

Mathematical notations

\mathbf{x} : column vector

e.g. $\mathbf{x} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \in \mathbb{R}^3$

\mathbf{x}^T : transpose of \mathbf{x} (row vector) $\mathbf{x}^T = (1, 2, 3)$

D : \mathbf{x} of dim (features) i th component

N : \mathbf{x} of samples (data), n th object

\mathbf{X} : data matrix with $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix}$$

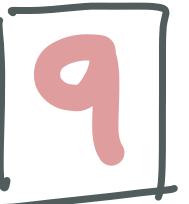
$N \times D$ matrix

Chapter 1.

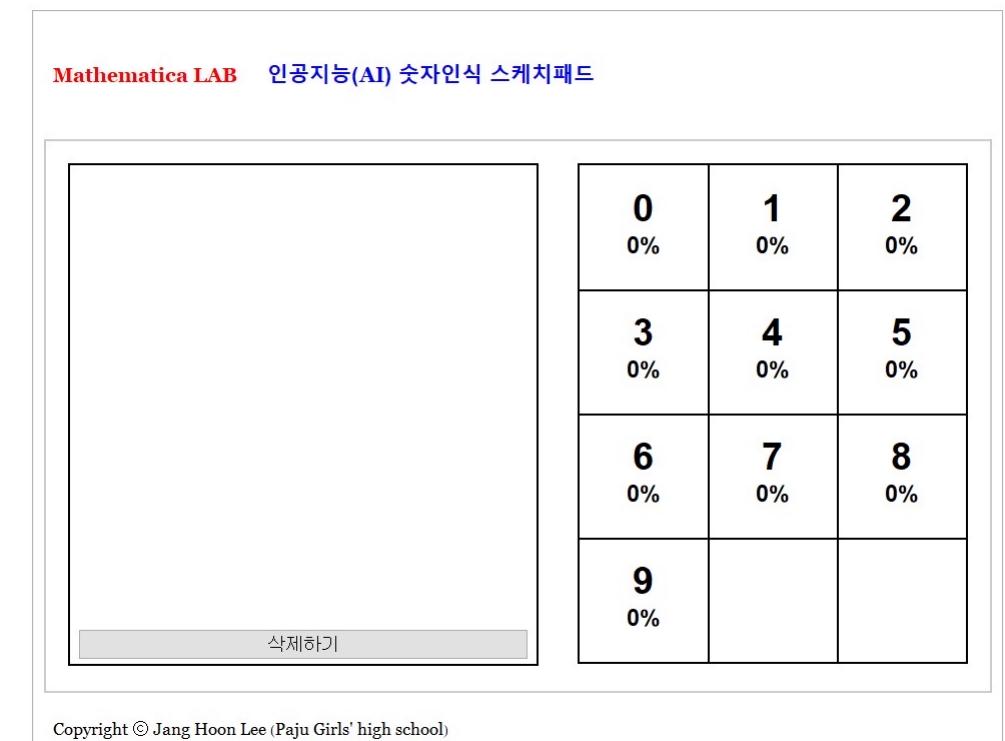
Pattern recognition

- Finding and analyzing features of data
- Decision making through prediction and classification.

Example hand - written image
28 x 28 pixels
784 - dim vector



number 9



Training

28×28 image

input vector (set) $\{x_1, x_2, \dots, x_N\}$ N 개의 손글씨 image

target vector (set) $t = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$ N -dim also called label set

$y(x)$: output of ML algorithm

e.g. mean of x , prediction of t (target)

training phase : process of formulation of $y(x)$

test set : data not used for training

- results of test set are important because we need to generalize our algorithm

Preprocessing

- null value
- transformation, scale conversation, normalization
- data type change
- feature extraction (특징 추출)
- dimensionality reduction (차원 축소)

Supervised learning (task) 지도학습

training data : input vectors with its target value
 x_1, \dots, x_N t_1, t_2, \dots, t_N

