

Lecture1. Learning from data

Outline

Introduction

- Learning from Data
- Problem Setup
- A Simple Learning Model
- Perceptron

Types of Learning

Summary

- Learning from data
 - 1) data
 - 2) pattern
 - 3) 수학적인 분석 불가능
- Learning model
 - hypothesis set + learning algorithm
 - 이때, hypothesis set H 를 functional form $h(x)$ 로 나타낸다 : parametric model
- Perceptron
 - $h(x) = \text{sign} (w^T x)$
- PLA : Perceptron learning algorithm
 - iterative, converge in linearly separable data
 - $w(t+1) = w(t) + yx, 0 | \text{ 때 } (x,y) \in \text{misclassified point}$
- Learning의 종류
 - 1) Supervised learning : (input, correct output) :classification, regression

- 2) Reinforcement learning : (input, some output, grade) : game learning
- 3) Unsupervised learning : (input) : clustering, HMMs, feature extraction(PCA, ICA, SVD)
: finding patterns, precursor to supervised learning, to create **high level representation**

=> 우리는 supervised learning에 대해서 이야기한다.

Summary

- learning is used when
 - ▶ we have data, and a pattern exists
 - ▶ but we can hardly pin it down mathematically
- learning model: hypothesis set and learning algorithm
 - ▶ our first example: perceptron and PLA
- types of machine learning
 - supervised: (input, correct output/label)
 - unsupervised: (input, no label)
 - reinforcement: (input, some output, grade for this output)
- supervised learning: our main theme for a while
 - ▶ unknown target function $y = f(x)$
 - ▶ known training data set $(x_1, y_1), \dots, (x_N, y_N)$
 - ▶ learning algorithm picks $g \approx f$ from hypothesis set \mathcal{H}

Lecture2. Feasibility of learning

Outline

Learning Unknown Target Function

Hoeffding Inequality

Connection to Learning

Generalizing the Hoeffding Bound

Feasibility of Learning

Error and Noise

Error Measures

Noisy Targets

Summary

- Hoeffding inequality

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

outside D에 위치한 f 를 probabilistic way로 learn

$u(\text{out of sample error}) = u^*(\text{in sample error})$ is PAC (probably approximately correct)

- Multiple hypothesis

H (hypothesis set)에 여러 개의 h 가 있는 경우

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

- Feasibility of learning

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

Hoeffding ineq, VC analysis, bias-variance analysis

2. $E_{in}(g)$ 를 충분히 작게 만들 수 있는가?

various learning paradigms will help

- complexity of $H = M$: M 이 크면, fits data well so smaller $E_{in}(g)$, but $E_{in} \sim E_{out}$ more

- complexity of f : f 가 복잡하면 hard to fit, higher $E_{in}(g)$

- Error and Noise (다르다)

Error: point-wise error 의 평균(E_{in}), 또는 Expectation(E_{out})

문제에 따라 error 다르게 정의 but

1) plausible measures : squared error, Gaussian noise

2) friendly measures : closed-form solution, convex optimization

Noise : f is not uniquely determined by the input : f 가 conditional probability

regression (noise is added), rather than interpolation

deterministic target as the average of noisy target : $f(\mathbf{x}) = \mathbb{E}(y|\mathbf{x})$

$f(x) = \mathbb{E}(P(y|x))$
 $f(x) = P(y|x)$
distribution.

- $P(y|x)$ and $P(x)$

$P(y|x)$: target distribution

$P(x)$: input distribution

$P(x,y) = P(y|x)P(x)$: model (x,y) s are generated by $P(x,y)$

Summary

- learning *unknown* target f : only possible in a probabilistic sense

▶ we need hypothesis $g \approx f$, which means $E_{out}(g) \approx 0$

- feasibility of learning is split into two questions

1. can we make sure that $E_{out}(g) \approx E_{in}(g)$?

2. can we make $E_{in}(g) \approx 0$?

- answering Q1: theoretical (first half of this course)

▶ Hoeffding's inequality: $\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$

▶ VC analysis and bias-variance analysis

- answering Q2: practical (second half)

▶ various learning paradigms will help us achieve that

- error measure quantifies the meaning of \approx in expression $g \approx f$
 - ▶ ideally, should be specified by the user
 - ▶ practical alternatives: *plausible/friendly* measures
- averaging *pointwise* error gives
 - ▶ in-sample: $E_{\text{in}}(h) = \frac{1}{N} \sum_n e(h(\mathbf{x}_n), f(\mathbf{x}_n))$
 - ▶ out-of-sample: $E_{\text{out}}(h) = \mathbb{E}_{\mathbf{x}}[e(h(\mathbf{x}), f(\mathbf{x}))]$
- noisy targets
 - ▶ $y \sim P(y|\mathbf{x})$ instead of $y = f(\mathbf{x})$ [y is a **random variable**]
 - ▶ $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ generated by $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$
 - ▶ $E_{\text{out}}(h) = \mathbb{E}_{\mathbf{x}, y}[e(h(\mathbf{x}), y)]$ *noisy*

Lecture3. Linear classification and regression

Outline

Linear Classification

- Checking Two Basic Criteria for Learning
- Non-Separable Data and Pocket Algorithm
- Example: Handwritten Digit Recognition

Linear Regression

- Introduction
- Linear Regression Algorithm
- Interpretation via Hat Matrix

Summary

- Linear model
 - small VC dimension : generalize well from Ein to Eout
 - use "signal" $w^T x$
 - 3 problems
 - 1) classification
 - 2) regression
 - 3) probability estimation (logistic regression)
- | | |
|--|---|
| ● linear classification | ● linear regression |
| $\mathcal{Y} = \{-1, +1\}$ | $\mathcal{Y} = \mathbb{R}$ |
| $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$ | $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ |
| $e(\hat{y}, y) = [\hat{y} \neq y]$ | $e(\hat{y}, y) = (\hat{y} - y)^2$ |
| ▶ NP-hard in general | ▶ efficient closed-form solution |

Linear classification

- Linear model : 1) Ein ~ Eout

VC dimension, dvc : d+1

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

$$E_{\text{out}} \leq E_{\text{in}} + \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}}$$

probability $\geq 1 - \delta$

$$m_{\mathcal{H}}(N) \leq N^{d_{\text{VC}}} + 1$$

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

따라서, 높은 확률로 이므로, large N에서 Ein~Eout

- Linear model : 2) Ein ~ 0

1) linearly separable data :

언제나 $Ein(W^*)=0$ 을 만족하는 hypothesis w^* 가 존재한다. : PLA 등

2) linearly in-separable data :

-> line 으로 나눌 수 없거나, outlier(noise)가 존재하는 경우

많은 경우 Ein 을 매우 작게 만들 수 있다. : pocket algorithm, SVM 등

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

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \underbrace{\frac{1}{N} [\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n]}_{E_{\text{in}}(\mathbf{w})}$$

in general, 위의 식은 N P-hard 이므로 approximately minimize Ein

ex1) Pocket algorithm $\rightarrow w(t+1) = w(t) + \alpha y$

run PLA for $w(t+1)$, $Ein(w(t+1))$ 계산해서 개선되었으면 $w^* = w(t+1)$, 이를 반복

// $Ein(w(t+1))$ 의해, much slower than PLA and good Ein 으로 빠르게 수렴한다는 보장 x

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, PLA 나 pocket algorithm 으로 해결

① 우리가 고려해야 하는 것
자주적으로 증가해.
② 전체에서 우리가 아는 부분 비율
자주적으로 감소해.

