'margin'

Classifier: max. margin classifier = SVM ← convex.

Primal problem
→ dual problem → Solution

comple. slackness

# Inclass 09: Support Vector Machine

[SCS4049] Machine Learning and Data Science

Seongsik Park (s.park@dgu.edu)

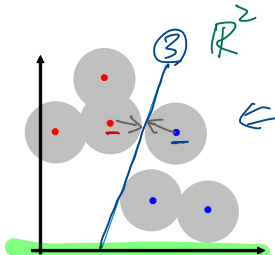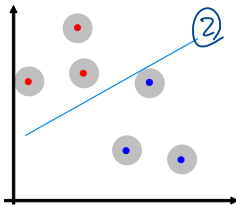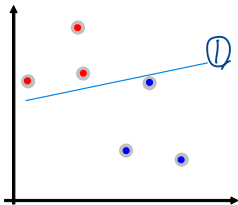School of AI Convergence & Department of Artificial Intelligence, Dongguk University

# Maximum Margin Classifier

■ 데이터 노이즈에 대한 강건성 (Robustness)
  ● 노이즈(측정 오차)에 대해서 강건한 것이 좋은 모델이다.

여유로운 것이 더 강건하다 ⇒ 넓은 통로가 좋다 ⇒ Large Margin Classification

margin
= boundary부터
가장가까운
샘플까지의
거리.

의사 결정은 경계의 데이터(support vectors)에 의해서 결정됨

*샘플의 일부, boundary를 형성, 예측할때 사용*



*Figure 5-1. Large margin classification*

*거리 기반, margin*

*normalize 필요.*

Input feature의 scale에 민감한 support vector machine



*Figure 5-2. Sensitivity to feature scales*

Hard margin classification (hard-SVM) ← 구속조건

- 모든 데이터들이 <u>margin 밖에 위치하도록</u> boundary를 설정
- 데이터가 <u>linearly separable</u>할 때만 적용 가능
- outlier에 매우 민감

정해진 margin안으로
못들어가는거
hard-SVM

Soft margin classification (soft-SVM)

- margin을 가능한 넓게 하면서도 margin 안에 들어오는 것도 허용
- hyperparameter C: 클수록 좁아짐 (엄격), 작을 수록 넓어짐 (허용)

Figure 5-3. Hard margin sensitivity to outliers



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

# A brief history of SVM

- SVM은 1992년 Boser, Guyon and Vapnik에 의해서 소개됨     *convex opt.*
- Statistical Learning Theory에 이론적 바탕을 둔 알고리즘
- 손글씨 숫자 인식에서 뛰어난 성능을 보이면서 널리 쓰이게 됨
- SVM으로 1.1% Test error rate ≈ 신중히 설계된 신경망과 맞먹음
- 실용적으로 우수한 성능
- Bioinformatics, text, image recognition 등 많은 성공 사례
- 선형/비선형 분류 뿐 아니라 회귀, outlier detection 도 수행
- 복잡한 소규모/중규모 데이터셋의 분류에 특히 잘 맞음
- Kernel 방법을 사용하는 대표적 알고리즘
- Liblinear & libsvm: Scikit-Learn 에서 liblinear 및 libsvm 을 사용

어떤 아이디어. → 수학적, 최적화.

maximum margin → 최적화 정의
$\downarrow$ convex

Support vector ← Solution
Complementary Slackness

# Hard SVM

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

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \tag{1}$$

↳nonlinear 현재는 가능.

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

The training data set comprises $N$ input vectors $\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N$, with corresponding target values $t_1, t_2, ..., t_N$ where $t_n \in \{-1, 1\}$, and new data points $\mathbf{x}$ are classified according to the sign of $y(\mathbf{x})$

$$y(x) > 0 \longrightarrow t = 1$$
$$y(x) < 0 \longrightarrow t = -1$$

$$y(x) = w^T x + b \leftarrow$$
$$= \theta^T x + \theta_0$$



$x_n$

거리
(양수)

$$\frac{|y(x_n)|}{||w||_2} = \frac{|w^T x_n + b|}{||w||_2}$$

$$\frac{|b|}{||w||_2}$$

$y(x) > 0$

$x \in \mathbb{R}^2$

$y(x) < 0$

$y(x) = 0$

학습데이터

입력 $x_1, x_2, x_3, \cdots, x_N$

출력 $t_1, t_2, t_3, \cdots, t_N$

boundary까지의 거리

$$\frac{|w^T x_1 + b|}{||w||}, \quad \frac{|w^T x_2 + b|}{||w||}, \quad \frac{|w^T x_3 + b|}{||w||}, \quad \cdots$$

margin

$$\min_{n \in} \left\{ \frac{|w^T x_n + b|}{||w||_2} \right\} = \frac{t_n y(x_n)}{||w||_2} \quad \oplus$$

$|y(x_n)|$

제일 가까운애.