Find an expression (function) describing input and target
 $(y(x) \sim \#)$

- classification : discrete target value
- regression : continuous

Unsupervised learning

비지도 학습

training data : only input vectors x_1, \dots, x_N

pattern recognition of input (feature analysis)

- clustering (군집화)
- density estimation (밀도(분포) 추정) $P(x)$
- visualization (dimensionality reduction)

1.1 Example: Polynomial curve fitting

input $x \in \mathbb{R}$

target $t \in \mathbb{R}$

N observations $\mathbf{x} := (x_1, x_2, \dots, x_N)^T$ input vector
 $\mathbf{t} := (t_1, t_2, \dots, t_N)^T$ target vector

generated by $\sin(2\pi x)$ with small noise (perturbation)

$$x_1 = 0, x_N = 1, \quad x_{j+1} - x_j = \frac{1}{N-1} \quad \forall j = 1, \dots, N-1$$

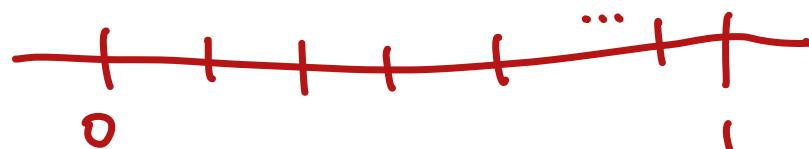
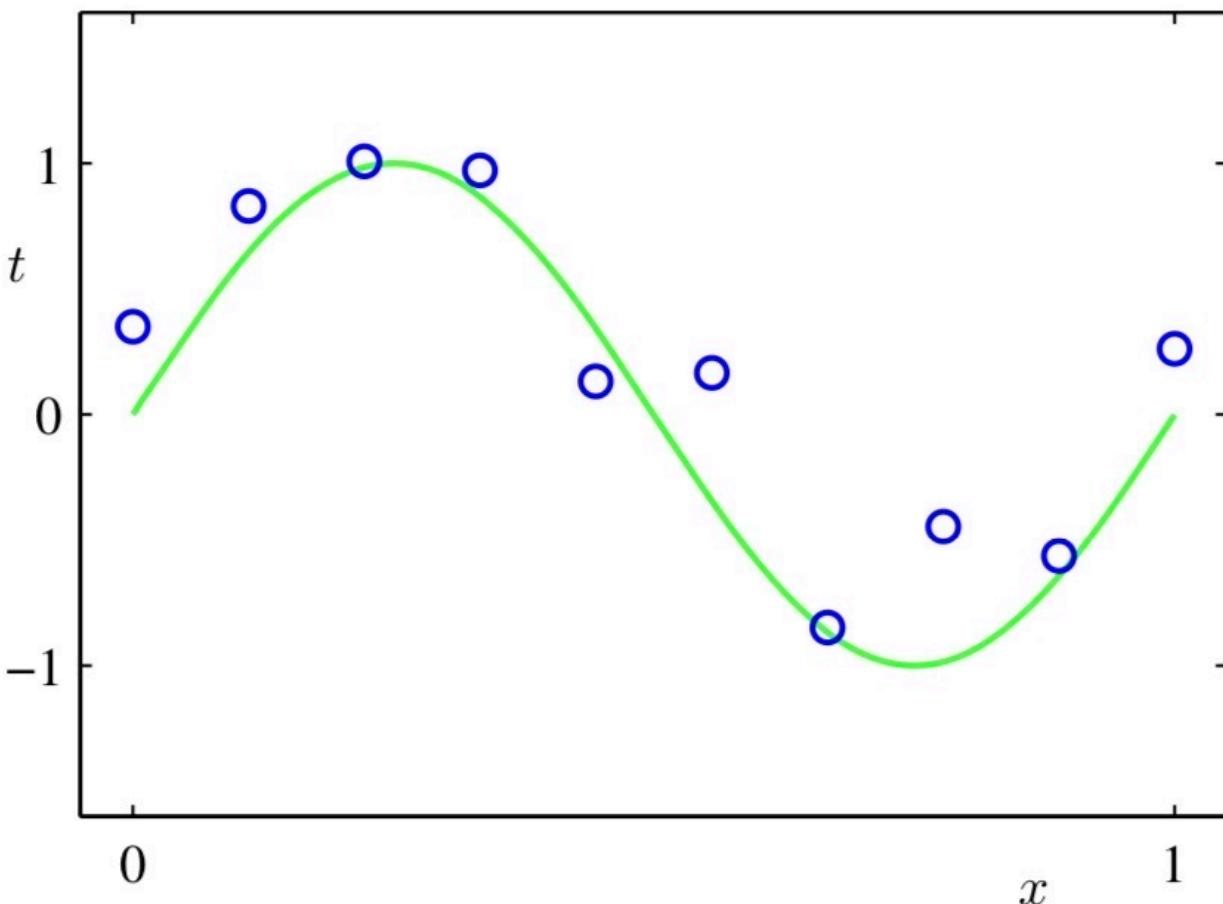


