

12. Support Vector Machines



12.1 Optimization objective

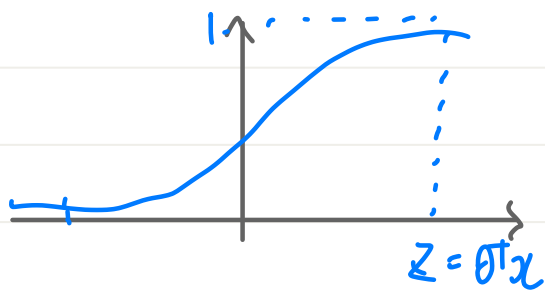
① Alternative view of logistic regression

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

한 개의 sample에 cost

$$\Rightarrow -(y \log h_{\theta}(x) + (1-y) \log (1-h_{\theta}(x)))$$

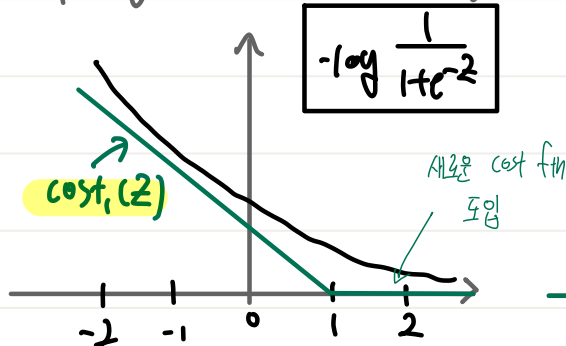
$$= -y \log \frac{1}{1 + e^{\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1 + e^{\theta^T x}}\right)$$



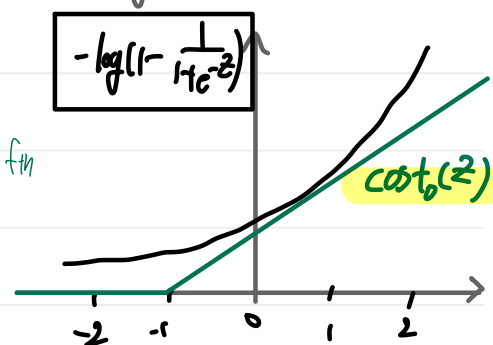
If $y=1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x > 0$

If $y=0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x < 0$

If $y=1$ (want $\theta^T x > 0$)



If $y=0$ (want $\theta^T x < 0$)



② Support vector machine

Logistic regression :

$$\min_{\theta} \frac{1}{n} \left[\sum_{i=1}^n y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1-y^{(i)}) (-\log (1-h_{\theta}(x^{(i)}))) \right] + \frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2$$

Support vector machine :

$$\min_{\theta} \frac{1}{n} \left[\sum_{i=1}^n y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2$$

편의를 위해 추가한 항임

→ 약제해도 상의 최적값 찾는데 영향X

$$\begin{aligned} & A + \lambda B \rightarrow \lambda \text{가 높으면 } B \text{에 큰 가중치} \\ \Rightarrow & CA + B \quad (C = \frac{1}{\lambda}) \rightarrow \lambda \text{가 높으면 } A \text{가 작아져 줄여 } B \text{에 큰 가중치} \end{aligned}$$

⇒ 두 사이 골라
않지만 같은
최적값 0
문은 < 같음

$$\Rightarrow \min_{\theta} C \sum_{i=1}^n [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

hypothesis : $h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{o.w} \end{cases}$ → SVM의 hypothesis는 y가 1이거나 0일 때 1 또는 0이다.

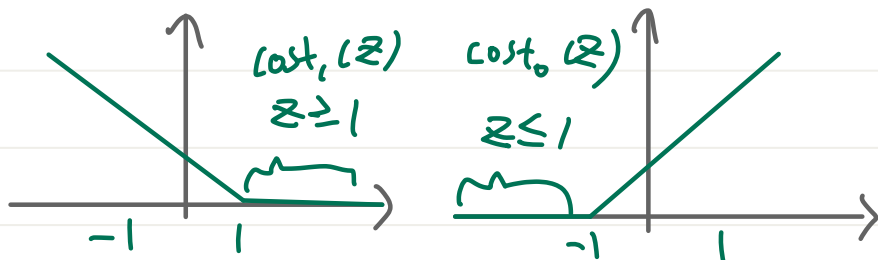
12. Support Vector Machines



12.2 Large Margin Intuition

① SVM Decision Boundary

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{i=1}^n \theta_i^2$$



If $y=1$, we want $\theta^T x \geq 1$ (not ≥ 0)