Linear regression : y is continuous

- Regression : Target distribution rather than function
statistical method
assumption : homoscedasticity : (same finite variance for each value of x)

- Parameter estimation : 평가 방법
 - least squares
 - OLS : ordinary least squares : closed-form solution exist
 - generalized least squares, iteratively reweighted least squares
 - maximum likelihood
 - ridge , lasso regression
 - least absolute deviation regression

- OLS solution w^* : minimizing $E_{in}(w)$

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2, \quad \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \quad (\mathbf{X} \text{는 } \mathbf{x}_{nT} \text{ 가 각 행에 위치})$$

OLS definition : linear regression에서, $h = \mathbf{w}^T \mathbf{x}$, ($\mathbf{w}^T \mathbf{x}$ called signal)

*Input matrix
 \mathbf{w} is a vector.*

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

$$\nabla E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y})$$

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

$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0.$$

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

invertible 하지 않으면 SVD 등 다양한 방법이 존재

$$\mathbf{w} = \mathbf{X}^{-1} \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

- OLS solution ($\mathbf{X}^T \mathbf{X}$ is invertible)

provides the best linear unbiased estimator (BLUE)

one step learning, analytic formula

- 참고 SVD

assume: the SVD of input matrix X is $X = U\Gamma V^T$

- ▷ $U \in \mathbb{R}^{N \times \rho}$ satisfies $U^T U = I_\rho$,
- ▷ $V \in \mathbb{R}^{(d+1) \times \rho}$ satisfies $V^T V = I_\rho$,
- ▷ $\Gamma \in \rho \times \rho$ is a positive diagonal matrix with $\rho = \text{rank}(X)$
- ▶ then $w_{\text{lin}} = V\Gamma^{-1}U^T y$ satisfies $X^T X w_{\text{lin}} = X^T y$

(observed) y ↗
 \hat{y} (learn)

- Hat matrix H (projection matrix)

H 는 E_{in} 과 E_{out} 의 관계를 나타낸다. (y 를 projection하면 \hat{y})

in-sample error에 의해 $y \neq \hat{y}$ (y 는 관찰 값, \hat{y} 는 learn한 값)

$$\hat{y} = \underbrace{X}_{\text{H}} \underbrace{(X^T X)^{-1} X^T y}_{H \in \mathbb{R}^{N \times N}}, H = X(X^T X)^{-1} X^T$$

즉 \hat{y} : orthogonal projection of y onto X 의 column space

- Projection review

- orthogonal projection in line, 결론 : $P=uu^T$
- orthonormal basis(u_1, \dots, u_k), space U , A 's column is u_1, \dots, u_k , 결론 : $P=AAT$
- non-orthonormal basis, 결론 : $P=A(ATA)^{-1}AT$.
 $(ATA)^{-1}$ is normalizing factor, $Py=\hat{y}$ P=A(ATA)^{-1}AT
normalizing factor

- 참고 H

$H^T = H$, $H^2 = H$, $1-H$ 역시 projection이고 H 와는 orthogonal space
 $\text{trace}(H)=d+1$, is sum for diagonal elements in H

- 정리

input x 로 실수인 y 를 예측한다. 이 때, best x hyper plane 만들자. Error는 OLS로.

- 1) 방법 1 : one-shot : $w = X^T y = (XTX)^{-1} X^T y$
- 2) 방법 2 : hat matrix, 투영과 관찰을 통해 차이가 최소가 되도록

generalization bound : $E_{\text{out}} = E_{\text{in}} + O(d/N)$

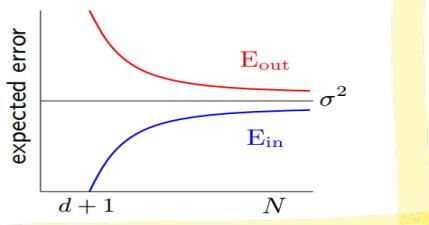
if we assume $y=f+e$, (noise e has 0 mean and finite variance,

we can also derive exact formulas for E_{in} and E_{out}

어디 어떻게 가능?

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

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



(E_{in} 은 in-sample noise 까지 같이 학습하기 때문에 수렴값보다 작다)

(E_{out} 은 w 가 fitted in in-sample error 이므로 수렴값보다 크다)

: 결론 : expected generalization error : $O(2(d+1)/N)$ 이므로 learning is feasible

Summary

- linear models

- ▶ used for classification, regression, probability estimation
- ▶ have small d_{VC} and generalize well \Rightarrow first practical choice
- ▶ use the 'signal': $\sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

- linear classification

$$\mathcal{Y} = \{-1, +1\}$$

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

$$e(\hat{y}, y) = \|\hat{y} - y\|$$

- ▶ NP-hard in general

- linear regression

$$\mathcal{Y} = \mathbb{R}$$

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

$$e(\hat{y}, y) = (\hat{y} - y)^2$$

- ▶ efficient closed-form solution

- linear regression algorithm (1-step learning): $\boxed{\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}}$

Lecture4. Logistic regression

Outline

Logistic Regression

- Predicting Probability
- Cross-Entropy Error Measure
- Training via Gradient Descent
- Logistic Regression Algorithm

Summary

Appendix: Entropy

- Logistic regression

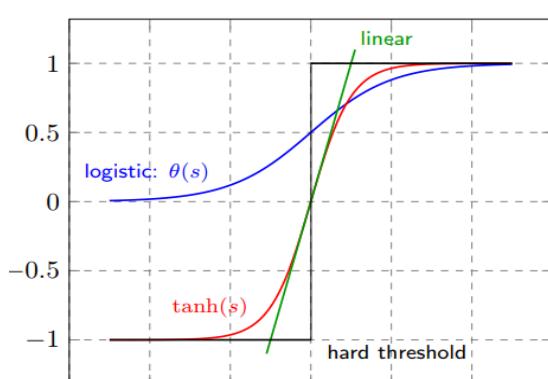
input : $h(x) = \theta(w^T x)$, θ is logistic function, sigmoid, soft threshold

output : real but bounded

θ 가 어떤 점에서 analytical and computational advantage 를 주는지 생각하면서 공부

θ 의 특징 : $1 - \theta(s) = \theta(-s)$

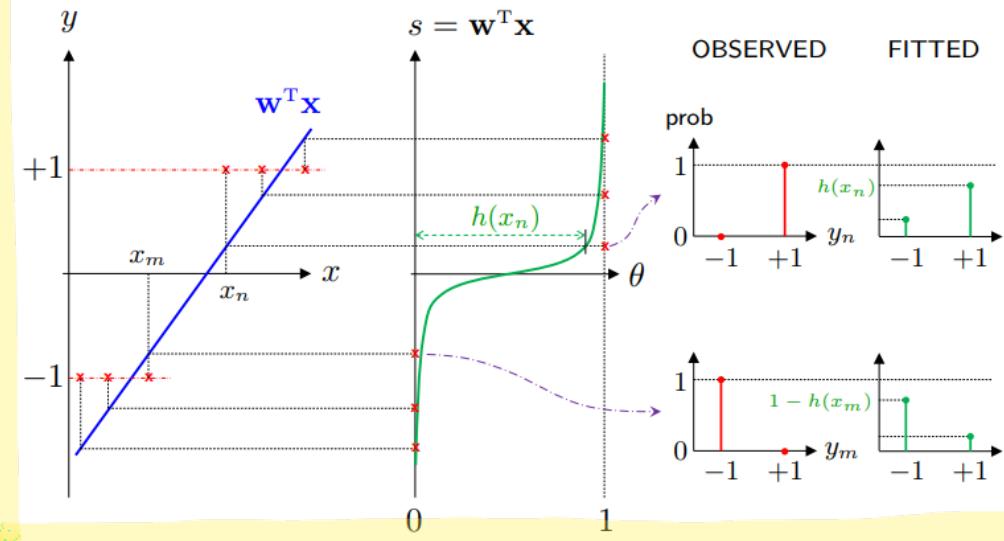
- Sigmoid function comparison



$$\theta(s) = \frac{1}{1+e^{-s}}$$

→ sigmoid function.

● Big picture



$$\mathbb{P}[y = +1 | \mathbf{x}] = f(\mathbf{x}) \sim h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$$

$$P(y|\mathbf{x}) = h(\mathbf{x})^{[y=+1]} (1-h(\mathbf{x}))^{[y=-1]} = \theta(y\mathbf{w}^T \mathbf{x})$$

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

$$h(-x) = \theta(-w^T x) = 1 - \theta(w^T x) = 1 - h(x).$$

<정의>

$$\begin{cases} h(x) = \theta(w^T x), \\ P(y|x) = \theta(yw^T x) \end{cases}$$

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

defn. $P(y|x) = \begin{cases} h(x) & \text{if } y=1 \\ 1-h(x) & \text{if } y=-1 \end{cases}$

$$\therefore P(y|x) = h(x)^{[y=1]} + (1 - h(x))^{[y=-1]}$$

$$= h(yx) = \theta(yw^T x)$$

$$= \frac{1}{1 + e^{-yw^T x}} \quad \therefore \text{성립.}$$

$$\begin{array}{ccc} & \swarrow & \searrow \\ y=1. & & y=-1. \\ \frac{1}{1 + e^{-w^T x}} & & \frac{1}{1 + e^{w^T x}} \end{array}$$

$$= 1 - \frac{1}{1 + e^{-w^T x}} = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$$

$$= \frac{1}{1 + e^{w^T x}}$$

Logistic regression 의 Error

- Learning target

$f(x) = P[y=+1 | x]$, $y=+1$ 일 확률, noisy label (단지 1 또는 -1 뿐)

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

$$= 1-f(x) \text{ for } y=-1$$

logistic regression에서
최적화할 수 있는 새로운 협력의 에러 필요.

