

Lecture1. Learning from data

- Learning from data

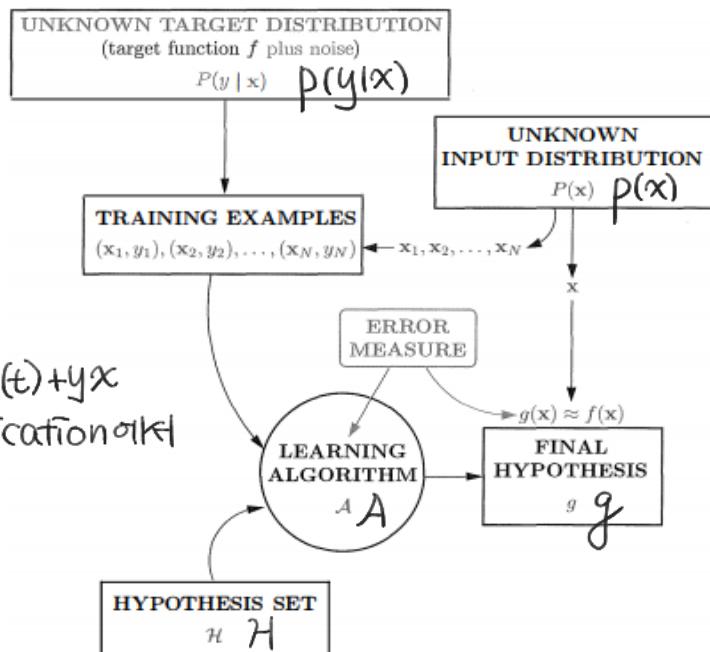
- data
- pattern
- 수학적인 분석 불가능

- Perceptron

$$h(x) = \text{sign}(w^T x)$$

- => 우리는 supervised learning

* PLA
 $w(t+1) = w(t) + yx$
 classification에서



Lecture2. Feasibility of learning

- Hoeffding inequality

$$\mathbb{P}[|\hat{\mu} - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

- Multiple hypothesis : H (hypothesis set)에 여러 개의 h 가 있는 경우

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2M e^{-2\epsilon^2 N} : M \text{ 빼 } looser, M \text{ must be finite}$$

- Feasibility of learning

1. $E_{out}(g) \sim E_{in}(g)$ 라고 할 수 있는가?

2. $E_{in}(g)$ 를 충분히 작게 만들 수 있는가?

문제의 E_{in}, E_{out}

- Error and Noise (다르다)
 $f \neq p(y|x)$

$$\rightarrow E_{out}(h) = \mathbb{E}_x [e(h(x), f(x))]$$

noisy하면 $E_{out}(h) = \mathbb{E}_{x,y} [e(h(x), y)]$

Lecture3. Linear classification and regression

→ small VC dim \Rightarrow generalization 퍼포먼스 좋아.

Linear classification

use signal: $w^T x$,

이 때 $x \in \mathbb{R}^d, w \in \mathbb{R}^d$

- Linear model : 1) $E_{in} \sim E_{out}$

VC dimension, dvc : $d+1$

VC generalization bound : (m: effective 한 hypothesis set 의 크기)

$$E_{out} \leq E_m + \sqrt{\frac{8}{N} \ln \frac{4mH(2N)}{\delta}} \Rightarrow mH(N) \leq N^{dvc+1}$$

probability $\geq 1 - \delta$

$$E_{out}(g) = E_{in}(g) + O\left(\sqrt{\frac{d}{N} \ln N}\right)$$

따라서, 높은 확률로

이므로, large N에서 $E_{in} \sim E_{out}$

- Linearly in-separable data 1. 전체적으로 linearly separable 하지만, outlier(noise)가 존재 minimum E_{in} 을 가지도록 문제를 풀자.

ex2) digit recognition

raw input space \rightarrow feature input space

feature space : reduce dimension, better handle curse of dimensionality,

but can cause loss of information

이후, PLA나 pocket algorithm으로 해결

- linearly separable
 - : $E_m(w^*) = 0$ 인즉하는 w^* 가 언제나 존재 \rightarrow PLA 등
- separable x
 - : E_m 작게 만들 수 있다. \rightarrow Pocket, SVM 등

- (1) 우리가 고려해야 하는 것
 자주적으로 쓰가 한다.
- (2) 전체에서 우리가 아는 부분이
 자주적으로 강조한다.

Linear regression: y is continuous

same bell shape

- assumption: homoscedasticity: (same finite variance for each value of x)

여러 평가방법 1) least squares, 2) maximum likelihood

- OLS solution w^* : minimizing $\text{Ein}(w)$

$$\text{Ein}(h) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2 = \frac{1}{N} \|Xw - y\|^2$$

OLS definition:

linear regression에서, $h = w^T x$, ($w^T x$ called signal)

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}_N, \quad w \in \mathbb{R}^d$$

이때, $\text{Ein}(w)$ is continuous, differentiable, convex, $\nabla \text{Ein}(w) = 0$

$$\nabla \text{Ein}(w) = \frac{2}{N} (X^T X w - X^T y) = 0 \text{을 만족하는 } w.$$

즉 normal equation ($X^T X w = X^T y$)을 만족하는 w 를 찾아야 한다.