Figure 1.2

Plot of a training data set of $N = 10$ points, shown as blue circles, each comprising an observation of the input variable x along with the corresponding target variable t . The green curve shows the function $\sin(2\pi x)$ used to generate the data. Our goal is to predict the value of t for some new value of x , without knowledge of the green curve.



Goal: predict \hat{t} for a given \hat{x} (new)

$$y(x, w) := w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

where M is the order of polynomial and x^j denotes x raised to the power of j

($y(x, w)$ is linear w.r.t w
not " x

error function: measurement of the misfit between $y(x, w)$ and target value t

SSE, sum of square error (function of w)

$$E(w) := \frac{1}{2} \sum_{n=1}^N \{ y(x_n, w) - t_n \}^2 \geq 0$$

$$E(w) = 0 \quad \text{iff} \quad \begin{array}{c} \text{예측값} \\ y(x_n, w) \end{array} = \begin{array}{c} \text{실제값} \\ t_n \end{array} \quad \forall n$$

$E(w)$ is a quadratic function of w so minimization of $E(w)$ has a unique solution ($\nabla_w E(w)$ is linear)

Let w^* be the sol of minimization of $E(w)$

$$E(w^*) = \min_w E(w)$$

The resulting polynomial is given by the function

$$\underline{y(x, w^*)}$$

$$w_0^* + x w_1^* + x^2 w_2^* \dots$$

How to determine M ?

model comparison (model selection)