If $y=0$, we want $\theta^T x \leq -1$ (not < 0)

whenever $y^{(i)}=1$: $\theta^T x^{(i)} \geq 1$

whenever $y^{(i)}=0$: $\theta^T x^{(i)} \leq -1$

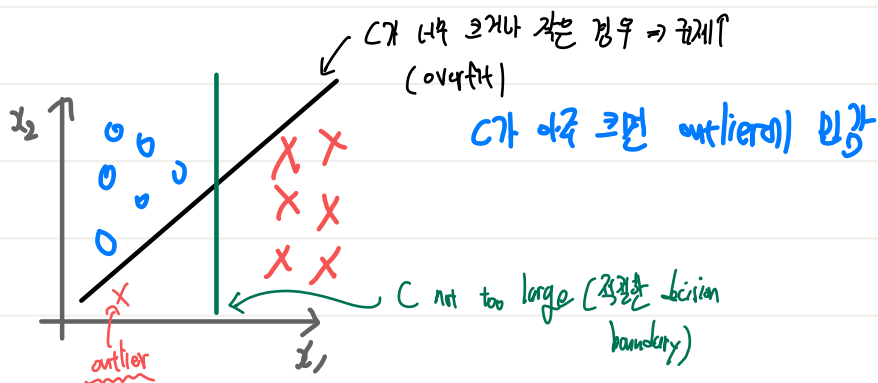
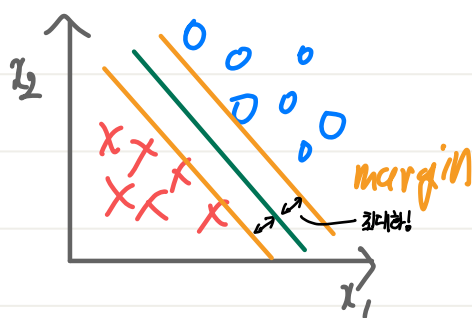
$$\Rightarrow \min_{\theta} \cancel{C \sum_{i=1}^m} + \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

↪ zero가 됨

$$\Rightarrow \min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 \quad \left(\begin{array}{ll} \theta^T x^{(i)} \geq 1 & \text{if } y^{(i)}=1 \\ \theta^T x^{(i)} \leq -1 & \text{if } y^{(i)}=0 \end{array} \right)$$

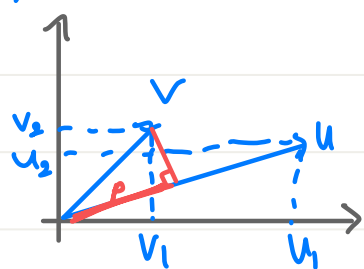
< linearly separable case >

⇒ large margin classifier



12.3 The mathematics behind large margin classification

① Vector Inner Product



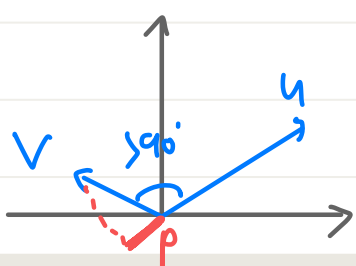
$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$\|u\|$ = length of vector u

$$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

p = length of projection of v onto u

$$u^T v = p \cdot \|u\| = u_1 v_1 + u_2 v_2 = v^T u$$



$$\rightarrow u^T v = p \cdot \|u\| \quad (p < 0)$$

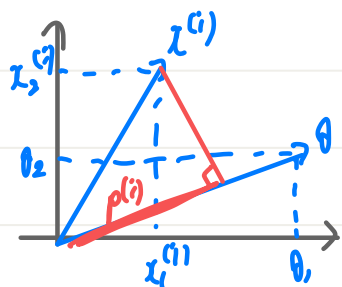
12. Support Vector Machines



② SVM Decision Boundary

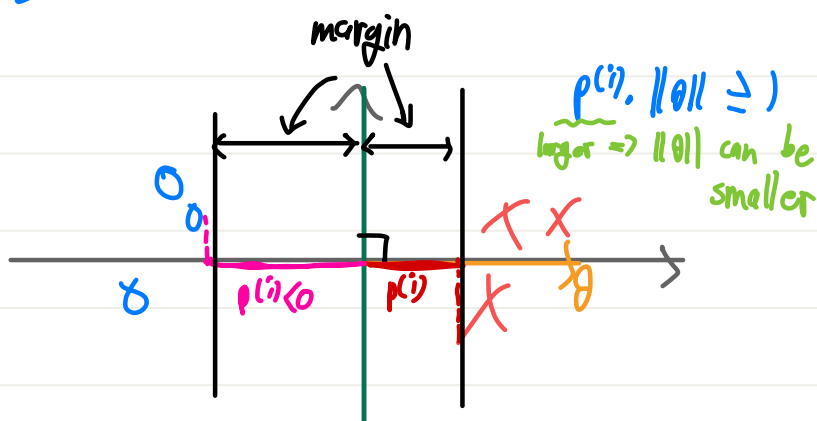
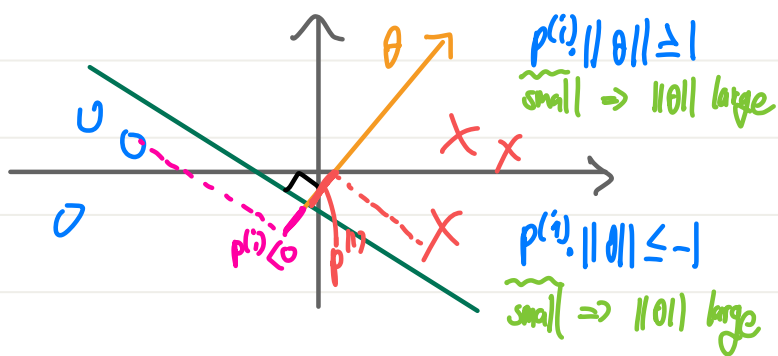
$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$