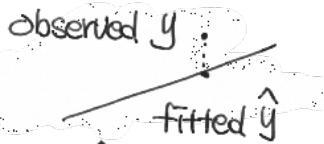
$$X^T X w = X^T y$$

- $X^T X$ is invertible, $w = (X^T X)^{-1} X^T y = X^{-1} y$ pseudo-inverse of X , most cases $\rightarrow 1$ step
- $X^T X$ is not invertible, X^{-1} still defined, but no unique solution

invertible 하지 않으면 SVD 등 다양한 방법이 존재

- Hat matrix H (projection matrix)

H 는 Ein 과 Eout 의 관계를 나타낸다.



in-sample error에 의해 $y \neq \hat{y}$

\hat{y} 는 y 를 X 의 column space로 orthogonal projection.

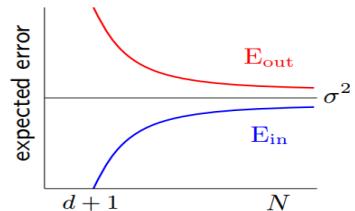
$$\text{fitted } \hat{y} = Xw^* = X(X^T X)^{-1} X^T y = Hy, \quad H = X(X^T X)^{-1} X^T$$

- generalization bound: $\text{Eout} = \text{Ein} + O(d/N)$

if assume $y = f(x) + \epsilon$ ↑
D mean. $0 < \infty$

$$\text{Ein} = \sigma^2 \left(1 - \frac{d+1}{N} \right) < \sigma^2$$

$$\text{Eout} = \sigma^2 \left(1 + \frac{d+1}{N} \right) > \sigma^2$$



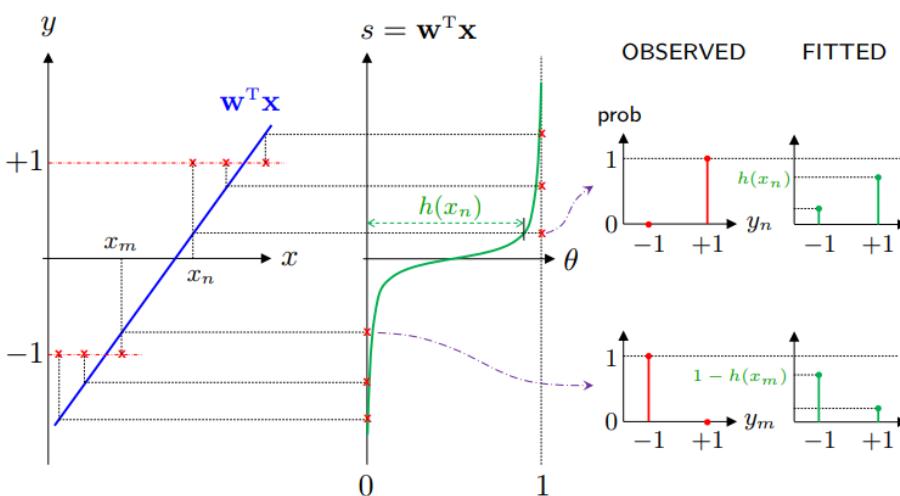
Lecture4. Logistic regression

- Logistic regression

input: $h(x) = \theta(w^T x)$, θ is sigmoid, θ 의 특징: $1 - \theta(s) = \theta(-s)$

- Big picture

$$\theta(w^T x) = \frac{1}{1 + e^{-w^T x}} : \text{sigmoid}$$



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

$$P(y|x) = \theta(yw^T x)$$

$$= \begin{cases} \theta(w^T x) & \text{if } y=1 \\ \theta(-w^T x) & \text{if } y=-1 \\ 1 - \theta(w^T x) \end{cases}$$

$$= h(x)^{[y=1]} + (1 - h(x))^{[y=-1]}$$

$$\begin{aligned} \mathbb{P}[y=+1|x] &= f(x) \sim h(x) = \theta(w^T x) \\ P(y|x) &= h(x)^{[y=+1]} (1 - h(x))^{[y=-1]} = \theta(yw^T x) \end{aligned}$$

Logistic regression 의 Error \rightarrow 유도과정에서 $\theta(-S) = 1 - \theta(S)$ 임을 이용.

target, noisy.

