

# Lecture 11: Support Vector Machine I

[SCS4049-02] Machine Learning and Data Science

---

Seongsik Park (s.park@dgu.edu)

AI Department, Dongguk University

# Tentative schedule

week	topic	date ( / )
1	Machine Learning Introduction & Basic Mathematics	09.02 / 09.07
2	Python Practice I & Regression	09.09 / 09.14
3	AI Department Seminar I & Clustering I	09.16 / 09.21
4	Clustering II & Classification I	09.23 / 09.28
5	Classification II	( ) / 10.05
6	Python Practice II & Support Vector Machine I	10.07 / 10.12
7	Support Vector Machine II & Ensemble Learning and Random Forest	10.14 / 10.19
8	( ) & <b>Mid-term exam</b>	10.21 / <b>10.26</b>
9	Neural networks	10.28 / 11.02
10	Backpropagation	11.04 / 11.09
11	Convolutional Neural Network	11.11 / 11.16
12	Model Optimization	11.18 / 11.23
13	Recurrent Neural network	11.25 / 11.30
14	Autoencoders	12.02 / 12.07
15	<b>Final exam</b>	( ) / <b>12.14</b>

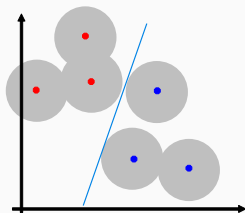
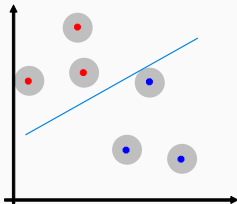
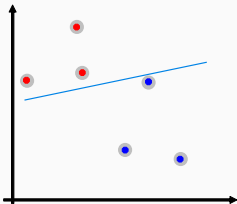
# SVM: Maximum Margin Classifier

---

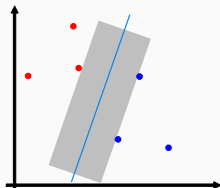
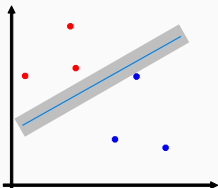
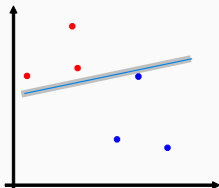
# What is a good decision boundary?

## ■ 데이터 노이즈에 대한 강건성 (Robustness)

- 노이즈(측정 오차)에 대해서 강건한 것이 좋은 모델이다.



■ 여유로운 것이 더 강건하다  $\Rightarrow$  넓은 통로가 좋다  $\Rightarrow$  Large Margin Classification



# What is a good decision boundary?

(support vectors)

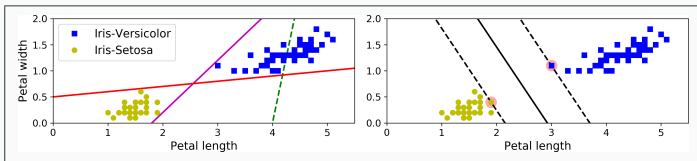


Figure 5-1. Large margin classification

Input feature scale support vector machine

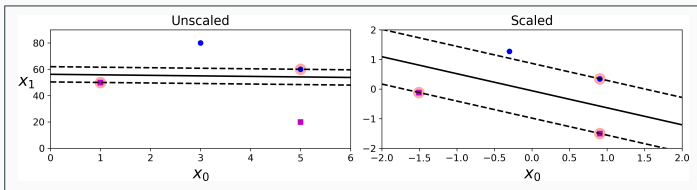


Figure 5-2. Sensitivity to feature scales

# Hard margin vs. soft margin

## Hard margin classification (hard-SVM)

- margin boundary
- (linearly separable)
- outlier

## Soft margin classification (soft-SVM)

- margin margin
- hyperparameter C: ( ), ( )

# Hard margin vs. soft margin

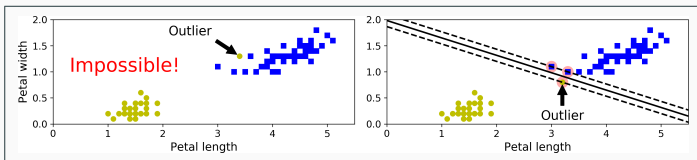


Figure 5-3. Hard margin sensitivity to outliers

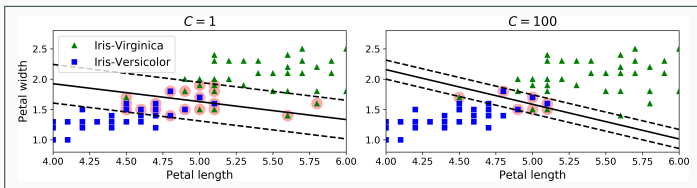


Figure 5-4. Large margin (left) versus fewer margin violations (right)