The number of parameters ($M+1$) determines model complexity

under-fitting vs over-fitting

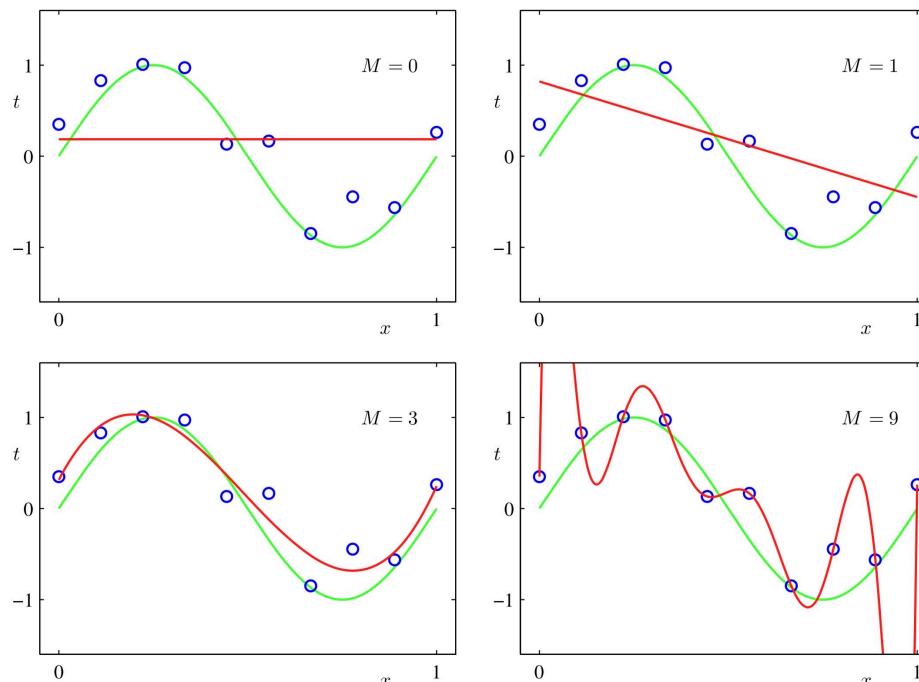


Figure 1.4 Plots of polynomials having various orders M , shown as red curves, fitted to the data set shown in Figure 1.2.

$N = 10$
9 차 다항식

RMSE (root mean square error)

$$E_{RMS} = \sqrt{\frac{2ECw^*)}{N}}$$

$M=9$

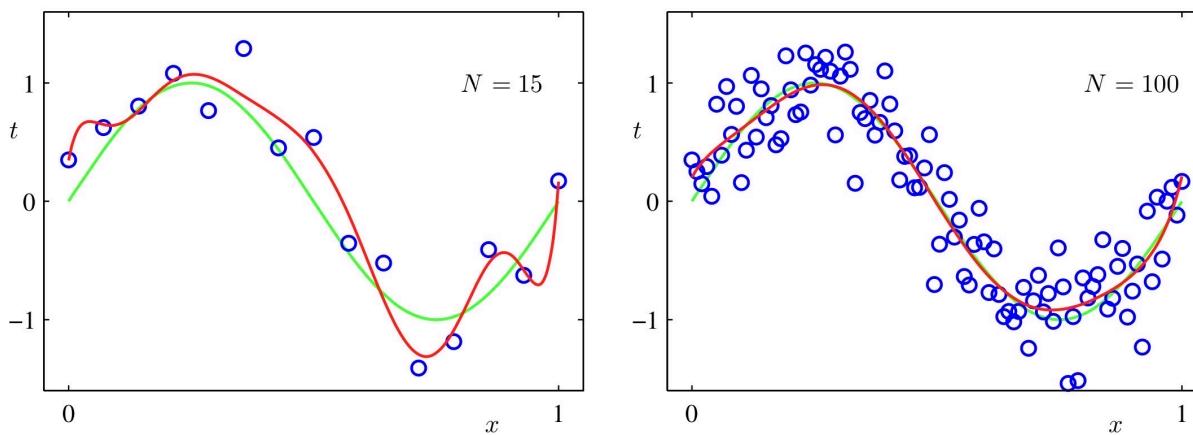


Figure 1.6 Plots of the solutions obtained by minimizing the sum-of-squares error function using the $M = 9$ polynomial for $N = 15$ data points (left plot) and $N = 100$ data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem.

Remark

$N=10$

- When $M = 9$, $E_{RMS} = 0$, 10 coefficients w_0, w_1, \dots, w_9 and 10 data points in training data.

Figure 1.5 Graphs of the root-mean-square error, defined by (1.3), evaluated on the training set and on an independent test set for various values of M .

