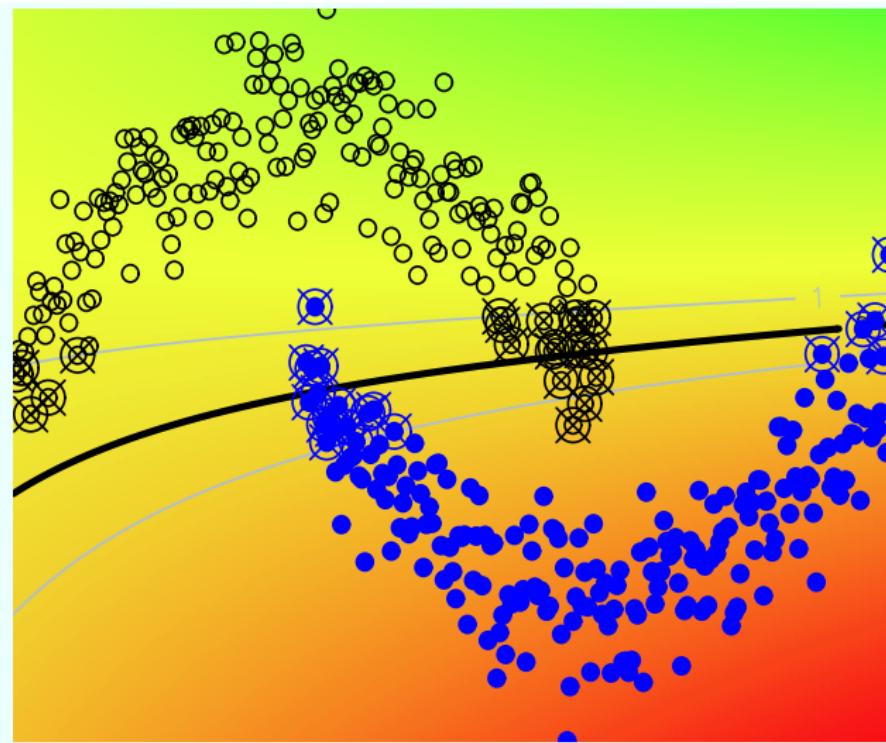
Applications of Kernel SVM



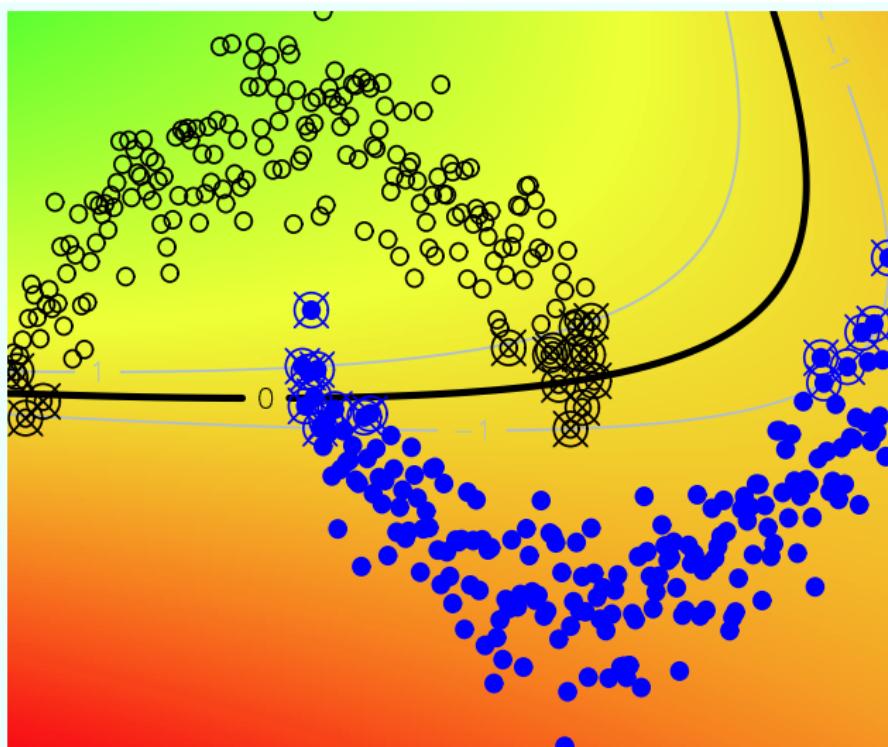
Linear Kernel SVM (46 SVs)



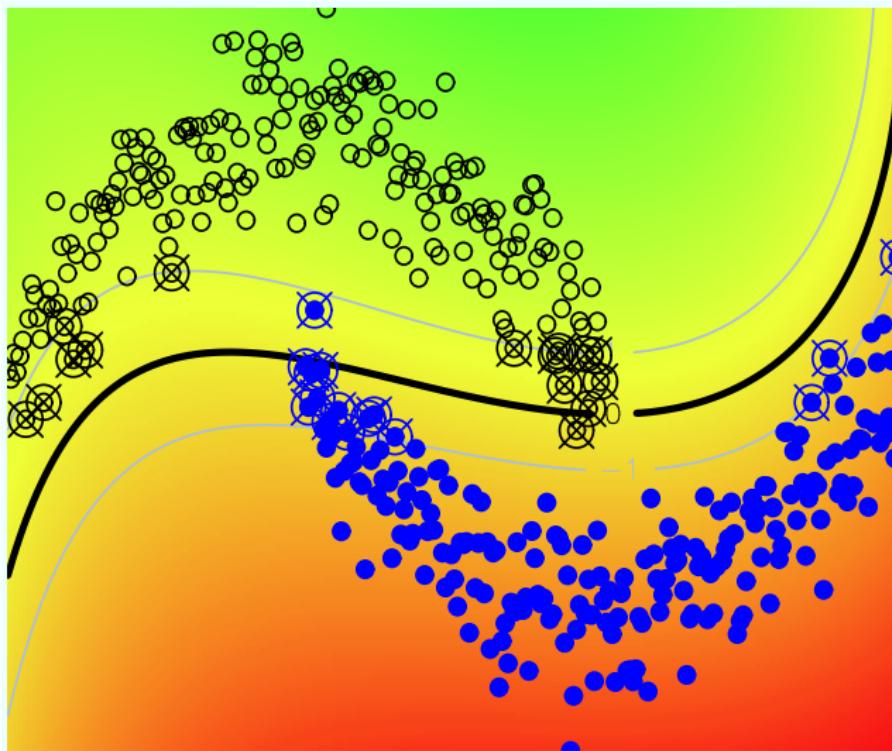
Kernel SVM (41 SVs, $k = 2$, $C = 1$)