- Error measure

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N \left(h(\mathbf{x}_n) - \frac{1}{2}(1+y_n) \right)^2$$

1) simple error measure

CROSS-entropy 에러 사용.
but very hard to minimize

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

2) cross-entropy measure, (error property and convex)

- based do intuitive probabilistic interpretation

- gradient-based optimization에서 좋은 성질을 가진다.

- Error measure 유도 방법 1: likelihood

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

(0 때, $1-\theta(x)=\theta(-x)$ 임을 이용)

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

likelihood of training data :

maximum likelihood \mathbf{h} 를 select

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

결론 : in-sample error :

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

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

consider two pmf $\{p, 1-p\}$ and $\{q, 1-q\}$ with binary outcomes

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

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

$$\begin{aligned} p &= \{y_n=1\} \rightarrow \frac{1}{2} \\ q &= h(\mathbf{x}_n) \rightarrow \frac{1}{2} \\ &: (y_n=1) \ln \frac{1}{h(\mathbf{x}_n)} \\ &+ (y_n=0) \ln \frac{1}{1-h(\mathbf{x}_n)} \end{aligned}$$

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

- $\frac{1}{N} \log \{\prod P(y|x)\}$.
 $= \frac{1}{N} \sum (\text{cross})$.
△. 같은 cross error를
상수로 치환.

: cross entropy error function

nth term is cross entropy for two pmfs measured on (\mathbf{x}_n, y)

observed : { $[y=+1]$, $1-[y=+1]$ }

fitted: { $h(\mathbf{x}_n)$, $1-h(\mathbf{x}_n)$ }

이 때, $p=[y=+1]$, $q=h(\mathbf{x}_n)$ 인 cross entropy

유도 과정에서 $1-\theta(x)=\theta(-x)$ 만 이용하여도 전개가 가능하다. Ein 은 다른 식으로 표현 돼.

$$\begin{aligned}
&\ast P(y|x) \\
&= h(\mathbf{x}_n)^{[y=+1]} + (1-h(\mathbf{x}_n))^{[y=-1]} \\
&\therefore \log \prod P \text{를 하면} \\
&\quad \sum [y=+1] \left(\log \frac{1}{h(\mathbf{x}_n)} \right) + \dots
\end{aligned}$$

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

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

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

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

no analytic solution \Rightarrow iterative solution

"gradient descent" : smooth 한 error measure 를 가지는 다양한 것을 train 가능

Iterative solution.

- Gradient descent : good at logistic regression

general for minimizing twice-differentiable function : logistic and neural network

initial location is crucial

- 1) 어디로? Negative gradient 방향, 2) 얼만큼? η learning rate 만큼

$$\Delta E_{in} = E_{in}(w(t+1)) - E_{in}(w(t))$$

$$= E_{in}(w(t) + \eta \hat{\mathbf{v}}) - E_{in}(w(t))$$

$$\approx \eta \nabla E_{in}(w(t))^T \hat{\mathbf{v}}$$

$$\geq -\eta \|\nabla E_{in}(w(t))\|$$

$$\hat{\mathbf{v}} = -\frac{\nabla E_{in}(w(t))}{\|\nabla E_{in}(w(t))\|}$$

ΔE_{in} 이

최적화 핵심 바탕자.

- Algorithm : gradient descent

for (iteratively)

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

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

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

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

$v = -\nabla Z_m(w(t))$ を 使う
이 때 $w(t+1) = w(t) - \eta v$
이를 여러번 반복.

- Gradient descent

- approximately minimize E_{in} (convex 이면 정확하게)

subset. 1회.

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

- initialization 은 normal distribution 에서 choose

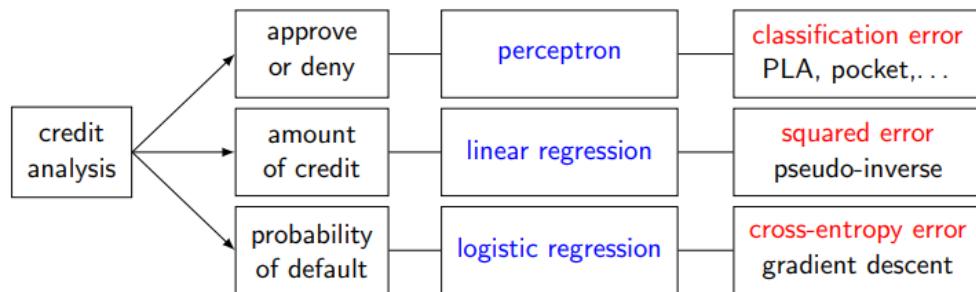
- when to terminate?

1) upper bound of large iterations

2) small size of gradient

3) error itself is small

Linear model review



(사실, 엄밀한 Perceptron 은 linear classification에서 쓰이는 unit을 이야기한다)

	linear classification	linear regression	logistic regression
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 $(\hat{y} - y)^2$	<u>cross-entropy error</u> $\llbracket y=+1 \rrbracket \ln \frac{1}{\hat{y}} + \llbracket y=-1 \rrbracket \ln \frac{1}{1-\hat{y}}$
$E_{\text{in}}(h)$	$\frac{1}{N} \sum_{n=1}^N \llbracket h(\mathbf{x}_n) \neq y_n \rrbracket$	$\frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$	$\frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$
opt.	combinatorial optimization (NP-hard)	set $\nabla E_{\text{in}}(\mathbf{w}) = 0$ (closed-form solution exists)	set $\nabla E_{\text{in}}(\mathbf{w}) = 0$ iterative optimization (e.g. gradient descent)

$$\begin{aligned}
 & \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{p(y|x)} \\
 &= \frac{1}{N} \sum_{n=1}^N \left[\llbracket y=+1 \rrbracket \ln \frac{1}{h(\mathbf{x}_n)} + \llbracket y=-1 \rrbracket \ln \frac{1}{1-h(\mathbf{x}_n)} \right]
 \end{aligned}$$

Summary

- logistic regression: estimates probabilities of binary events
 - ▶ approximates $\mathbb{P}[y = +1|\mathbf{x}] = f(\mathbf{x})$ by $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$
 - ▶ likelihood: $P(y|\mathbf{x}) = h(\mathbf{x})^{\llbracket y=+1 \rrbracket} (1 - h(\mathbf{x}))^{\llbracket y=-1 \rrbracket} = \theta(y \mathbf{w}^T \mathbf{x})$
 - ▶ training: minimize cross-entropy error iteratively
- gradient descent: iteratively solves unconstrained optimization
 - ▶ general technique for minimizing twice-differentiable functions
 - ▶ batch gradient descent: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla E_{\text{in}}(\mathbf{w})$ [over all data]
 - ▶ stochastic gradient descent: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla e_n(\mathbf{w})$ [pointwise]
 - ▶ mini-batch gradient descent, batch normalization
- linear models should be a first try: simple/robust/work & generalize well

Lecture5. Artificial neural networks

Outline

Introduction

Multilayer Perceptrons

Representations

Network Architectures

Training Multilayer Perceptrons → 허설: 배프로퍼게이션

Back-Propagation Algorithm

Backward Phase (Output Layer)) 정통

Backward Phase (Hidden Layer)) 간접

Backprop Demystified) 간접한 방법

시험: 복잡한 네트워크
정통법 아님 간단한
방법으로 풀기

Deep Learning: Motivation

Summary

- Artificial neural network (ANN) = Multilayer perceptron (MLP)

Relu : allows spares propagation of activations and gradients

better for handling vanishing gradient problem

- 분포를 구현하기 위한 것이 함수

	K class	Binary
함수	Softmax function	Logistic function
1 번의 event	멀티눌리 distribution = categorical distribution (주사위 한번 던지는 것)	베르눌리 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