- Learning target

$$f(x) = P(y|x) = (h(x) \text{ for } y=+1), (1-h(x) \text{ for } y=-1)$$

- Error measure : Simple 한 것 말고, cross-entropy 사용

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$$

2) cross-entropy measure

- Error measure 유도 방법 1: likelihood

likelihood : $P(y|x) = \theta(y \mathbf{w}^T \mathbf{x})$

(이 때, $1-\theta(x) = \theta(-x)$ 임을 이용)

$$\ln \frac{1}{P(y_n|\mathbf{x}_n)} = [y=+1] \ln \frac{1}{h(\mathbf{x}_n)} + [y=-1] \ln \frac{1}{1-h(\mathbf{x}_n)}$$

↑ ↑ ↑
observed fitted

$$P(y_1|\mathbf{x}_1)P(y_2|\mathbf{x}_2) \cdots P(y_N|\mathbf{x}_N) = \prod_{n=1}^N P(y_n|\mathbf{x}_n)$$

likelihood of training data :

maximum likelihood h 를 select

$$\text{수식을 변형 : minimizing } -\frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n|\mathbf{x}_n) \right) = \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n|\mathbf{x}_n)}, \text{ wrt } \mathbf{w}$$

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$$

결론 : in-sample error :

$P(y|x) = \theta(y \mathbf{w}^T \mathbf{x})$ 를 대입하고, sigmoid 함수를 대입해 정리한 식

implied point-wise error : $\ln(1+e^{-(y \mathbf{w}^T \mathbf{x})})$

- Error measure 유도 방법 2: cross-entropy

cross entropy : $(p)\log(1/q)+(1-p)\log(1/(1-q))$

: 'error' for 'observed' pmf $\{p, 1-p\}$ by 'fitted' pmf $\{q, 1-q\}$

$$\begin{aligned} NLL(\mathbf{w}) &\triangleq -\log \left\{ \prod_{n=1}^N P(y_n|\mathbf{x}_n) \right\} \\ &= -\log \left\{ \prod_{n=1}^N h(\mathbf{x}_n)^{[y_n=+1]} (1-h(\mathbf{x}_n))^{[y_n=-1]} \right\} \\ &= \sum_{n=1}^N \left\{ [y_n = +1] \log \frac{1}{h(\mathbf{x}_n)} + [y_n = -1] \log \frac{1}{1-h(\mathbf{x}_n)} \right\} \end{aligned}$$

: cross entropy error function

Logistic regression 의 optimization : iterative (e.g. gradient descent)

- set $\nabla E_{in}=0$ (in error surface)

Iteratively solve
unconstrained optimization

$$\nabla_{\mathbf{w}} E_{in} = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}} = -\frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n)$$

no analytic solution: Iteratively 풀자.

어디로? $-\nabla E_{in}$ 방향

얼만큼? 가만히. $|\nabla E_{in}|$ 만큼.

- Gradient descent

- Algorithm : gradient descent

매번 모든 \mathbf{x}_n 을 E_{in} 을 계산,

for (iteratively)

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla E_{in}(\mathbf{w}(t))$$

compute $\nabla E_{in}(\mathbf{w}(t))$

B는 no이 아니라

set v^{\wedge} = negative gradient = $-\nabla E_{in}(\mathbf{w}(t))$,

- v^{\wedge} is due to cancellation of $|\nabla E_{in}|$

subset, 또 1개의 point 를

update weights $\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla E_{in}(\mathbf{w}(t))$

update 하는 방법 사용하기로.

- ∇E_{in} 의 계산이 $O(N)$ time if batch-> 다 볼 필요 없는 mini batch, stochastic 사용하기도

언제 끝낼까? 그레이디언트 디센트

- ① Iteration 수로 bound
- ② gradient 작고
- ③ error itself 작고

Linear model review

	linear classification	linear regression	logistic regression	= probability estimation
y	$\{-1, +1\}$	\mathbb{R}	$\{-1, +1\}$	
$\hat{y} = h(\mathbf{x})$	$\text{sign}(\mathbf{w}^T \mathbf{x})$	$\mathbf{w}^T \mathbf{x}$	$\theta^*(\mathbf{w}^T \mathbf{x})$	
$e(\hat{y}, y)$	0-1 loss $\llbracket \hat{y} \neq y \rrbracket$	squared error $(\hat{y} - y)^2$	cross-entropy error $\llbracket y=+1 \rrbracket \ln \frac{1}{\hat{y}} + \llbracket y=-1 \rrbracket \ln \frac{1}{1-\hat{y}}$	
$E_{in}(h)$ opt.	$\frac{1}{N} \sum_{n=1}^N \llbracket h(\mathbf{x}_n) \neq y_n \rrbracket$ combinatorial optimization (NP-hard) PLA.	$\frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$ set $\nabla E_{in}(\mathbf{w}) = 0$ (closed-form solution exists) $\nabla \neq 0$	$\frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})$ set $\nabla E_{in}(\mathbf{w}) = 0$ iterative optimization (e.g. gradient descent)	