# A brief history of SVM

- SVM 1992 Boser, Guyon and Vapnik
- Statistical Learning Theory (Vapnik Chervonenkis)
- 
- SVM 1.1% Test error rate  $\approx$  (e.g., LeNet 4)
- 
- bioinformatics, text, image recognition
- 
- / † outlier detection
- /
- Kernel
- 
- Liblinear libsvm: Scikit-Learn liblinear libsvm



# Convex Optimization and Duality

---

## MoG: Lagrange multiplier

The method can be summarized as follows: in order to find the maximum or minimum of a function  $f(x)$  subjected to the equality constraint  $g(x) = 0$ , form the Lagrangian function

$$L(x; \lambda) = f(x) - \lambda g(x) \quad (1)$$

and find the stationary points of  $L$  considered as a function of  $x$  and the Lagrange multiplier  $\lambda$ . The solution corresponding to the original constrained optimization is always a saddle point of the Lagrangian function, which can be identified among the stationary points from the definiteness of the bordered Hessian matrix.

## MoG: Lagrange multiplier

Minimize  $f(x, y) = x^2 + y^2$  subject to the constraint  $x^2 + y^2 = 1$ , i.e.,

$$g(x, y) = x^2 + y^2 - 1 = 0 \quad (2)$$

Hence,

$$L(x, y, \lambda) = f(x, y) + \lambda g(x, y) = x^2 + y^2 + \lambda (x^2 + y^2 - 1) \quad (3)$$

Gradient

$$\nabla_{x,y,\lambda} L(x, y, \lambda) = (1 + 2\lambda x, 1 + 2\lambda y, x^2 + y^2 - 1) \quad (4)$$

and therefore,

$$\nabla_{x,y,\lambda} L(x, y, \lambda) = \begin{cases} 1 + 2\lambda x = 0 \\ 1 + 2\lambda y = 0 \\ x^2 + y^2 - 1 = 0 \end{cases} \quad (5)$$

## MoG: Lagrange multiplier

$$r ; ; L( ; ; ) \quad ( ) \quad \begin{matrix} \geq & 1 + 2 & = 0 \\ & 1 + 2 & = 0 \\ > & 2 + 2 & 1 = 0 \end{matrix} \quad (6)$$

This yields

$$= = \frac{1}{2}; \quad \notin 0 \quad (7)$$

$$\frac{1}{4} + \frac{1}{4} = 1 = 0 \quad (8)$$

So,

$$= \frac{1}{2} \quad (9)$$

which implies that the stationary points of  $L$  are

$$\frac{P}{2}; \frac{P}{2}; \frac{P}{2}; \quad ; \quad \frac{P}{2}; \quad \frac{P}{2}; \frac{P}{2} \quad (10)$$

# Optimization problem in standard form

$$\text{minimize } f_0(x) \tag{11}$$

$$\text{subject to } f_i(x) \leq 0; \quad i = 1, 2, \dots, m \tag{12}$$

$$h_j(x) = 0; \quad j = 1, 2, \dots, p \tag{13}$$

- $x \in \mathbb{R}^n$  is the optimization variable
- $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$  is the objective or cost function
- $f_i : \mathbb{R}^n \rightarrow \mathbb{R}; \quad i = 1, 2, \dots, m$  are the inequality constraint functions
- $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$  are the equality constraint functions