Error measures

설계하는 목표에 대한

- e_k : error signal on k -th output

$$e_k \triangleq \text{error} \quad \text{NN output} - \text{correct label}$$

$$e_k = h_k - y_k.$$

설계하는 목표에 대한

- \mathcal{E}_n : error energy on example $(x_n, y_n) \leftarrow \text{sum of squared errors}$

$$\mathcal{E}_n = \frac{1}{2} \sum_{k=1}^K e_{k,n}^2 \quad \mathcal{E}_n = \frac{1}{2} \sum_{k=1}^K e_{k,n}^2$$

- \mathcal{E}_D : mean-squared error on data $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$

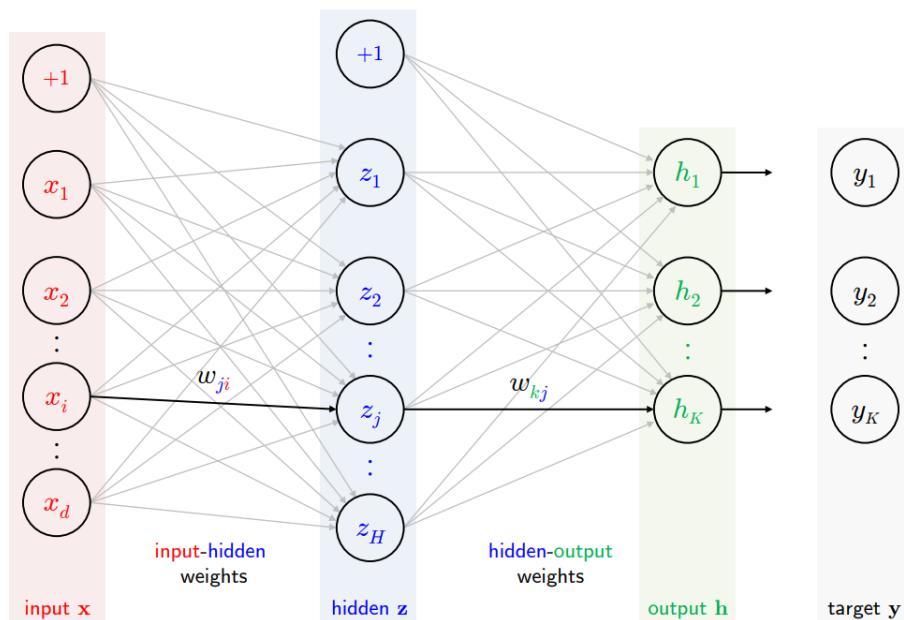
$$\mathcal{E}_D = \frac{1}{N} \sum_{n=1}^N \mathcal{E}_n = \frac{1}{2N} \sum_{n=1}^N \sum_{k=1}^K e_{k,n}^2 \quad \mathcal{E}_D = \frac{1}{2N} \sum_{n=1}^N \sum_{k=1}^K e_{k,n}^2$$

-> 우리는 \mathcal{E}_D (Ein)을 최소화하는 W 를 찾고 싶다.

- Optimization method

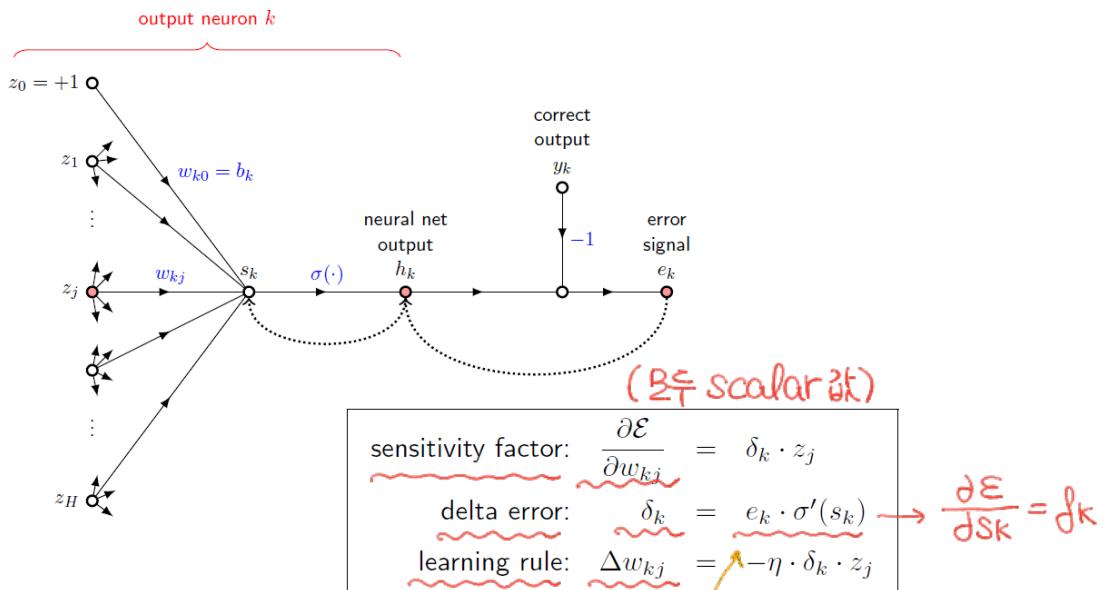
- second order method (Newton's method) but not attractive here
- first order method (gradient descent)

- Back propagation



- Back propagation in output layer

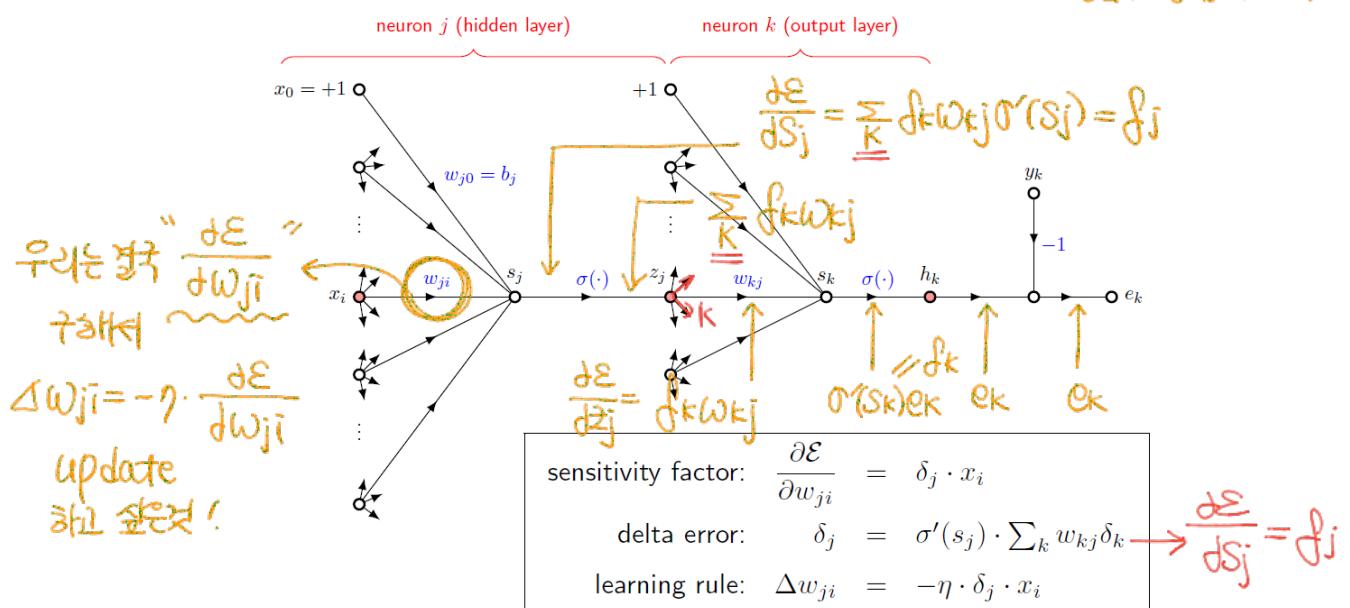
Summary: output layer



(이때, error, input 이 클수록 w 를 많이 변화시킨다)

- Back propagation in input layer

Summary: hidden layer



- Back prop Algorithm

repeat

pick n (하나의 x 를 pick)

* Back Prop 알고리즘

forward : compute all z, h
 backward : compute all delta
 update weights
 until it is time to stop, update final weights

하나의 input을 선택 ~> 여러개의 subset도 가능
 forward $\rightarrow z, h$ 계산
 backward $\rightarrow \text{delta}$ 계산
 이를 바탕으로 전부 $\frac{\partial \mathcal{E}}{\partial w_{ij}}$ 를 계산 가능.
 그래서 weight update \rightarrow 평균을 구해서 update 값 결정.
 이를 여러번 반복
 : "전체 N/subset 수" 만큼
 반복해서 W 를 update.