* VC generalization bound

① linear classification

$$d_{VC} = d+1, \quad m_H(N) \leq N^{d_{VC}+1}$$

$$\therefore E_{out}(g) = E_{in}(g) + O\left(\sqrt{\frac{d}{N} \ln N}\right)$$

② linear regression

$$\therefore E_{out}(g) = E_{in}(g) + O\left(\frac{d}{N}\right)$$

→ linear은 first try로 좋아.
Simple, robust, work well, generalize well

Lecture 5. Artificial neural networks

- 분포를 구현하기 위한 것이 함수

	K class	Binary
함수	Softmax function	Logistic function
1 번의 event	멀티플리 distribution = categorical distribution (주사위 한번 던지는 것)	베르누리 distribution (동전 한번 던지는 것)
N 번의 event	Multinomial distribution	Binomial distribution

softmax function $\sigma : \mathbb{R}^K \mapsto \mathbb{R}^K$

$$\sigma(\mathbf{h})_j = \frac{e^{h_j}}{\sum_{k=1}^K e^{h_k}}$$

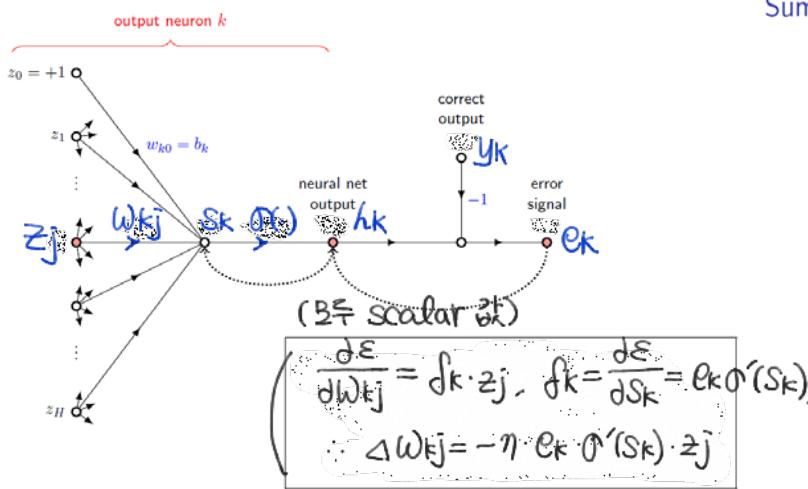
- Error measures

- 1) ϵ_k : 하나의 원소에 대해: $\epsilon_k = h_k - y_k$
 - 2) ϵ_n : 하나의 vector에 대해: $\epsilon_n = \frac{1}{2} \sum_{k=1}^d \epsilon_k^2$
 - 3) ϵ_D : 전체 vector set, data에 대해: $\epsilon_D = \frac{1}{2N} \sum_{n=1}^N \sum_{k=1}^d \epsilon_{kn}^2$
- 우리는 ϵ_D 를 최소화하는 w 찾고 싶다.

back prop.
first order method

Back propagation in output layer

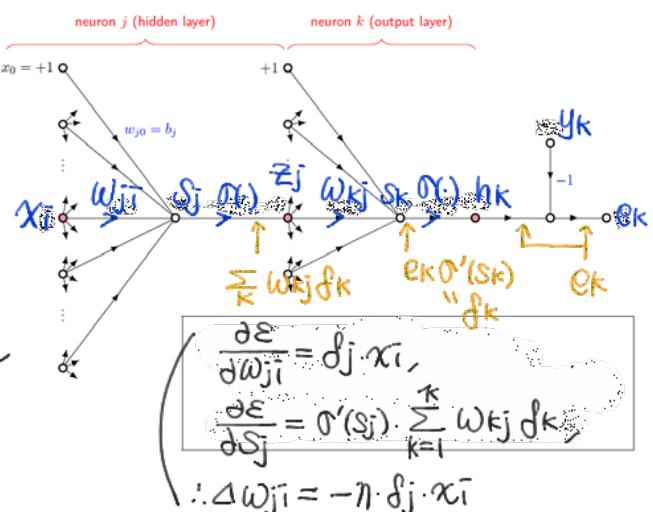
Summary: output layer



(이때, error, input이 클수록 w 를 많이 변화시킨다)

Back propagation in input layer

Summary: hidden layer



* 이때, 각각 다른 layer들은 fan out,
∴ backpropagate sum을 해주는건 규칙!

한번에 하나의 n pick하지 않고
subset을 pick해서
평균으로 update할 수도 있다.

- Back prop Algorithm (다음을 repeat)

pick n (하나의 x 를 pick)

forward : compute all z, h

backward : compute all delta

update weights

until it is time to stop, update final weights

$\frac{\partial \epsilon}{\partial w_{kj}}$ 를 모두 계산,
모든 w 를 update.