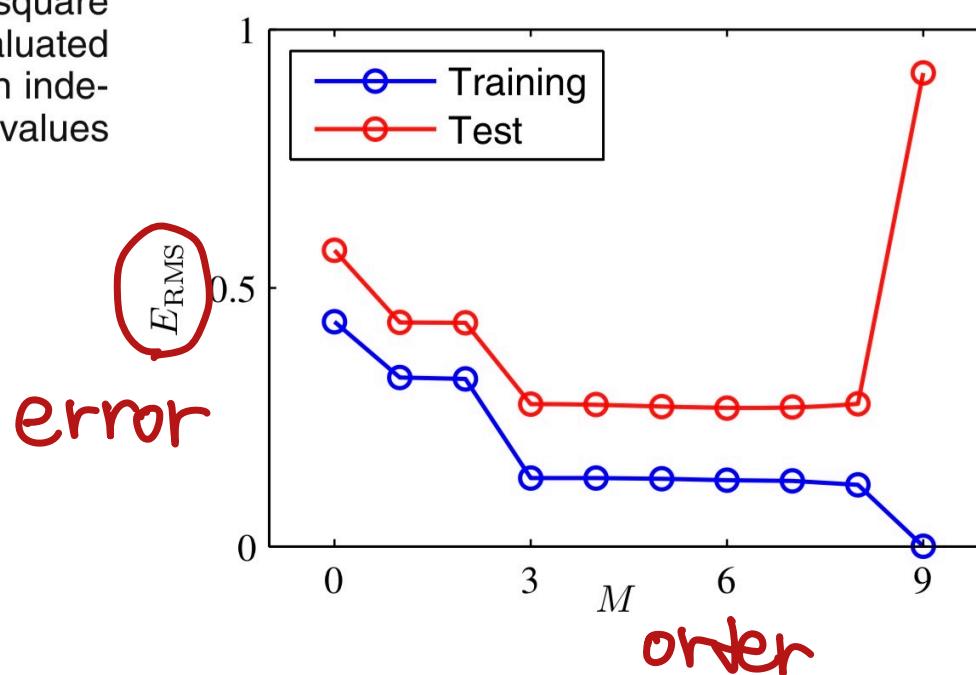


Table 1.1 Table of the coefficients w^* for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19		0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

Remark

solution

- When M increases, magnitude of w_i^* becomes larger
- The size of the training data set (fig. 1.6) over-fitting problem less serve as the size of the data set increases.