Back prop 간단한 계산

- 간단한 계산 공식

- 1) $\text{out} = f(\text{in})$: $f'(\text{in})$ 을 곱한다.
- 2) $\text{out} = \text{in}_1 * \text{in}_2$: in_1 에 대한 식을 구하는 경우 in_2 를 곱한다.
- 3) $\text{out} = \text{sum}(\text{in})$: 그대로 fanout 한다.
- 4) $\text{out} = \max(\text{in})$: ini 가 max 이면 fanout, 아니면 0이다.
- 5) $\text{out} = \text{in}, (\text{fanout})$: sum 으로 모두 더한다.

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

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

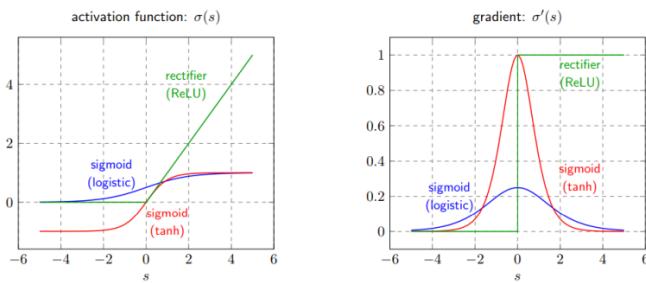
$$\delta_j = \left(\sum_k \delta_k \cdot w_{kj} \right) \cdot \sigma'(s_j)$$

Deep learning에 대한 간단한 overview

- DL: deep하게 쌓을 경우 exponential 한 무엇인가를 polynomial 에 나타내는 것이 가능해

- Vanishing gradient problem

(non-linearity 를 위해서는 양 끝단에서 작동해야 하지만 이때는 gradient 가 죽는다)



Summary

- artificial neural network: universal approximator of functions
 - neuron model: synapse (weights), adder, activation functions
 - conventional: feed-forward multilayer perceptrons (w/ backprop)
 - recent: RBM, AE, CNN, RNN, GAN and many other variants

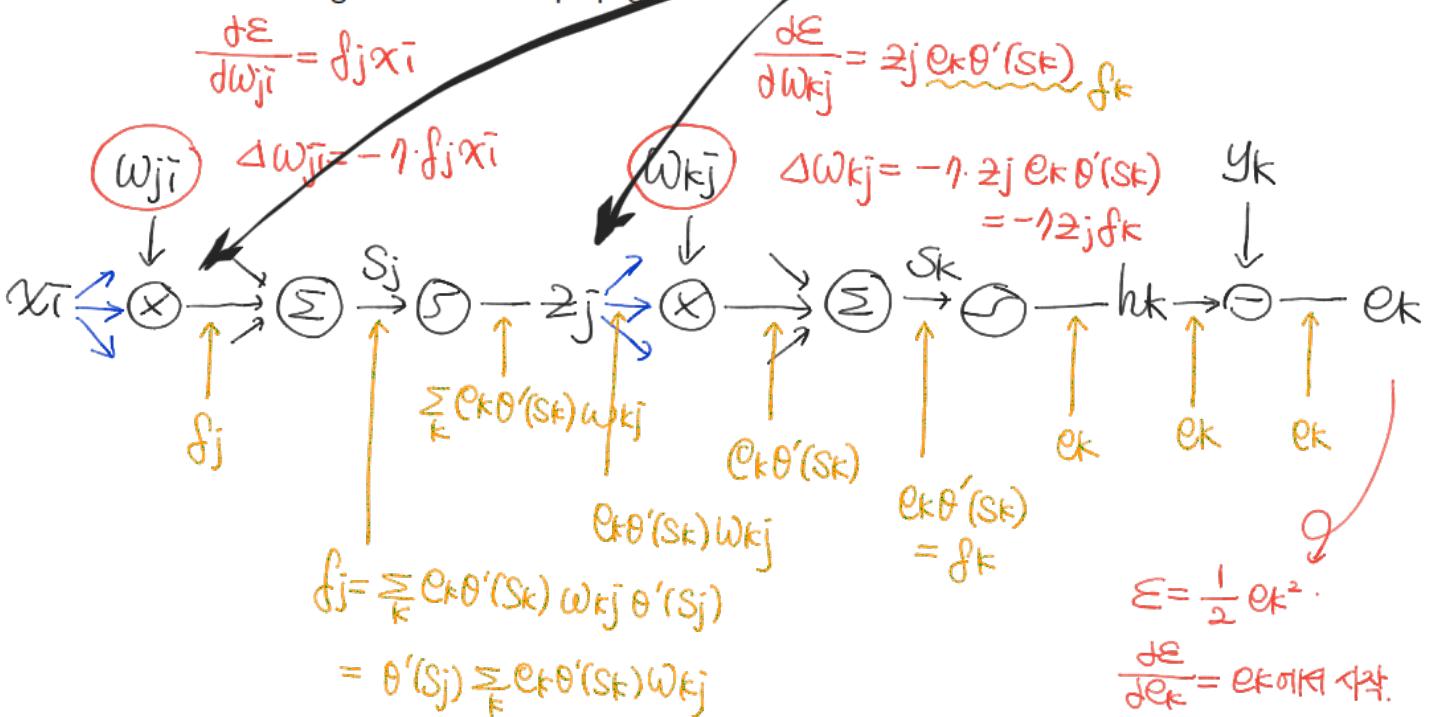
- hidden layers act as feature detectors

- # hidden layers $\uparrow \Rightarrow$ intelligence \uparrow but training difficulty \uparrow
- practical solutions exist for effective training of deep neural networks

- two types of signals for training feed-forward multilayer perceptrons

- function signals: forward propagation

- error signals: backward propagation



Lecture6. Support vector machines : SVM

Outline

Linear Discriminant Functions

Support Vector Machines

Introduction

Problem Formulation

Derivation of Maximum Margin Classifier

Summary

Appendix

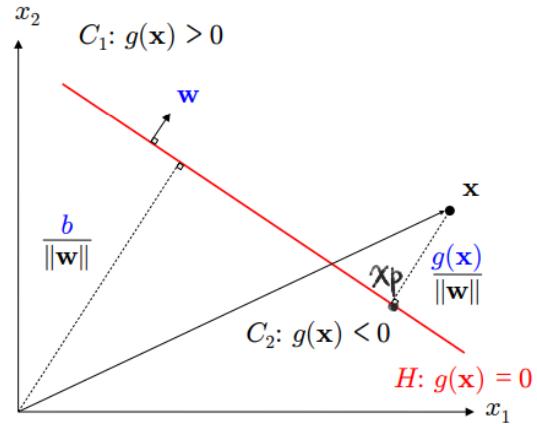
- **Discriminant function : g**
 - representing classifiers
 - divide the feature space into K decision regions (이 때, reject region 을 두기도)
 $g_1(\text{vector } x) > g_2(\text{vector } x)$: class1 로 분류
ex) $g_i(x) = p(C_i|x)$: Bayes error 를 minimize
- **Linear discriminant function : $g(x) = w^T x + b$**
지금까지는 b 를 x 에 포함하여 같이 보았지만 이제는 따로 생각한다.

2 category case : $g(x)=0$ defines decision boundary

for linear $g(x)$, this decision surface is a hyperplane

- Geometric interpretation

- big picture:



$g(\mathbf{x})$ gives distance from \mathbf{x} to H

$$\begin{aligned}
 g(\mathbf{x}) &= \mathbf{w}^\top \mathbf{x} + b \\
 &= \mathbf{w}^\top \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + b \\
 &= (\mathbf{w}^\top \mathbf{x}_p + b) + r \|\mathbf{w}\| \\
 &= r \|\mathbf{w}\|
 \end{aligned}$$

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^\top \mathbf{x} + b}{\|\mathbf{w}\|}$$

proof:

(이때, w 는 surface 의 방향을 지정, b 는 위치를 지정)

- SVM : maximum margin (hyperplane)classifier

maximum margin : best generalization performance 을 위해

앞으로의 overview

• first, assume linear separability (hard-margin SVM)

- ▶ steps A–E: derive & solve “primal” form
- ▶ steps 1–5: derive & solve “dual” form

• second) extend to non-separable cases (next lecture)

- ▶ almost linearly separable (with outliers) \Rightarrow soft-margin SVM
- ▶ linearly non-separable \Rightarrow nonlinear transform & kernel trick
→ 주제 병렬적으로 적용 가능.
(non-separable 하거나 outlier 있는)

● Optimization

<http://www.csc.kth.se/utbildning/kth/kurser/DD3364/Lectures/KKT.pdf>

constraint 가 존재하는 min/max 문제

equality : 라그랑지안, inequality : KKT

최적화에 대한 내용

1) 우위가 풀고자 하는 문제: $\max / \min f$,

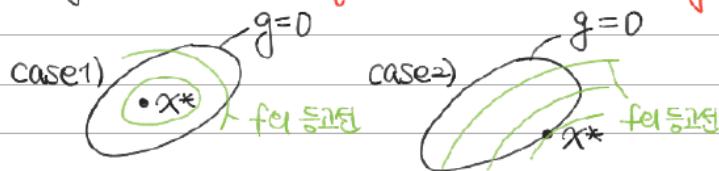
문제의 제약조건: $g \geq 0$ 또는 $g \leq 0$.