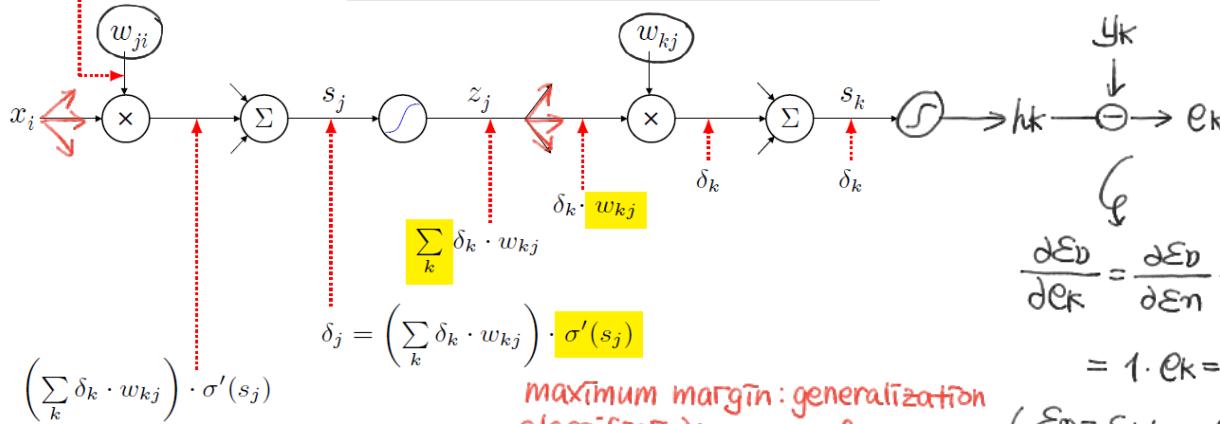
Back prop 간단한 계산

- 간단한 계산 공식

- 1) $out = f(in)$: $f'(in)$ 을 곱한다.
- 2) $out = in1 * in2$: $in1$ 에 대한 식을 구하는 경우 $in2$ 를 곱한다.
- 3) $out = \text{sum}(in)$: 그대로 fanout 한다.
- 4) $out = \text{max}(in)$: in_i 가 max 이면 fanout, 아니면 0이다.
- 5) $out = \text{in}, (\text{fanout})$: sum으로 모두 더한다.

- computing $\frac{\partial E}{\partial w_{ji}}$

$$\left(\sum_k \delta_k \cdot w_{kj} \right) \cdot \sigma'(s_j) \cdot x_i \Rightarrow \frac{\partial E}{\partial w_{ji}} = \delta_j \cdot x_i = \left(\sum_k \delta_k \cdot w_{kj} \right) \cdot \sigma'(s_j) \cdot x_i$$



$$\begin{aligned} \frac{\partial E}{\partial e_k} &= \frac{\partial E}{\partial \varepsilon_n} \cdot \frac{\partial \varepsilon_n}{\partial e_k} \\ &= 1 \cdot e_k = e_k \end{aligned}$$

$$\begin{aligned} E_D &= E_1 + \dots + E_N \\ \varepsilon_n &= \frac{1}{2}(e_1^2 + \dots + e_N^2) \end{aligned}$$

Lecture6. Support vector machines : SVM

- Linear discriminant function: $g(x) = w^T x + b$
 - big picture:
- 이제는 b 를 따로 생각한다.

최적화에 대한 내용

풀고자 하는 문제: Constraint optimization
 $(\max, \min f) (g \leq 0, g \geq 0) \rightarrow L = f \pm \lambda g$ 이용.

$$\rightarrow L = f \pm \lambda g : \begin{cases} \textcircled{1} \lambda \geq 0 & \textcircled{2} \text{제약 } g \geq 0, \forall g \leq 0 \\ \textcircled{3} \nabla_x L = \nabla_x f \pm \lambda \nabla_x g = 0 \\ \textcircled{4} \lambda g = 0 \quad (\text{둘다 } 0 \text{은 아닙니다}) \end{cases}$$

4) 4가지 case.

$$\begin{aligned} (\text{f}) &\max \text{ 문제. } \min \text{ 문제} \\ (\text{g}) &\text{ } g \geq 0, \quad g \leq 0. \end{aligned}$$

case A) $\min f, \quad g \geq 0$.

$$\begin{aligned} x^* \text{보다 크라} & \quad \nabla g \downarrow \quad \nabla f \downarrow \\ g > 0 & \quad \therefore \nabla f - \lambda \nabla g = 0 \\ \therefore \nabla f - \lambda \nabla g &= 0 \end{aligned}$$

case B) $\min f, \quad g \leq 0$

$$\begin{aligned} x^* \text{보다 작라} & \quad \nabla g \uparrow \quad \nabla f \downarrow \\ g < 0 & \quad \therefore \nabla f + \lambda \nabla g = 0 \\ \therefore \nabla f + \lambda \nabla g &= 0 \end{aligned}$$

case 1) $\lambda = 0, g \neq 0$

$$\begin{aligned} x^* & \quad g = 0 \\ \text{f 등고선} & \\ g: \text{inactive constraint} & \end{aligned}$$

case C) $\max f, \quad g \geq 0$.