margin

$$\max_{w,b} \left( \min_n \frac{t_n y(x_n)}{\|w\|} = \frac{t_n (w^T x_n + b)}{\|w\|} \right)$$

max margin

margin

$$\max_{w,b} \left\{ \frac{1}{\|w\|} \min_n t_n y(x_n) \right\}$$

$$\Rightarrow \text{objective}$$

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 $\mathbf{w}$ and $b$ such that a function satisfies $y(\mathbf{x}_n) > 0$ for points having $t_n = +1$ and $y(\mathbf{x}_n) < 0$ for points having $t_n = -1$, so that $t_n y(\mathbf{x}_n) > 0$ for all training data points.

$$t_n \, y(x_n) > 0 \qquad \Longrightarrow \qquad \text{모든 sample이 제대로 분류.}$$
$$n = 1, 2, \cdots, N \qquad \qquad (\text{linearly separable})$$

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

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n(\mathbf{w}^T\mathbf{x}_n + b)}{\|\mathbf{w}\|}. \tag{2}$$

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

$$\arg\max_{\mathbf{w},b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^T\mathbf{x}_n + b)] \right\} \tag{3}$$

where we have taken the factor $1/\|\mathbf{w}\|$ outside the optimization over $n$ because $\mathbf{w}$ does not depend on $n$.

Primal optimization problem for hard SVM

*objective* (handwritten)

margin을 가능한 (handwritten)
최대로 만들자. (handwritten)

$$\text{maximize} \quad \frac{1}{\|\mathbf{w}\|} \min[t_n(\mathbf{w}^T\mathbf{x}_n + b)] \qquad (4)$$

*constraint* (handwritten)

$$\text{subject to} \quad t_n(\mathbf{w}^T\mathbf{x}_n + b) \geq 0 \qquad (5)$$

$$(6)$$

→ 모든 sample을 지키는 판별함수 (handwritten)

Equivalently,

$$\left( \begin{array}{ll} \text{minimize} & \frac{1}{2}\|\mathbf{w}\|^2 \quad \textit{objective} \qquad (7) \\ \text{subject to} & t_n(\mathbf{w}^T\mathbf{x}_n + b) \geq 1 \quad \textit{constraint} \qquad (8) \end{array} \right)$$

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.

$$\min_{w,b} \quad \frac{1}{2} \|w\|^2$$

$$\mathcal{L}(w, b, a_1, a_2, a_3, \cdots, a_N)$$

$$\text{s.t.} \quad t_1(w^T x_1 + b) \geq 1 \leftarrow \boxed{a_1}$$

$$t_2(w^T x_1 + b) \geq 1 \leftarrow a_2$$

$$\vdots$$

$$= \frac{1}{2}\|w\|^2$$

$$+ \boxed{a_1}(1 - t_1(w^T x_1 + b))$$

$$t_N(w^T x_N + b) \geq 1 \leftarrow a_N$$

$$+ \boxed{a_2}(1 - t_2(w^T x_2 + b))$$

$$\vdots$$

$$f_i(\;) \leq 0.$$

$$1 - t_n(w^T x_n + b) \leq 0.$$

$$\max \frac{1}{||w||} \quad \min \quad t_n(w^T x_n + b)$$

$$\text{s.t.} \qquad t_n(w^T x_n + b) \geq 0.$$

$$y(x) = w^T x + b$$

Scalar $\theta n$.

① $w = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b = 1$

$y(x) = 1 + x_1 + x_2 = 0$

② $w = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \quad b = 2$

$y(x) = 2 + 2x_1 + 2x_2 = 0$

$X_n$ 거리

$$\Rightarrow \frac{1}{\|w\|_2} \, t_n(w^T x_n + b)$$

$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \, b = 1$

$w = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \, b = 2$

① $\frac{1}{\sqrt{3}}$ $t_n(1 + x_1 + x_2)$

$\times 2$

② $\frac{1}{2\sqrt{3}}$ $t_n(2 + 2x_1 + 2x_2)$

$\times 2$

$$\max \quad \frac{1}{||w||} \quad \min \quad t_n(w^T x_n + b) \to 1$$

$$\text{s.t.} \quad t_n(w^T x_n + b) \geq \min t_n(w^T x_n + b) \geq 0$$
$$\hookrightarrow 1$$

$$\Rightarrow \quad \max \quad \frac{1}{||w||} \, 0 \Rightarrow \quad \min \frac{1}{2} ||w||_2^2$$

$$\text{s.t.} \quad t_n(w^T x_n + b) \geq 1$$

Setting the derivatives of $L(\mathbf{w}, b, \mathbf{a})$ with respect to $\mathbf{w}$ and $b$
equal to zero, we obtain the following two conditions

$$\mathbf{w} = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n \qquad = \qquad \begin{array}{l} a_1 t_1 x_1 \\ + a_2 t_2 x_2 \\ + a_3 t_3 x_3 \\ + \\ \vdots \end{array} \qquad (9)$$

$$0 = \sum_{n=1}^{N} a_n t_n \qquad (10)$$

① $a_n = 0$    $t_n(\mathbf{w}^T \mathbf{x}_n + b) > 1$ $\Rightarrow$ 얘 $x_n$은 margin위가 아니고
                                                           $\mathbf{w}$ 영역이 안겹침.