Set, $\theta^T x^{(i)} \geq 1$ if $y^{(i)} = 1$ (일반하게 $\theta_0 = 0, n=2$)
 $\theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$ ($p^{(i)}$ 는 $x^{(i)}$ 에 θ 의 수직거리)



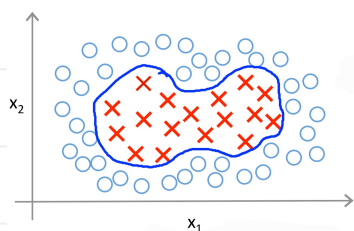
$$\theta^T x^{(i)} = p^{(i)} \cdot \|\theta\|$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$



12.4 Kernels 1

① Non-linear Decision Boundary



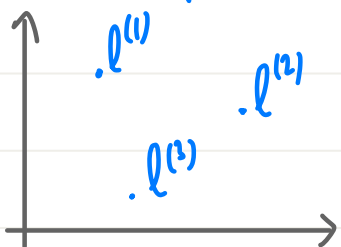
predict $y=1$ if

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$$

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0 & \text{o.w.} \end{cases} \Rightarrow f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2$$

$$f_4 = x_1^2, f_5 = x_2^2, \dots$$

② Kernel



주어진 x 에서 A 로 feature를 계산한다.

landmarks $l^{(1)}, l^{(2)}, l^{(3)}$ 의 근접성이 각각의 점

Given x :

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp(\dots)$$

kernel (Gaussian kernel)

ex)

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{i=1}^n (x_i - l_i^{(1)})^2}{2\sigma^2}\right)$$

if $x \approx l^{(1)}$: $f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$

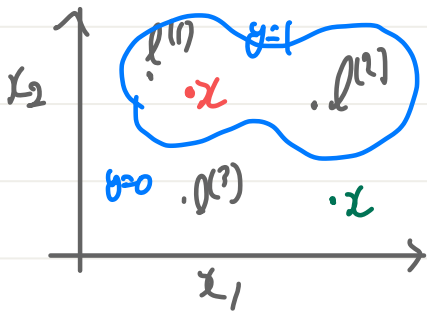
If x far from $l^{(1)}$: $f_1 \approx \exp\left(-\frac{(\text{large num})^2}{2\sigma^2}\right) \approx 0$

* gaussian 함수는 정규분포다 같이 bell-curve 형태고 landmark가 가까울수록 큰 값, 멀수록 작아짐

12. Support Vector Machines



③ Example



predict "y" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

Set, $\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$

red x is $l^{(1)}$ or $l^{(2)}$ $f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= -0.5 + 1 = 0.5 \geq 0 \Rightarrow \text{red } x \text{ is } y=1 \text{ 영역에 포함}$$

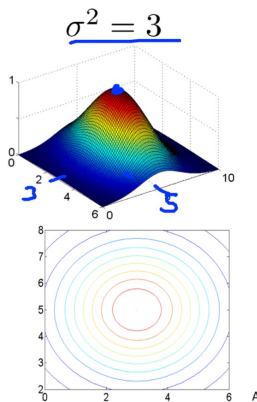
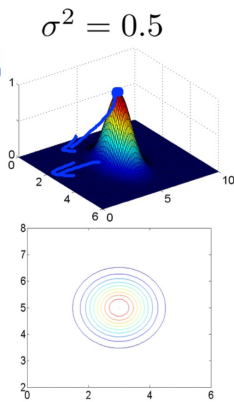
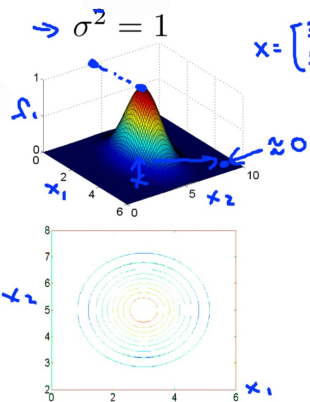
green x is $l^{(3)}$ or $l^{(4)}$ $f_1, f_2, f_3 \approx 0$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 = -0.5 < 0 \Rightarrow \text{green } x \text{ is } y=0 \text{ 영역에 포함}$$

④ r^2 function

ex) $l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$

$$f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$



σ^2 커질수록 그래프의 기울기 낮아진다!