$\therefore L = f \pm \lambda g$ 을 이용해 풀어나간다.

2) $L = f \pm \lambda g$.

- ① $\lambda \geq 0$.
- ② $g \leq 0$ 또는 $g \geq 0$. 문제의 조건
- ③ $\nabla_x L = \nabla_x f \pm \lambda \nabla_x g = 0$.
- ④ $\lambda g = 0$ (이제, 둘다 0은 아님)

3) f 와 g 의 관계



$\lambda = 0$, g 는 필요x

g : inactive constraint

$\lambda > 0$, $g = 0$.

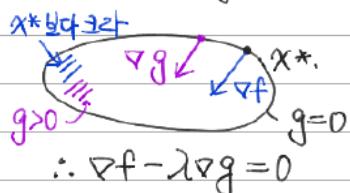
g : active constraint

4) 4가지 case.

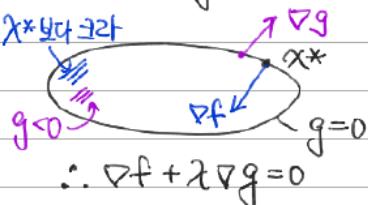
(f) \max 문제, \min 문제

(g) $g \geq 0$, $g \leq 0$.

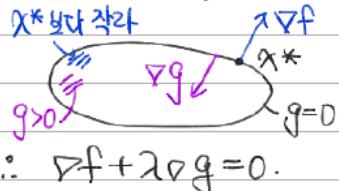
case A) $\min f$, $g \geq 0$.



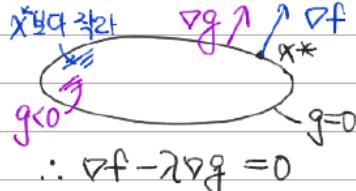
case B) $\min f$, $g \leq 0$



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



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



Primal form : 이해하기 좋다. But 의미를 파악하기 힘들다.

- Step A

$g(x+) = +1, g(-x) = -1$ 로 scale 조정

-> called support vectors

- Step B : setting constraints

$$y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq +1 \quad (y=+1, -1 \text{의 경우를 모두 포함})$$

- Step C : determining objective function

margin : $2/\|\mathbf{w}\|$

therefore, want to minimize $\|\mathbf{w}\|$ -> 편의상 제곱한 것을 minimize

- Step D : SVM optimization problem

SVM (primal) optimization problem

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to} && y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq 1, \forall t \in [1, N] \end{aligned}$$

- constrained optimization problem

convex quadratic problem with $d+1$ parameters, N constraints

- 특징

cost function sis convex function of \mathbf{w}

constraints are linear in \mathbf{w}

特点: \mathbf{w}, b :: d+1 차원
 NIH (\mathbf{x}_n, y_n)

- Step E : QP (Quadratic programming)

% QP

$$\begin{aligned} & \underset{\mathbf{z}}{\text{minimize}} && \frac{1}{2} \mathbf{z}^\top Q \mathbf{z} + \mathbf{c}^\top \mathbf{z} \\ & \text{subject to} && A\mathbf{z} \geq \mathbf{a} \end{aligned} \quad , \text{solution : } \mathbf{z}^* \leftarrow \text{QP}(Q, \mathbf{c}, A, \mathbf{a})$$

solving primal by QP

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to} && y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq 1, \forall t \in [1, N] \end{aligned}$$

$$\mathbf{z} = \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \in \mathbb{R}^{d+1}$$

$$\mathbf{w}^\top \mathbf{w} = [b \quad \mathbf{w}^\top] \begin{bmatrix} 0 & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} = \mathbf{z}^\top \begin{bmatrix} 0 & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} \mathbf{z} \Rightarrow Q = \begin{bmatrix} 0 & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix}, \mathbf{c} = \mathbf{0}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \equiv [y_n \quad y_n \mathbf{x}_n^\top] \mathbf{z} \geq 1 \Rightarrow \begin{bmatrix} y_1 & y_1 \mathbf{x}_1^\top \\ \vdots & \vdots \\ y_N & y_N \mathbf{x}_N^\top \end{bmatrix} \mathbf{z} \geq \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

이제, N개의 조건을 한번에 고려.

여기서 QP solver에 넣을 값을 지정해 줄 수 있어야 한다.

- Issues on primal problem

단점 1)

d가 무한으로 가는 경우, 문제를 푸는 것이 불가능하다.

|w|의 제곱 계산량이 너무 많다.

단점 2)

constraint가 너무 많다. (N개)

단점 3)

insight를 주지 못한다. Support vector의 의미를 찾을 수 없다.

$$A = \begin{bmatrix} y_1 & y_1 x_{11} & \dots & y_1 x_{1d} \\ \vdots & \vdots & & \vdots \\ y_N & y_N x_{N1} & \dots & y_N x_{Nd} \end{bmatrix} = \begin{bmatrix} y_1 & y_1 x_{1T} \\ \vdots & \vdots \\ y_N & y_N x_{NT} \end{bmatrix}_N^{d+1}$$

* QP에 SVM 대입.

만약 \mathbf{z} 를 다르게

QP식: $\frac{1}{2} \mathbf{z}^\top Q \mathbf{z} + \mathbf{c}^\top \mathbf{z}$ 를 최소화.

설정한다면 Q.C.A.a가

조건: $A \mathbf{z} \geq \mathbf{a}$.

어떻게 변환하는지 보자.

SVM primal식: $\frac{1}{2} \mathbf{w}^\top \mathbf{w}$ 를 최소화.

조건: $y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq 1, \forall t \in [1, N]$.

* $\mathbf{z} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$ 라 하자.

$$\Rightarrow N \begin{bmatrix} y_1 x_{11} & \dots & y_1 x_{1d} & y_1 \\ \vdots & & \vdots & \vdots \\ y_N x_{N1} & \dots & y_N x_{Nd} & y_N \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \geq \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq 1$.

$$i) [\mathbf{w} \ b] \underbrace{\begin{bmatrix} I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}}_{[\mathbf{w}]} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \mathbf{w}^\top \mathbf{w}$$

$$\therefore A = \begin{bmatrix} y_1 x_{1T} & y_1 \\ \vdots & \vdots \\ y_N x_{NT} & y_N \end{bmatrix}, a = \mathbf{1}$$

$$\therefore Q = \begin{bmatrix} I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, C = \mathbf{0}$$

Dual form : 이해하기 힘들다. But support vector, kernel trick에 대한 이해 가능

- Dual formulation

optimization theory
if cost, constraint functions are strictly convex, dual form is exist

and solving dual form is equivalent to solving primal

- Step 1 : compute Lagrangian

주어진 primal form :

$$\begin{aligned} \text{minimize}_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \quad & y_t (\mathbf{w}^\top \mathbf{x}_t + b) \geq 1, \quad \forall t \in [1, N] \end{aligned}$$

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

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

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

-> solution 은 saddle point for Lagrangian

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

- Step 2 : KKT conditions

KKT condition 을 이용한다.

(nonlinear programming 이 optimal 하기 위한 필요 조건이다)

- w,b 를 a로 대체하기 위해 사용한다.

- Step1 의 식을 dual form 으로 변환할 수 있다.

Lagrangian duality

Constrained **minimum** of convex programming problem can also be obtained as **maximization** task applied on Lagrangian.