$$\begin{aligned} x^* \text{보다 작라} & \quad \nabla g \downarrow \quad \nabla f \uparrow \\ g > 0 & \quad \therefore \nabla f + \lambda \nabla g = 0 \\ \therefore \nabla f + \lambda \nabla g &= 0 \end{aligned}$$

case D) $\max f, \quad g \leq 0$

$$\begin{aligned} x^* \text{보다 크라} & \quad \nabla g \uparrow \quad \nabla f \uparrow \\ g < 0 & \quad \therefore \nabla f - \lambda \nabla g = 0 \\ \therefore \nabla f - \lambda \nabla g &= 0 \end{aligned}$$

case 2) $\lambda > 0, g = 0$

$$\begin{aligned} x^* & \quad g = 0 \\ \text{f 등고선} & \\ g: \text{active constraint} & \end{aligned}$$

Primal form : 이해하기 좋다. But 의미를 파악하기 힘들다.

- SVM primal (linearly classified)

SVM (primal) optimization problem

$$\begin{array}{ll} \text{minimize}_{\mathbf{w}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} & y_t (\mathbf{w}^\top \mathbf{x}_t + b) \geq 1, \forall t \in [1, N] \end{array} \rightarrow \begin{array}{l} \mathbf{w}, b \text{에 대해 convex} \\ \mathbf{w}, b \text{에 대해 linear} \end{array}$$

- QP (Quadratic programming)

$$\begin{array}{ll} \text{minimize}_{\mathbf{z}} & \frac{1}{2} \mathbf{z}^\top Q \mathbf{z} + \mathbf{c}^\top \mathbf{z} \\ \text{subject to} & A \mathbf{z} \geq \mathbf{a} \end{array}$$

solution : $\mathbf{z}^* \leftarrow \text{QP}(Q, \mathbf{c}, A, \mathbf{a})$

$$\left(\begin{array}{l} \mathbf{z} = \begin{bmatrix} \mathbf{b} \\ \mathbf{w} \end{bmatrix}, Q = \begin{bmatrix} 0 & 0 \\ 0 & I_d \end{bmatrix}, \mathbf{c} = \mathbf{0}_{d+1} \\ A = \begin{bmatrix} \mathbf{y}_1 \mathbf{y}_1^\top \\ \vdots \\ \mathbf{y}_N \mathbf{y}_N^\top \end{bmatrix}_N, \mathbf{a} = \mathbf{1}_N \end{array} \right) \text{ or } \left(\begin{array}{l} \mathbf{z} = \begin{bmatrix} \mathbf{b} \\ \mathbf{w} \end{bmatrix}, Q = \begin{bmatrix} I_d & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{c} = \mathbf{0}_{d+1} \\ A = \begin{bmatrix} \mathbf{y}_1 \mathbf{y}_1^\top & \mathbf{y}_1 \\ \vdots & \vdots \\ \mathbf{y}_N \mathbf{y}_N^\top & \mathbf{y}_N \end{bmatrix}_N, \mathbf{a} = \mathbf{1}_N \end{array} \right)$$

Dual form : support vector, kernel trick에 대한 이해 가능

- Dual formulation

if cost, constraint functions are strictly convex, dual form is exist
and solving dual form is equivalent to solving primal

- Step 1 : compute Lagrangian

위의 primal form 을 constraint 가 없는 Lagrangian optimization 문제로 변환 :

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha_t [y_t (\mathbf{w}^\top \mathbf{x}_t + b) - 1]$$

α (알파) : Lagrangian multipliers = dual variables, 0 이상의 값을 가진다.

-> solution 은 saddle point for Lagrangian

L 은 w,b 에 대해서는 min, a 에 대해서는 max : saddle point

- Step 2 : KKT conditions

condition 1: $\frac{\partial \mathcal{L}(\mathbf{w}^*, b^*, \alpha^*)}{\partial \mathbf{w}} = 0$) 미분해서 0

condition 2: $\frac{\partial \mathcal{L}(\mathbf{w}^*, b^*, \alpha^*)}{\partial b} = 0$) stationarity $\nabla f - \lambda \nabla g = 0$

condition 3: $\underbrace{\alpha_t^* [y_t (\mathbf{w}^{*\top} \mathbf{x}_t + b^*) - 1]}_{\alpha_t^* = 0 = 0} = 0, \forall t$) $\lambda g = 0$

y_t 가 흑자여서

존재하니 \rightarrow Support vector

active constraint if S.V.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{t=1}^N \alpha_t y_t \mathbf{x}_t = 0$$

$$\frac{\partial \mathcal{L}}{\partial b} = \sum_{t=1}^N \alpha_t y_t = 0$$

$$\mathbf{w} = \sum_{t=1}^N \alpha_t y_t \mathbf{x}_t = \sum_{t=1}^N \alpha_t y_t \mathbf{x}_t$$

$\mathbf{w} \neq 0$ 할 때 적은 수의 S.V.로 구할 수.