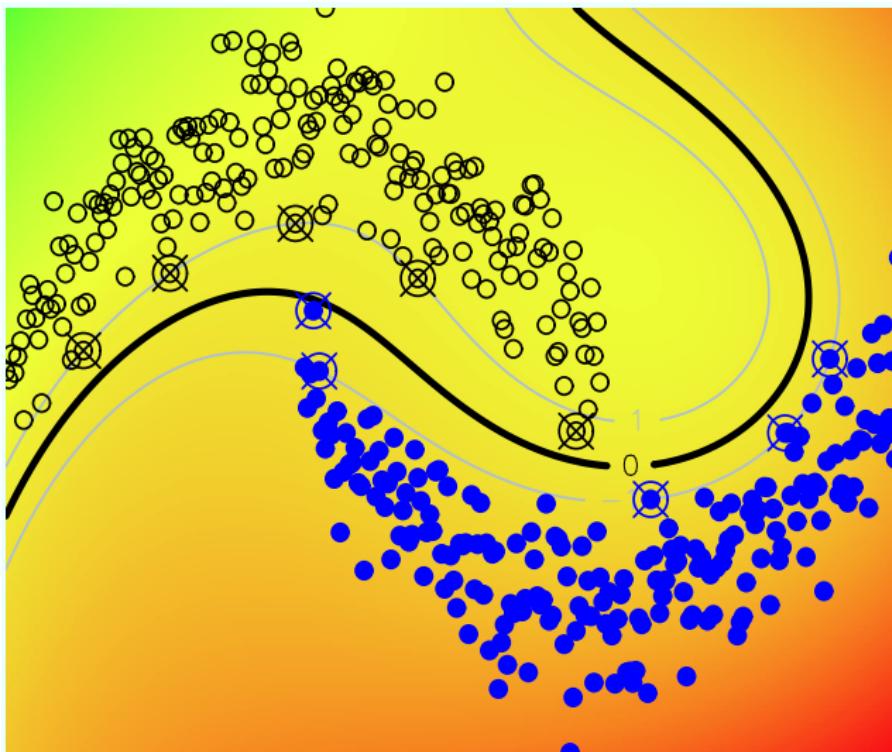
Kernel SVM (32 SVs, $k = 2$, $C = 50$)



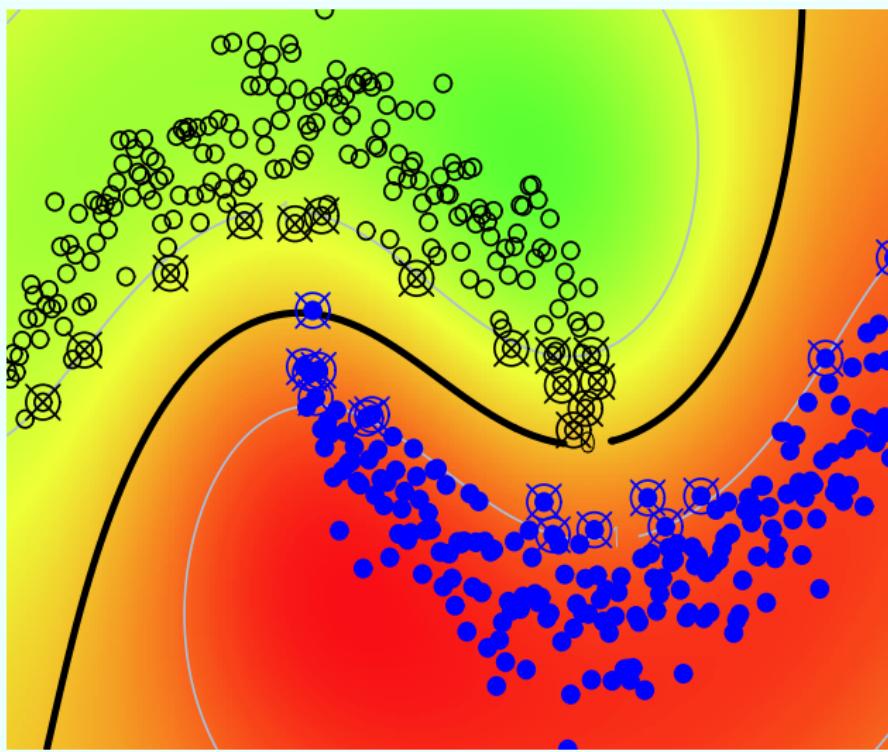
Kernel SVM (30 SVs, $k = 3$, $C = 1$)



Kernel SVM (10 SVs, $k = 3$, $C = 50$)



Kernel SVM (29 SVs, RBF, $\sigma = 1$, $C = 1$)



Kernel SVM (18 SVs, RBF, $\sigma = 0.5$, $C = 50$)