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: $\underbrace{\alpha_t^* [y_t(\mathbf{w}^{*\top} \mathbf{x}_t + b^*) - 1]}_{\alpha_t^* = 0 = 0} = 0, \forall t$ } $\lambda g = 0$
 $\Rightarrow \alpha_t^* = 0$ $\Rightarrow \mathbf{g}$ 가 힙不排除. 즉 boundary에 존재해. \Rightarrow support vector

KKT condition 을 Lagrangian 식에 대입

(condition 1,2)

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

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

(insight 1 : \mathbf{w} 를 a 의 식으로 표현 가능, \mathbf{w} 는 \mathbf{x} 와 y 의 weighted sum)

(insight 2 : ay 의 비율 합은 0)

(condition3)

condition3 에서, 둘 중 하나만 0 이다. (둘 다 0 이 될 수 없다)

$$\alpha_t > 0 \Rightarrow y_t(\mathbf{w}^{*\top} \mathbf{x}_t + b) = 1 \rightarrow \alpha_t > 0. 즉 S.V이면$$

$$y_t(\mathbf{w}^{*\top} \mathbf{x}_t + b) > 1 \Rightarrow \alpha_t = 0$$

active constraint

(만일, a 가 0 이 아니면, active(=) constraint $y_t(\mathbf{w}^{*\top} \mathbf{x}_t + b) \geq 1, \forall t \in [1, N]$ 을 의미,

즉 constraint에서 =등호가 성립한다는 것을 의미. 경계에 존재)

-> Lagrangian multipliers a 는 0 또는 양수

※: 서포트 베터수

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

$N_s \ll N$

weight 정할때 S.V만 쓰면 된다!

유한한 적은 수의 ($\mathbf{x}_t, y_t, \alpha_t$)로 \mathbf{w} 를 구할수 있다.

-> dual 을 사용하니 s.v.의 intuition 이 살아난다.

-> primal 에서는 \mathbf{w} 가 d 차원이므로 d 가 무한대이면 문제를 풀 수 없었다.

- Step 3 : formulate dual problem

기존의 primal form 의 Lagrangian :

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

식을 조금 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 form 을 위의 KKT condition 에서 얻은 3 가지 condition 을 이용해 변형한다.

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha_t y_t \mathbf{w}^\top \mathbf{x}_t + \sum_{t=1}^N \alpha_t \\ &= \frac{1}{2} \left\| \sum_{t=1}^N \alpha_t y_t \mathbf{x}_t \right\|^2 - \sum_{t=1}^N \alpha_t y_t \left(\sum_{s=1}^N \alpha_s y_s \mathbf{x}_s \right)^\top \mathbf{x}_t + \sum_{t=1}^N \alpha_t \\ &= -\frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \mathbf{x}_s^\top \mathbf{x}_t + \sum_{t=1}^N \alpha_t \end{aligned}$$

SVM (dual) optimization problem

이제 max 하는 문제로 !!!

$\underset{\alpha}{\text{maximize}}$

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

subject to

$$\sum_{t=0}^N \alpha_t y_t = 0 \quad \leftarrow \text{(나머지는 문제에서 이미 포함되어 있음.)}$$

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

SVM (primal) optimization problem

dual variable에 대한 조건이 추가됨. $0 \leq \alpha_t \leq 1, \forall t \in [1, N]$

$$\begin{aligned} \underset{\mathbf{w}}{\text{minimize}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \quad & y_t (\mathbf{w}^\top \mathbf{x}_t + b) \geq 1, \forall t \in [1, N] \quad (\text{to compare}) \end{aligned}$$

primal 이 w 에 대한 식이었다면 dual 은 a 에 대한 식

(첫 번째 식에서, $X\mathbf{X}$ 를 제외하면 d 에 대한 식이 아닌, N 에 대한 식으로,

d 가 무한으로 가더라도 풀 수 있다.

뒤에서 $(X\mathbf{X})$ 가 d 차원과 관련이 있지만, 내적의 형태이므로 kernel trick 이용 가능)

$$\underset{\mathbf{z}}{\text{minimize}} \quad \frac{1}{2} \mathbf{z}^\top Q \mathbf{z} + \mathbf{c}^\top \mathbf{z}$$

이제, QP solver 를 이용해 문제를 풀 수 있다.

$$\text{subject to} \quad A\mathbf{z} \geq \mathbf{a}$$

(d 에 scale 하지 않고 N 에 scale. but 여전히 large N 에 대한 문제가 있다)

dual 의 중요한 성질

- 1) 오직 terms of training data
- 2) function depends only on input ad inner products $\mathbf{x}^\top \mathbf{x}$ -> can use kernel trick
따라서, nonlinearly separable case도 고려할 수 있다.
- 3) SVM에서, primal and dual is equivalent

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

dual에서 \mathbf{a} 는 \mathbf{w} 로 directly 변환될 수 있다.

● Step 4 : solve dual problem as QP

- QP solver ~ 등호조건이 1개 추가된 QP solver를 이용

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \mathbf{c}^\top \mathbf{x} \quad (\text{d1} \vdash \mathbf{x}) \\ \text{subject to} & \underline{A} \mathbf{x} \leq \underline{b} \\ & \underline{E} \mathbf{x} = \underline{d} \\ & \maximizex \text{가 } -1\text{배수} \\ & \text{통해 minimizex 문제!} \\ \text{- Dual form} & \maximize{\alpha} \\ & \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \langle \mathbf{x}_s, \mathbf{x}_t \rangle \\ & \sum_{t=1}^N \alpha_t y_t = 0 \\ \text{subject to} & \alpha_t \geq 0, \forall t \in [1, N] \\ & \sum_{t=1}^N \alpha_t y_t = 0 \end{array}$$

$\therefore \mathbf{C} = -\mathbf{1} \quad \mathbf{x}_i^\top \mathbf{x}_j$
 $\therefore Q = \begin{bmatrix} +y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \end{bmatrix}$
 $\therefore A = -I_N, b = 0_N$
 $\therefore E = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}^\top = [y_1 \dots y_N], d = 0$
 이때 E 는 vector, d 는 scalar.

- Dual form을 QP에 적용하면

$$\begin{array}{ll} \underset{\alpha}{\text{minimize}} & \frac{1}{2} \alpha^\top Q \alpha - \mathbf{1}^\top \alpha \\ \text{subject to} & -\alpha \leq 0 \\ & \mathbf{y}^\top \alpha = 0 \end{array}$$

$$Q = \begin{bmatrix} y_1 y_1 \mathbf{x}_1^\top \mathbf{x}_1 & y_1 y_2 \mathbf{x}_1^\top \mathbf{x}_2 & \cdots & y_1 y_N \mathbf{x}_1^\top \mathbf{x}_N \\ y_2 y_1 \mathbf{x}_2^\top \mathbf{x}_1 & y_2 y_2 \mathbf{x}_2^\top \mathbf{x}_2 & \cdots & y_2 y_N \mathbf{x}_2^\top \mathbf{x}_N \\ \cdots & \cdots & \cdots & \cdots \\ y_N y_1 \mathbf{x}_N^\top \mathbf{x}_1 & y_N y_2 \mathbf{x}_N^\top \mathbf{x}_2 & \cdots & y_N y_N \mathbf{x}_N^\top \mathbf{x}_N \end{bmatrix}$$

- final form

• final form

$$\underset{\alpha}{\text{minimize}} \frac{1}{2} \alpha^T \underbrace{\begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \cdots & y_1 y_N x_1^T x_N \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 & \cdots & y_2 y_N x_2^T x_N \\ \cdots & \cdots & \cdots & \cdots \\ y_N y_1 x_N^T x_1 & y_N y_2 x_N^T x_2 & \cdots & y_N y_N x_N^T x_N \end{bmatrix}}_{\text{quadratic coefficients } Q} \alpha + \underbrace{(-1^T) \alpha}_{\text{linear coefficients } C}$$

subject to $\underbrace{y^T \alpha = 0}_{\text{linear constraint}}$

$$\underbrace{0}_{\text{lower bounds}} \leq \alpha \leq \underbrace{\infty}_{\text{upper bounds}}$$

↑ 허용 soft margin
SVM 할 때 부여하고.