# Convex optimization problem

Standard form convex optimization problem

$$\text{minimize } f_0(x) \quad (14)$$

$$\text{subject to } f_i(x) \leq 0; \quad i = 1, 2, \dots, m \quad (15)$$

$$A_i x + b_i = c_i; \quad i = 1, 2, \dots, p \quad (16)$$

- $f_0, f_1, \dots, f_m$  are convex
- equality constraints are affine

Often written as

$$\text{minimize } f_0(x) \quad (17)$$

$$\text{subject to } f_i(x) \leq 0; \quad i = 1, 2, \dots, m \quad (18)$$

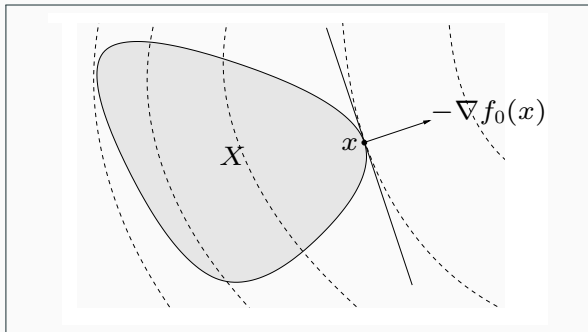
$$A_i x + b_i = c_i \quad (19)$$

Important property: feasible set of a convex optimization problem is convex

## Optimality criterion for differentiable $f_0$

$x^*$  is optimal if and only if it is feasible and

$$r_0(x^*) = \begin{cases} 0 & \text{if } x^* \text{ is optimal} \\ > 0 & \text{if } x^* \text{ is not optimal} \end{cases} \quad (20)$$



if nonzero,  $r_0(x)$  defines a supporting hyperplane to feasible set at  $x$

# Lagrangian

standard form problem

$$\text{minimize } f_0(x) \quad (21)$$

$$\text{subject to } f_i(x) \leq 0; \quad i = 1, 2, \dots, m \quad (22)$$

$$h_i(x) = 0; \quad i = 1, 2, \dots, p \quad (23)$$

variable  $x \in \mathbb{R}^n$ , domain  $D$ , optimal value

Lagrangian:  $\mathcal{L}(x; \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i (f_i(x) - 0) + \sum_{i=1}^p \mu_i (h_i(x) - 0)$

$$\mathcal{L}(x; \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i (f_i(x) - 0) + \sum_{i=1}^p \mu_i (h_i(x) - 0) \quad (24)$$

- weighted sum of objective and constraint functions
- $\lambda_i$  is Lagrange multiplier associated with  $f_i(x) \leq 0$
- $\mu_i$  is Lagrange multiplier associated with  $h_i(x) = 0$



# Lagrange dual function

Lagrange dual function:  $\mathcal{L} : \mathcal{R} \times \mathcal{R} \rightarrow \mathcal{R}$

$$\mathcal{L}(\lambda; \mathbf{b}) = \underset{\mathbf{z}}{\text{SH}} \mathcal{Q}(\mathbf{z}; \lambda; \mathbf{b}) \quad (25)$$

$$= \underset{\mathbf{z}}{\text{SH}} \left( \mathbf{z}^T \mathbf{A} \mathbf{z} + \sum_{i=1}^m \lambda_i (\mathbf{z}^T \mathbf{a}_i - b_i) \right) \quad (26)$$

$\mathcal{L}$  is concave, can be  $-\infty$  for some  $\lambda$ ;

lower bound property: if  $\lambda \geq \mathbf{0}$ , then  $\mathcal{L}(\lambda; \mathbf{b}) \leq \mathbf{0}$

proof: if  $\tilde{\mathbf{z}}$  is feasible and  $\lambda \geq \mathbf{0}$ , then

$$\mathbf{0}(\tilde{\mathbf{z}}) = \mathcal{Q}(\tilde{\mathbf{z}}; \lambda; \mathbf{b}) = \underset{\mathbf{z}}{\text{SH}} \mathcal{Q}(\mathbf{z}; \lambda; \mathbf{b}) = \mathcal{L}(\lambda; \mathbf{b}) \quad (27)$$

minimizing over all feasible  $\tilde{\mathbf{z}}$  gives  $\mathbf{0}(\tilde{\mathbf{z}}) \leq \mathcal{L}(\lambda; \mathbf{b})$

# The dual problem

Lagrange dual problem

$$\text{maximize} \quad ( ; ) \quad (28)$$

$$\text{subject to} \quad \mathbf{0} \quad (29)$$

- finds best lower bound on  $z^*$ , obtained from Lagrange dual function
- a convex optimization problem; optimal value denoted  $d^*$
- $( ; )$  are dual feasible if  $\mathbf{0}, ( ; ) \in \mathcal{D}$
- often simplified by making implicit constraint  $( ; ) \in \mathcal{D}$  explicit

# Weak and strong duality

weak duality:

- always holds (for convex and nonconvex problems)
- can be used to find nontrivial lower bounds for difficult problems

strong duality: =

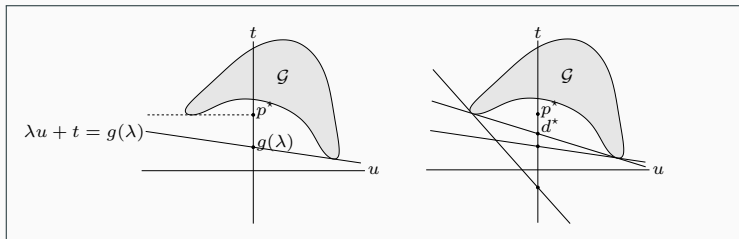
- does not hold in general
- holds for convex problems
- conditions that guarantee strong duality in convex problems are called constraint qualifications

# Geometric interpretation

for simplicity, consider problem with one constraint  $\lambda_1(\lambda) = 0$

interpretation of dual function

$$g(\lambda) = \inf_{(u,t) \in G} (\lambda u + t) \quad \text{where } G = \{(u,t) \mid f_1(u) \leq 0, g_0(u) \leq t\} \quad (30)$$



- $\lambda u + t = g(\lambda)$  is supporting hyperplane to  $G$
- hyperplane intersects  $u$ -axis at  $u = -g(\lambda)/\lambda$

# Geometric interpretation

## ■ Primal problem:

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } g(\mathbf{x}) \leq 0, \\ & \mathbf{x} \in \mathcal{X} \end{aligned}$$

where  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$

- $n = 2$  에 대해서, 집합  $\mathcal{G}$  를 다음과 같이 정의하자.

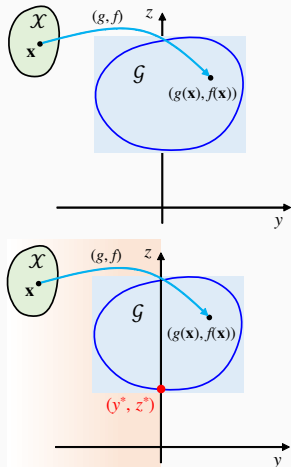
$$\mathcal{G} = \{(y, z) \mid y = g(\mathbf{x}), z = f(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathcal{X}\}$$

⇒ 그러면,  $\mathcal{G}$  는 사상  $(g, f)$  하에서  $\mathcal{X}$  의 치역(image)이다.

$$(g, f): \mathcal{X} \rightarrow \mathcal{G}$$

- 그러면 Primal solution은 최소 세로 좌표값  $z$  를 갖는  $y \leq 0$  인  $\mathcal{G}$  내의 점이다.

⇒ clearly  $(y^*, z^*)$



# Geometric interpretation

## ■ Lagrange Dual Problem

$$\text{maximize}_u \theta(u)$$

$$\text{subject to } u \geq 0$$

where (Lagrangian subproblem):

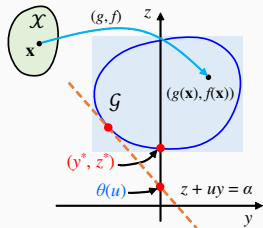
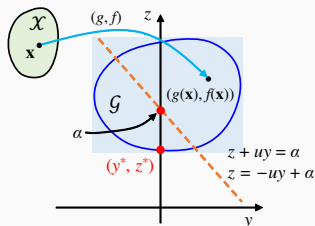
$$\theta(u) = \inf\{f(\mathbf{x}) + u g(\mathbf{x}) : \mathbf{x} \in \mathcal{X}\}.$$

1.  $u \geq 0$  일 때, Lagrangian dual subproblem 은 다음과 동등하다.

minimize  $z + uy$  over points  $(y, z)$  in  $\mathcal{G}$ ,

여기서  $z + uy = \alpha$  는 기울기가  $-u$  이고  $z$  축과 만나는 점이  $\alpha$  인 직선식이다.

2.  $\mathcal{G}$  에 대해서  $\theta(u) = z + uy$  를 최소화하기 위해서  $\mathcal{G}$  와의 접촉을 유지하면서 직선  $z + uy = \alpha$  를 평행 이동해 내려가면, 최후로 얻어지는  $z$ -축 절편값은 주어진  $u \geq 0$  에 대응하는  $\theta(u)$  값이다.