S.V.의 intuition이 살아나면 dual.

- Step 3 : formulate dual problem

기존의 Lagrangian에서 식을 변형하면

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha_t y_t \mathbf{w}^\top \mathbf{x}_t - b \sum_{t=1}^N \alpha_t y_t + \sum_{t=1}^N \alpha_t$$

primal 의 w, dual 은 a에 대한 식

inner products $\mathbf{X}^\top \mathbf{X} \rightarrow$ can use kernel trick

따라서, nonlinearily separable case 도 고려할 수.

dual a에 대한 max

maximize α

$$\sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \mathbf{x}_s^\top \mathbf{x}_t$$

subject to

$$\sum_{t=0}^N \alpha_t y_t = 0$$

$$\alpha_t \geq 0, \forall t \in [1, N]$$

kernel trick (inner product)

dual a에 대한 N개 조건 추가해.

- Step 4 : solve dual problem as QP

- QP solver (凸优化問題 QP)

$$\begin{aligned} \text{minimize}_{\mathbf{x}} \quad & \frac{1}{2} \mathbf{x}^T Q \mathbf{x} + \mathbf{c}^T \mathbf{x} \quad (\alpha \text{의 } \mathbf{x}) \\ \text{subject to} \quad & A \mathbf{x} \leq \mathbf{b} \\ & E \mathbf{x} = \mathbf{d} \end{aligned}$$

$\max \rightarrow$ \rightarrow 배수 $\min \alpha$ 로

- Dual form 을 QP에 적용하면

$$\begin{aligned} \text{minimize}_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \boldsymbol{\alpha}^T Q \boldsymbol{\alpha} - \mathbf{1}^T \boldsymbol{\alpha} \\ \text{subject to} \quad & -\boldsymbol{\alpha} \leq \mathbf{0} \\ & \mathbf{y}^T \boldsymbol{\alpha} = \mathbf{0} \end{aligned}$$

$$\mathbf{x} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{bmatrix} = \boldsymbol{\alpha}$$

$$\mathbf{C} = -\mathbf{1}_N \quad \frac{x_i^T x_j}{N}$$

$$\boldsymbol{\alpha} = \left[\frac{y_i y_j \langle x_i, x_j \rangle}{N} \right]_N$$

$$A = -I_N, \quad b = \mathbf{0}_N$$

$$E = [y_1 \dots y_N], \quad d = 0$$

- Step 5 : complete

SVM classifier $g(\mathbf{x})$: linearly separable case

* Q 매우 dense 한 matrix.
α 매우 sparse 한 vector.

generalization bound of SVM (Vapnik, 1995):

$$\begin{aligned} g(\mathbf{x}) &= \text{sign} \left(\sum_{t=1}^N \alpha_t^* y_t \mathbf{x}_t^T \mathbf{x} + b^* \right) \\ &= \text{sign} (\mathbf{w}^* \mathbf{x} + b^*) \quad \mathbf{w}^* = \sum \alpha_t^* y_t \mathbf{x}_t \end{aligned}$$

$$\mathbb{E}[E_{\text{out}}] \leq \frac{\mathbb{E}[\# \text{ of support vectors }]}{N-1} \quad (\text{매우 좋다})$$

E_{out} 매우 좋다. * S.V. pt depend

Lecture7. Support vector machines (linearly non-separable)

Linearly non-separable case 1) kernel trick 으로 해결

- Non-linear transform

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \mathbf{x}_s^T \mathbf{x}_t$$

kernel trick, inner product.
 $\mathbf{z}_s^T \mathbf{z}_t = k(\mathbf{x}_s, \mathbf{x}_t)$

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \mathbf{z}_s^T \mathbf{z}_t$$

- Kernel: non-linear similarity measure

- polynomials (of degree $q > 0$):

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{a} \mathbf{x}^T \mathbf{x}' + b)^q$$

- hyperbolic tangent (aka sigmoid, multilayer perceptron kernel):

$$K(\mathbf{x}, \mathbf{x}') = \tanh(\beta \mathbf{x}^T \mathbf{x}' + \gamma)$$

- radial-basis functions ($\gamma > 0$):

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

- multiquadric:

$$K(\mathbf{x}, \mathbf{x}') = \sqrt{\|\mathbf{x} - \mathbf{x}'\|^2 + c^2}$$

- Gaussian RBF:

$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right)$$

- wave:

$$K(\mathbf{x}, \mathbf{x}') = \frac{\theta}{\|\mathbf{x} - \mathbf{x}'\|} \sin \frac{\|\mathbf{x} - \mathbf{x}'\|}{\theta}$$

- equivalent kernel :

poly에서 $a=b=0$

- linear kernel : poly에서 $q=1$

- RBF kernel :

infinite feature dim,
ROI 크면 over fitting.