② $a_n > 0$    $t_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$ $\Rightarrow$ 이 $x_n$은 margin을 만든 개

                                                  $\downarrow$

                                                  Support rector.

10/18

## Lagrangian function with constraint

Eliminating $\mathbf{w}$ and $b$ from $L(\mathbf{w}, b, \mathbf{a})$ using these conditions then gives the *dual representation* of the maximum margin problem in which we maximize

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \tag{11}$$

with respect to $\mathbf{a}$ subject to the constraints

$$a_n \geq 0, \qquad n = 1, ..., N \tag{12}$$

$$\sum_{n=1}^{N} a_n t_n = 0. \tag{13}$$

Here the kernel function is defined by $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$.

In order to classify new data points using the trained model, we evaluate the sign of $y(\mathbf{x})$. This can be expressed in terms of the parameter $\{a_n\}$ and the kernel function by substituting for $\mathbf{w}$ to give

$$\mathbf{w}^T\mathbf{x} + b = \sum_n a_n t_n \mathbf{x}_n^T\mathbf{x} + b$$

$$> 0 \qquad \text{득베터} $$

$$\text{입이 버젼에?}$$

$$y(\mathbf{x}) = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n^T\mathbf{x} + b \tag{14}$$

$$= \sum_{n=1}^{N} a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \tag{15}$$

$$\mathbf{x}^T\mathbf{x}_1 \quad a_1 > 0 \qquad a_1 t_1 \mathbf{x}^T\mathbf{x}_1 \quad \Leftarrow$$

$$\mathbf{x} \rightarrow \quad \mathbf{x}^T\mathbf{x}_n \quad a_n > 0 \qquad + a_n t_n \mathbf{x}^T\mathbf{x}_n \quad \Leftarrow$$

$$\mathbf{x}^T\mathbf{x}_{10} \quad a_{10} > 0 \qquad + a_{10} t_{10} \mathbf{x}^T\mathbf{x}_{10} \quad \Leftarrow$$

We show that a constrained optimization of this form satisfies the *Karush-Kuhn-Tucker* (KKT) conditions, which in this case require that the following three properties hold

$$a_n \geq 0 \qquad \text{Dual const.} \tag{16}$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0 \qquad \text{Primal const} \tag{17}$$

$$a_n \{ t_n y(\mathbf{x}_n) - 1 \} = 0 \qquad \text{com. Slackness} \tag{18}$$

Thus for every data point, either $a_n = 0$ or $t_n y(\mathbf{x}_n) = 1$. Any data point for which $a_n = 0$ 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 $t_n y(\mathbf{x}_n) = 1$, they correspond to points that lie on the maximum margin hyperplanes in feature space.
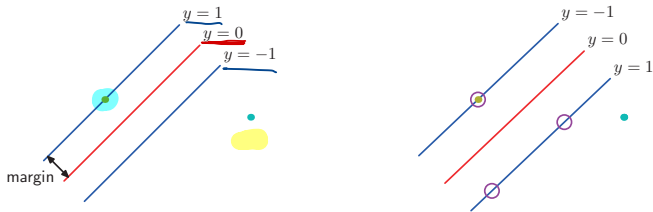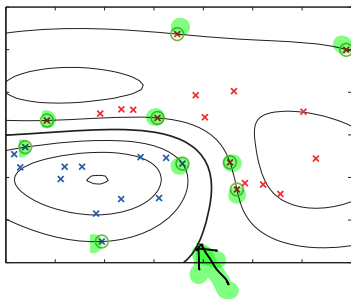
**Figure 7.1** The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points, as shown on the left figure. Maximizing the margin leads to a particular choice of decision boundary, as shown on the right. The location of this boundary is determined 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.

장점

_Convex opt._

- 강력하고 우수한 이론을 바탕으로 함
- 많은 블랙박스 알고리즘과는 대조적으로
  비교적 직관적인 해석과 이해가 가능
- 학습이 상대적으로 쉬움 ← convex opt → global 방법.
- 신경망처럼 지역 최적값에 빠지는 일이 없음
- 학습 시간이 차원에 의존하지 않으며
  kernel trick 덕분에 고정된 입력에만 의존
- 과적합이 잘 조절되는 경향
- 많은 분야에서 신경망 및 기타 알고리즘과 필적하는 성능 ∝
- 데이터가 작은 조건이나 고차원 공간에서도 잘 일반화

단점

구속조건 : boundary, input.

- 노이즈에 민감      hard SVM, Soft SVM
- 비교적 적은 수의 잘못된 label로 성능이 심각하게 악화
- 커널 함수를 선택하는 방법에 대한 정리된 원칙이 없음
- hyperparameter $C$의 적정값을 정하기 위한 원칙이 없음
- 컴퓨터의 메모리와 계산 시간 측면에서 비용이 높은 편이며
  multiclass에서 더욱더 심화됨

① convex opt. 풀어야.
② 예측. 데이터개수↑

# 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)