# Geometric interpretation

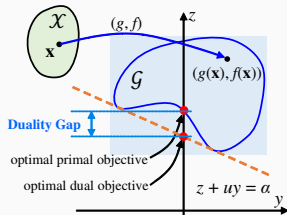
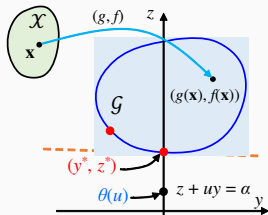
3. 마지막으로 dual problem 을 풀기 위해서는, 최후의 z-축 절편 값  $\theta(u)$  가 최대값이 되는 기울기  $-u$  ( $u \geq 0$ ) 를 찾아야 한다.

이러한 직선은 기울기가  $-u^*$  이고 점  $(y^*, z^*)$  에서 집합  $\mathcal{G}$  를 지지(support)한다.

그러므로, dual problem 의 해는  $-u^*$  이며 최적 dual objective value는  $z^*$  이다.

- Primal problem 과 dual problem 의 최적해가 동일한 경우, duality gap 이 없다고 말한다 (**strong duality**).
- 동일하지 않은 경우는 duality gap 이 존재한다고 말한다 (**weak duality**).
- 적절한 convexity condition 들이 충족되면 primal 과 dual 최적화 문제에 duality gap 이 없다.
- [우측 그림] 집합  $\mathcal{G}$  의 nonconvexity 로 인한 Duality Gap

$$f(\mathbf{x}) \geq \theta(\mathbf{u})$$



# Karush-Kuhn-Tucker (KKT) conditions

the following four conditions are called KKT conditions (for a problem with differentiable  $f, g_i$ )

1. primal constraints:  $g_i(x) \leq 0, i = 1, \dots, m, f(x) = 0, i = 1, \dots, p$
2. dual constraints:  $\lambda_i \geq 0$
3. complementary slackness:  $\lambda_i g_i(x) = 0; i = 1, \dots, m$
4. gradient of Lagrangian with respect to  $x$  vanishes:

$$\nabla_x f(x) + \sum_{i=1}^m \lambda_i \nabla_x g_i(x) + \sum_{i=1}^p \mu_i \nabla_x h_i(x) = 0 \quad (31)$$

if strong duality holds and  $x^*, \lambda^*, \mu^*$  are optimal, then they must satisfy the KKT conditions



## SVM: Theory

---

# Maximum margin classifier

We begin our discussion of support vector machines to the two-class classification problem using linear models of the form

$$f(\mathbf{x}) = \dots + b + \mathbf{w} \cdot \phi(\mathbf{x}) \quad (32)$$

where  $\phi(\mathbf{x})$  denotes a fixed feature-space transformation, and we have made the bias parameter  $b$  explicit.

The training data set comprises  $n$  input vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , with corresponding target values  $y_1, y_2, \dots, y_n$  where  $y_i \in \{-1, 1\}$ , and new data points  $\mathbf{x}$  are classified according to the sign of  $f(\mathbf{x})$ .

# Maximum margin classifier

We shall assume that the training data set is linearly separable in feature space, so that by definition there exists at least one choice of the parameters  $w$  and  $b$  such that a function satisfies  $f(x) > 0$  for points having  $y = +1$  and  $f(x) < 0$  for points having  $y = -1$ , so that  $y f(x) > 0$  for all training data points.

## Maximum margin classifier: optimality criterion

Thus the distance of a point  $\mathbf{x}$  to the decision surface is given by

$$\frac{(\mathbf{x})}{k..k} = \frac{(\dots (\mathbf{x}) + )}{k..k}; \quad (33)$$

The margin is given by the perpendicular distance to the closest point  $\mathbf{x}$  from the data set, and we wish to optimize the parameters ...and in order to maximize this distance. Thus the maximum margin solution is found by solving

$$-\mathbf{q} \setminus - \mathbf{x} \quad \frac{1}{k..k} \setminus \mathcal{S} [ (\dots (\mathbf{x}) + ) ] \quad (34)$$

where we have taken the factor  $1/k..k$  outside the optimization over because ...does not depend on .

## Dual problem for convex optimization

Direct solution of this optimization problem would be very complex, so we shall convert it into an equivalent problem that is much easier to solve.

## Lagrangian function with constraint

In order to solve this constrained optimization problem, we introduce Lagrange multipliers  $\lambda_1, \dots, \lambda_m$ , with one multiplier  $\lambda_i$  for each of the constraints, giving the Lagrangian function

$$\mathcal{L}(\mathbf{x}; \lambda) = \frac{1}{2} \mathbf{x}^T \mathbf{K} \mathbf{x} - \sum_{i=1}^m \lambda_i (f_i(\mathbf{x}) - c_i) \quad (35)$$

where  $\mathbf{x} = (x_1, \dots, x_n)^T$ . Note the minus sign in front of the Lagrange multiplier term, because we are minimizing with respect to  $\mathbf{x}$  and  $\lambda_i$ , and maximizing with respect to  $\lambda_i$ .