- Valid kernel 인가? : Mercer's theorem

Let function $K : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$ be given. For K to be a valid (Mercer) kernel, it is necessary and sufficient that, for any finite set of instances, the corresponding kernel matrix is

$$k_{ij} = K(x_i, x_j) \text{ if }$$

: 1) symmetric 2) positive semi-definite ← 어떤 finite set에 대해서도 성립.

kernel matrix : $K = (m, m)$ matrix, $K_{ij} = K(x_i, x_j)$

positive semi-definite : PSD

all eigen values, determinant is non-negative, real symmetric K 가 $t^T K t \geq 0$ for all vector

Linearly non-separable case2) soft-margin SVM : C가 클수록 hard margin에 가까워

- Reformulating SVM problem : primal form \rightarrow 원래 primal에서 $+C\sum \xi_t$, 조건식 변형과 추가.

minimize	$\frac{1}{2} \ \mathbf{w}\ ^2 + C \sum_{t=1}^N \xi_t$: primal value
subject to	$y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq 1 - \xi_t, \quad \forall t$: $\mathbf{w}, b, \xi_t \forall t$.
	$\xi_t \geq 0, \quad \forall t$	

- Reformulating SVM problem : dual form

- Lagrangian $\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta)$

$= \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{t=1}^N \xi_t - \sum_{t=1}^N \alpha_t (y_t(\mathbf{w}^\top \mathbf{x}_t + b) - 1 + \xi_t) - \sum_{t=1}^N \beta_t \xi_t$: dual value : $\alpha_t, \beta_t \forall t$.
--	---

- minimize wrt \mathbf{w}, b , and ξ
- maximize wrt each $\alpha_t \geq 0$ and $\beta_t \geq 0$

- KKT conditions are obtained as follows:

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{t=1}^N \alpha_t y_t \mathbf{x}_t = \mathbf{0} \Rightarrow \mathbf{w} = \sum_{t=1}^N \alpha_t y_t \mathbf{x}_t$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{t=1}^N \alpha_t y_t = 0 \Rightarrow \sum_{t=1}^N \alpha_t y_t = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_t} = C - \alpha_t - \beta_t = 0 \Rightarrow \alpha_t + \beta_t = C$$

$$\alpha_t [y_t(\mathbf{w}^\top \mathbf{x}_t + b) - 1 + \xi_t] = 0 \Rightarrow \underbrace{\alpha_t = 0}_{\beta_t \xi_t = 0} \vee \underbrace{y_t(\mathbf{w}^\top \mathbf{x}_t + b) = 1 - \xi_t}_{\beta_t = 0}$$

$$\beta_t \xi_t = 0 \Rightarrow \underbrace{\beta_t = 0}_{\xi_t = 0} \vee \underbrace{\xi_t = 0}_{\beta_t = 0}$$

Final form, 아래의 식을 QP에 적용

final solution

Soft-margin SVM classifier

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

where

$$\mathbf{w} = \sum_{0 < \alpha_t \leq C} \alpha_t y_t \mathbf{x}_t, \quad \mathbf{w}^* = \sum_{0 < \alpha_t \leq C} \alpha_t y_t \mathbf{x}_t$$

$$b = y_{sv} - \sum_{0 < \alpha_t \leq C} \alpha_t y_t \mathbf{x}_t^\top \mathbf{x}_{sv}, \quad b = y_{sv} - \mathbf{w}^* \mathbf{x}_{sv}$$

$w = \text{sum}(ayx)$ 에서 $a=C$ 이므로 w 에서 비중이 크다.

$\alpha = C$ 인 S.V를

- 모두 정리

	모범생	원래 S.V	군사분계선	간접
$y_t(\mathbf{w}^\top \mathbf{x}_t + b)$	$y(\mathbf{w}^\top \mathbf{x} + b) > 1$	$y(\mathbf{w}^\top \mathbf{x} + b) = 1$	$0 \leq y(\mathbf{w}^\top \mathbf{x} + b) < 1$	$y(\mathbf{w}^\top \mathbf{x} + b) < 0$
ξ_t	$E=0$	$E=0$	$0 < E \leq 1$	$1 < E$
a (alpha)	$a=0$	$0 < a \leq C$	$a=C$	$a=C$
b (beta)	$b \neq 0$	$b \neq 0$	$b=0$	$b=0$
Vector x lies	Out of margin	On margin	In margin	The other area
분류	right	right	right	wrong