Regularization

- prevent over-fitting
- discourage the coefficients from reaching large values

SSE, RMSE, MSE

$$\tilde{E}(w) = E(w) + \frac{\lambda}{2} \|w\|^2$$

where $\|w\|^2 = w^T w = \sum_{i=0}^M w_i^2$

- Ridge regression, weight decay

$N = 10$
 $M = 9$

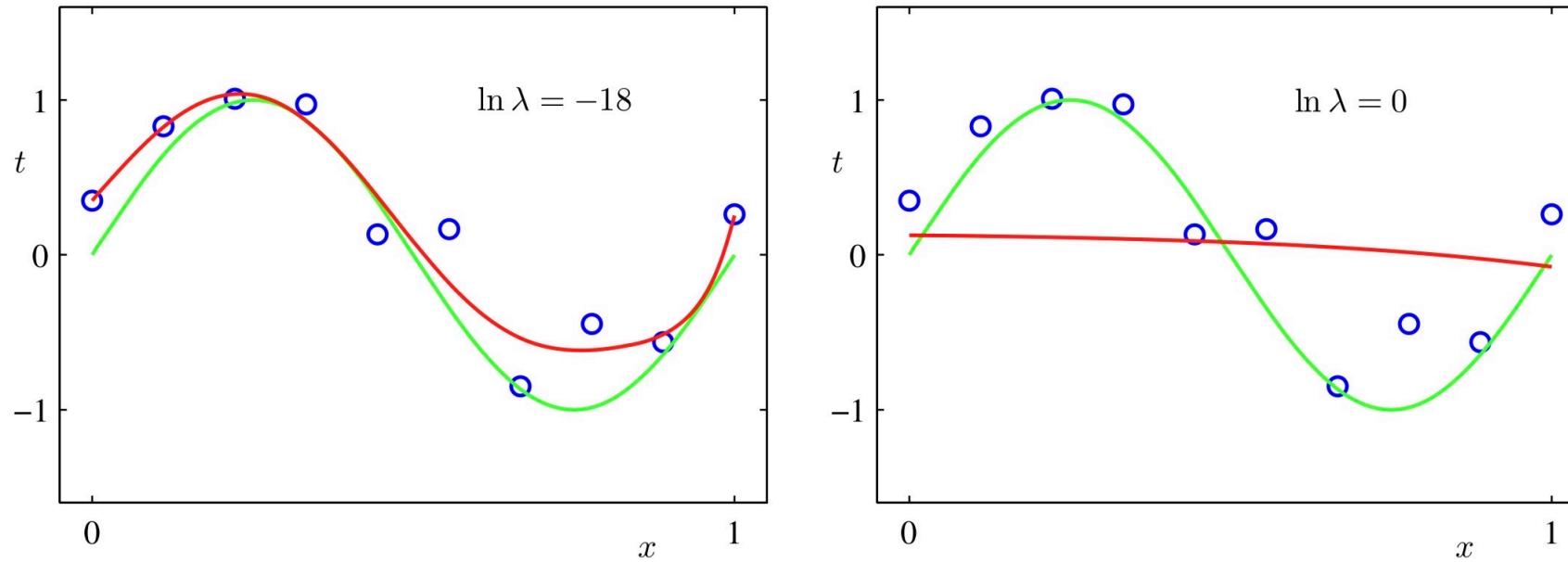


Figure 1.7 Plots of $M = 9$ polynomials fitted to the data set shown in Figure 1.2 using the regularized error function (1.4) for two values of the regularization parameter λ corresponding to $\ln \lambda = -18$ and $\ln \lambda = 0$. The case of no regularizer, i.e., $\lambda = 0$, corresponding to $\ln \lambda = -\infty$, is shown at the bottom right of Figure 1.4.

Validation set (hold - out set)

- prevent over-fitting



1.2 Probability Theory (elementary)

"uncertainty"

Sample space : set of all possible outcomes of experiment
 Ω

Random variable : mathematical formalization of quantity or
 X
object which depends on random events

$$X: \underline{\Omega} \rightarrow \underline{E}, \mathbb{R}^d \text{ measurable space } E$$

probability measure P on (Ω, \mathcal{F}) σ -field

$$(\Omega, \mathcal{F}, P)$$

$$\Omega = \{1, 2, 3, 4, 5, 6\} \quad P(\{\omega\}) = \frac{1}{6} \quad n=1, 2, \dots, 6$$

win + 1000 when n is even

lose - 1000 when n is odd

$$X(\omega) = \begin{cases} 1000 & \omega = 2, 4, 6 \\ -1000 & \omega = 1, 3, 5 \end{cases}$$