Setting the derivatives of  $\mathcal{L}(\mathbf{x}; \lambda)$  with respect to  $\mathbf{x}$  and  $\lambda_i$  equal to zero, we obtain the following two conditions

$$\mathbf{0} = \mathbf{K} \mathbf{x} - \sum_{i=1}^m \lambda_i \nabla f_i(\mathbf{x}) \quad (36)$$

$$\mathbf{0} = \sum_{i=1}^m \lambda_i \nabla f_i(\mathbf{x}) - \mathbf{g} \quad (37)$$

# Lagrangian function with constraint

Eliminating ...and from  $\mathcal{Q}(\dots; \dots)$  using these conditions then gives the ... of the maximum margin problem in which we maximize

$$\mathcal{Q}(\dots) = \sum_{i=1}^X \dots - \frac{1}{2} \sum_{i=1}^X \sum_{j=1}^X \dots (\dots; \dots) \quad (38)$$

with respect to  $\dots$  subject to the constraints

$$\dots = 0; \dots = 1; \dots; \dots^1 \quad (39)$$

$$\sum_{i=1}^X \dots = 0 \quad (40)$$

Here the kernel function is defined by  $(\dots; \dots^\theta) = (\dots) (\dots)^\theta$ .

## Prediction for a new sample: support vector machine

In order to classify new data points using the trained model, we evaluate the sign of  $f(\ddagger)$ . This can be expressed in terms of the parameter  $f^* g$  and the kernel function by substituting for ... to give

$$f(\ddagger) = \sum_{i=1}^n \alpha_i (f^*; \ddagger) + \dots \quad (41)$$



## KKT condition: complementary slackness

We show that a constrained optimization of this form satisfies the KKT conditions, which in this case require that the following three properties hold

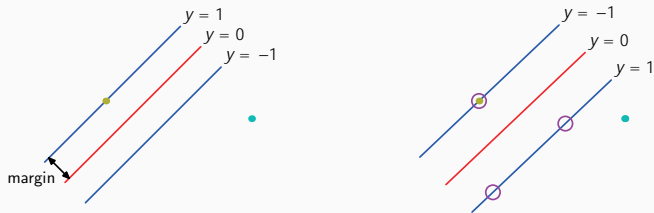
$$g(x) = 0 \quad (42)$$

$$\alpha_i g(x_i) = 0 \quad (43)$$

$$f(x) - g(x) = 0 \quad (44)$$

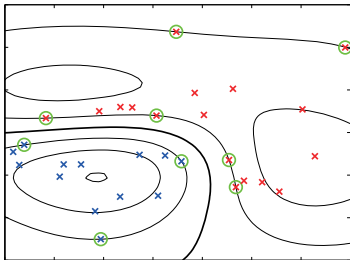
Thus for every data point, either  $g(x) = 0$  or  $\alpha_i = 0$ . Any data point for which  $g(x) = 0$  will not appear in the sum and hence plays no role in making predictions for new data points. The remaining data points are called **support vectors**, and because they satisfy  $\alpha_i = 1$ , they correspond to points that lie on the maximum margin hyperplanes in feature space.

# KKT condition: complementary slackness



**Figure 7.1** The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points, as shown on the left figure. Maximizing the margin leads to a particular choice of decision boundary, as shown on the right. The location of this boundary is determined by a subset of the data points, known as support vectors, which are indicated by the circles.

**Figure 7.2** Example of synthetic data from two classes in two dimensions showing contours of constant  $y(\mathbf{x})$  obtained from a support vector machine having a Gaussian kernel function. Also shown are the decision boundary, the margin boundaries, and the support vectors.



# Appendix

---

## Reference and further reading

- “Chap 7 | Sparse Kernel Machines” of C. Bishop, Pattern Recognition and Machine Learning
- “Chap 5 | Support Vector Machines” of A. Geron, Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow
- “Chap 4 | Convex Optimization Problems”, “Chap 5 | Duality” of S. Boyd, Convex Optimization
- “Lecture 6 | Support Vector Machines” of Kwang Il Kim, Machine Learning (2019)