- practical considerations

Q 는 (N, N) 의 dense 한 matrix, large N 에서는 solution 을 찾기 힘들다.

worst case 의 경우, $O(N^3)$ time, $O(N^2)$ space

- QP solver 는 return solution vector α

대부분 (s.v.가 아닌 경우) $\alpha=0$ 이므로 매우 sparse 한 vector α

● Step 5 : complete

final solution

α^* = QP 에서 나온 solution

w^* = α^* 와 training data 로 계산한 결과

$$b^* = -\frac{1}{2} (w^{*\top} x_+ + w^{*\top} x_-) \quad \begin{aligned} w^{*\top} x_+ + b &= +1 \\ w^{*\top} x_- + b &= -1 \end{aligned}$$

b^* = for any s.v. x_+ and x_- ,

SVM classifier $g(x)$: linearly separable case

$$g(x) = \text{sign} \left(\sum_{t=1}^N \alpha_t^* y_t x_t^T x + b^* \right)$$

$$\text{sign}(w^{*\top} x + b^*)$$

$$w^* = \sum_{t=1}^N \alpha_t^* y_t x_t$$

generalization bound

generalization bound of SVM (Vapnik, 1995):

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

E_{out} 매우 좋다.

- Issues on dual problem

장점 1)

d 가 아니라 N 에 대해서 scale, x 의 inner product

-> 따라서 high d (심지어 무한 dimension)도 handle 가능 \rightarrow non-linear

장점 2)

SVM 의 key properties 의미를 충분히 전달한다.

Summary

- support vector machines: max-margin linear discriminant

▶ larger margin \Rightarrow lower d_{VC} \Rightarrow better generalization

▶ # support vectors \approx effective # parameters $\approx d_{VC}$

▶ typically much fewer support vectors than data points

▶ one of the best off-the-shelf classification algorithms

내중에 VC dimension
다 배우고 나서 다시보기.

- derivation of support vector machines

▶ formulated in mathematical optimization framework

▶ primal/dual forms can be solved by quadratic programming

▶ dual form: allows more efficient solution (also more insightful)

Lecture7. Support vector machines (linearly non-separable)

Outline

Kernel Methods

- Nonlinear Transforms for SVM
- Kernel Trick
- Examples of SVM Kernels
- Mercer's Theorem

Soft-Margin SVM

- Handling Outliers and Noise
- Derivation of Soft-Margin SVM

Summary

Linearly non-separable case1) kernel trick 으로 해결

- Non-linear transform

feature space Z, feature vector z

$\langle z, z' \rangle$ 을 $\langle x, x' \rangle$ 으로 나타낼 수 있다. Z의 차원이 무한이어도 계산 가능하다.

$$\begin{aligned} \mathcal{L}(\alpha) &= \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \mathbf{x}_s^\top \mathbf{x}_t \\ &\downarrow \\ \mathcal{L}(\alpha) &= \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{s=1}^N \sum_{t=1}^N \alpha_s \alpha_t y_s y_t \mathbf{z}_s^\top \mathbf{z}_t \end{aligned}$$

$$\begin{aligned} \mathbf{z}^\top \mathbf{z}' &= x_1 x_1' x_1' x_1 \\ &\quad + x_1 x_2 x_1' x_2' + x_2 x_1 x_2' x_1' \\ &\quad + x_2 x_2 x_2' x_2' \end{aligned}$$

$$\begin{aligned} \mathbf{z}_s^\top \mathbf{z}_t &= K(x_s, x_t) \end{aligned}$$

1) 문제를 푸는 cost 는 x, z 에서 동일하다.

2) generalization bound 는 d 가 아니라 #s.v.이므로 d 가 커져도 변하지 않는다.

- Kernel trick (or kernel substitution)

-> inner product 만을 이용하는 어떤 algorithm 에서도 적용 가능하다.

Kernel : inner product in feature space, $z^\top z' = K(x, x')$,

Kernel trick : feature space 로 가지 않고 input space 에서 kernel 을 계산할 수 있다.

K(x, x') : x, x' 또는 z, z' 이 얼마나 similar 한가, that is nonlinear similarity measure

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

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{a}\mathbf{x}^\top \mathbf{x}' + b)^q$$

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

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

- Gaussian RBF:

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

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

$$K(\mathbf{x}, \mathbf{x}') = \tanh(\beta \mathbf{x}^\top \mathbf{x}' + \gamma)$$

- multiquadric:

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

- wave:

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

- equivalent kernel : polynomial 에서 $a=b=0$

- linear kernel : polynomial 에서 $q=1$

- RBF kernel 은 infinite feature dimension, r 이 너무 크면 over fitting 일어난다.

- 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

: 1) symmetric 2) positive semi-definite

kernel matrix

m finite points $\{x_1, \dots, x_m\}$

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

$K_{ij} = K(x_i, x_j)$ 가 성질을 만족.

이 때, 어떤 finite set

$\{x_1, \dots, x_m\}$ 에 대해서도 성립해야 한다.

positive semi-definite (PSD) ≥ 0

all eigen values, determinant is non-negative

real symmetric K 가 $t^\top K t$ 는 0 이상 for all vector t

$t^\top K t \geq 0$ for $\forall t$

Linearly non-separable case 2) soft-margin SVM

- Reformulating SVM problem : primal form

primal form (C : positive constant, C 가 클수록 hard margin에 가까워)

위와 primal form에서 $C \leq \xi_t$ 를 해.
또한 조건식 변형 추가

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^N \xi_t \\ & \text{subject to} && y_t(\mathbf{w}^\top \mathbf{x}_t + b) \geq 1 - \xi_t, \quad \forall t \\ & && \xi_t \geq 0, \quad \forall t \end{aligned}$$

primal variable: (\mathbf{w}, b, ξ_t)

- Reformulating SVM problem: dual form 으로 변형하는 것

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

$$= \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{t=1}^N \xi_t - \underbrace{\sum_{t=1}^N \alpha_t (y_t(\mathbf{w}^\top \mathbf{x}_t + b) - 1 + \xi_t)}_{\text{dual objective}} - \sum_{t=1}^N \beta_t \xi_t$$

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

(primal variable에 대해서는 minimum, dual variable에 대해서는 maximum)

- KKT conditions are obtained as follows:

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

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

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

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

$$\beta_t \xi_t = 0 \Rightarrow \beta_t = 0 \vee \xi_t = 0$$

Final form, 아래의 식을 QP에 적용

final solution

Soft-margin SVM classifier

$$\begin{aligned} & \text{maximize} && \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ & \text{subject to} && \sum_{t=1}^N \alpha_t y_t = 0 \\ & && 0 \leq \alpha_t \leq C, \quad \forall t \end{aligned}$$

only change

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

where

$$\begin{aligned} \mathbf{w} &= \sum_{0 < \alpha_t \leq C} \alpha_t y_t \mathbf{x}_t \\ b &= y_{sv} - \sum_{0 < \alpha_t \leq C} \alpha_t y_t \mathbf{x}_t^\top \mathbf{x}_{sv} \end{aligned}$$

- solution \mathbf{a} 를 이용해 optimal \mathbf{w} 를 구한다. 이후 from any s.v.에서 b 를 구한다.

- ξ_t 는 식에 나타나지 않지만 C 의 형태로 존재한다.

- $\xi_t > 0$ 인 vectors (linearly separable에서는 없었던) :

outliers.
 margin 내에 있거나,
 잘못 분류된

w = sum (ayx)에서 a=C 이므로 w에서 비중이 크다.

- 모두 정리

	모범생	원래 S.V	군사분계선	간첩
$y_t(\mathbf{w}^\top \mathbf{x}_t + b)$	$y(\mathbf{w}^\top \mathbf{x} + b) > 1$	$y(\mathbf{w}^\top \mathbf{x} + b) = 1$	$0 \leq y(\mathbf{w}^\top \mathbf{x} + b) < 1$	$y(\mathbf{w}^\top \mathbf{x} + b) < 0$
ξ_t	$\xi_t = 0$	$\xi_t = 0$	$0 < \xi_t \leq 1$	$1 < \xi_t$
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

Summary

- SVM advantages:

- systematic implementation through quadratic programming \Rightarrow very efficient implementations exist
- excellent data-dependent generalization bounds exist
- regularization built into cost function
- statistical performance independent of dim of feature space

매우 잘 개발된
QP solver 사용.

이며Independent.

- drawbacks:

- treatment of non-separable case somewhat heuristic
- number of support vectors may depend strongly on
 - kernel type and hyper-parameters
- systematic choice of kernels is difficult (need prior information)
- optimization may require clever heuristics for large problems

어때?

*S.V는 어떤 커널, 어떤 hyper 파라미터에 따라 달라진다.