12.5 Kernels 2

① choosing the landmarks

\Rightarrow 한 가지 방법은 landmark는 training example과 같은 위치에 놓는 것

Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})$

choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(n)} = x^{(n)}$

Given example x :

$$f_1 = \text{similarity}(x, l^{(1)} = x^{(1)})$$

$$f_2 = \text{similarity}(x, l^{(2)} = x^{(2)})$$

\vdots

For training example $(x^{(i)}, y^{(i)})$:

$$f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)})$$

$$f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)})$$

\vdots

$$f_n^{(i)} = \text{sim}(x^{(i)}, l^{(n)})$$

$$f^{(i)} = \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_n^{(i)} \end{bmatrix}, f^{(i)} = 1$$

← actually $f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)} = x^{(i)}) = \exp(-\frac{0}{2\sigma^2}) = 1$

12. Support Vector Machines



② SVM with kernels

hypothesis: given x , compute features $f \in \mathbb{R}^{m+1}$
→ predict "y=1" if $\theta^T f \geq 0$

Training:

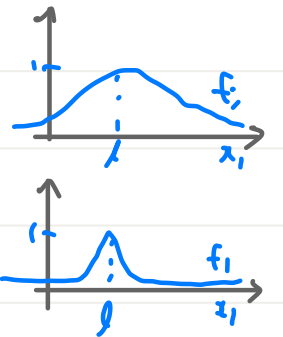
$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^m \theta_j^2 \quad (x^{(i)} \text{ 대신 } f^{(i)})$$

✗ 꼭 SVM이 아니더라도 kernel 개념 도입할 수 있지만 SVM의 classification 성능 향상도 매우 크기에 \Rightarrow kernel은 거의 SVM에서만 사용

③ SVM parameters

$C (= \frac{1}{\lambda})$ { Large C : Lower bias, high variance (small λ) overfit
Small C : Higher bias, low variance (large λ) underfit

σ^2 { Large σ^2 : features f_i vary more smoothly \Rightarrow Higher bias, lower variance
Small σ^2 : features f_i vary less smoothly. \Rightarrow Lower bias, higher variance



12.6 Using an SVM

• SVM은 잘 알려지고, 잘 구현된 library가 많다 직접 구축하기 보다 잘 만들어진 package 이용하는 편이 낫다 eg) liblinear, libsvm

• C , kernel의 종류 명시해 주어야 함

ex) No kernel (linear kernel): $\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$
Predict "y=1" if $\theta^T x \geq 0$ $\rightarrow n$ large, m small $x \in \mathbb{R}^n$

ex) Gaussian kernel:

$f_i = \exp\left(-\frac{\|x - x^{(i)}\|^2}{2\sigma^2}\right)$, where $x^{(i)} = x^{(i)}$ Note) feature scaling 필요 (가우시안 kernel 사용전)
Need to choose σ^2

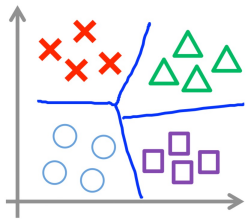
Note) 모든 similarity function " $\text{sim}(x, l)$ "은 kernel로 사용할 수 있는 것은 아님.

"Mercer Theorem"이라는 조건 만족해야 SVM package의 optimization 과정이 보장되지 않고 제대로 동작함

12. Support Vector Machines



① Multi-class classification



$$y \in \{1, 2, 3, \dots, k\}$$

많은 SVM package 등이 이미 multi-class classification을 가지고 있다.

"one vs all" method

K개의 SVM train 한다. 구분된 $y = i$ 에 대해 ($i = 1, 2, \dots, k$)

$\Rightarrow \theta^{(1)}, \theta^{(2)}, \dots, \theta^{(k)}$ 중 $\theta^{(i)}$ 중의 $\text{largest } (\theta^{(i)})^T x$ 를 만족하는 class i 선택

② Logistic regression vs SVM

$n = \#$ features ($x \in \mathbb{R}^n$)

$m = \#$ training examples

① if n is large (relative to m) $n \geq m$ ex) $n = 10,000$, $m = 10 \sim 100$

\Rightarrow Use logistic regression or SVM without a kernel (linear kernel)

② If n is small, m is intermediate $n = 1 \sim 1000$, $m = 10 \sim 10,000$

\Rightarrow Use SVM with Gaussian kernel

③ If n is small, m is large $n = 1 \sim 1000$, $m = 50,000^+$

\Rightarrow Create/add more features, then use logistic regression or SVM w/o kernel

④ Neural Network likely to work well for most of these settings, but may be slower to train