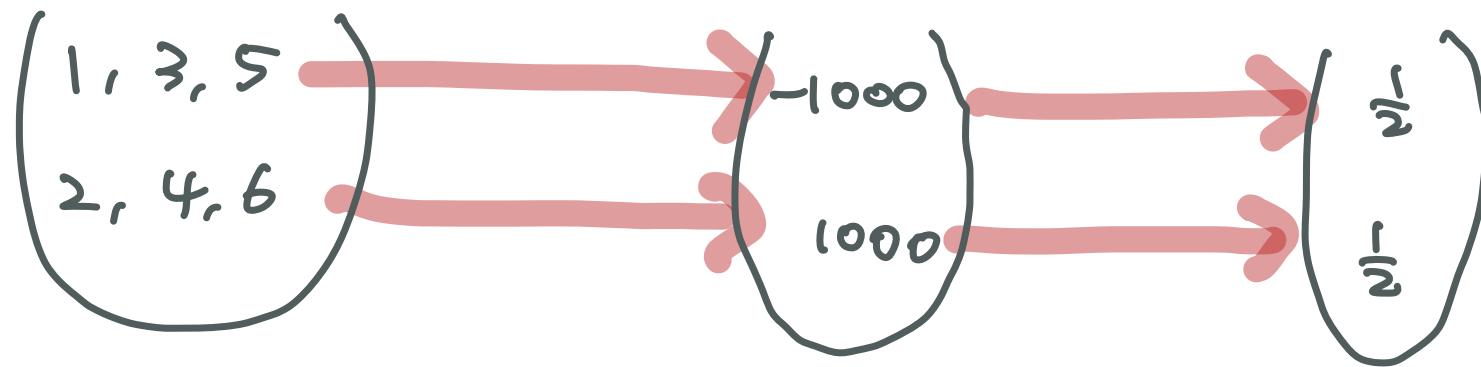
random variable

$$\begin{aligned} P(X = 1000) &= P(\{\omega \in \Omega \mid X(\omega) = 1000\}) \\ &= P(\{2, 4, 6\}) = \frac{1}{2} \end{aligned}$$

Sample space

r.v.

probability



Sum rule and product rule

two random variables X, Y

$$X = \{x_1, \dots, x_m\}$$

$$Y = \{y_1, \dots, y_L\}$$

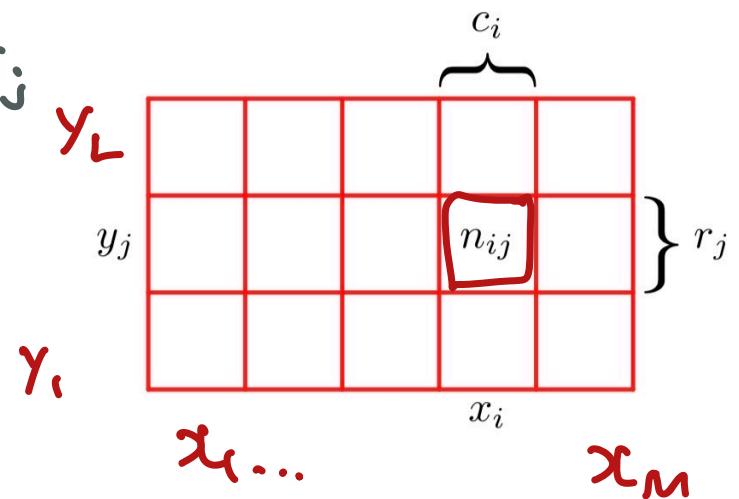
$X \quad Y$
1 (x_i, y_j)
2 (\dots)
3 \vdots
 $N(\quad)$

Consider N trials in which we sample both X, Y

Let $n_{ij} := *$ of trials in which $X = x_i$ and $Y = y_j$

Joint probability of $X = x_i$ and $Y = y_j$

$$P(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$



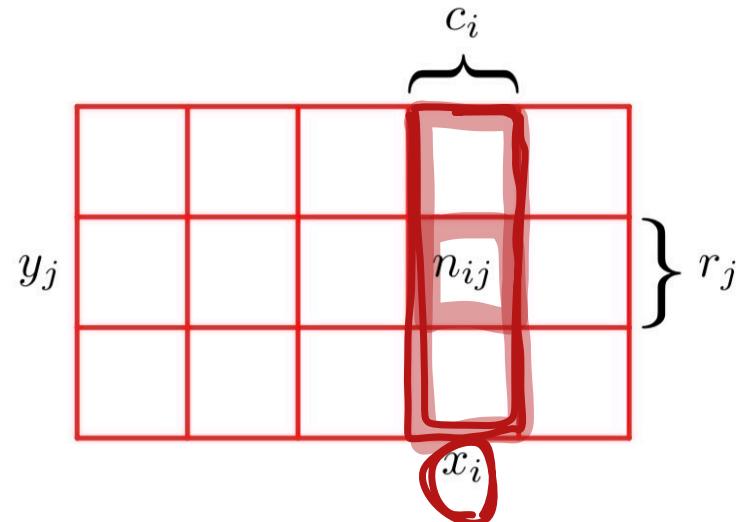
Let $c_i := \#$ of trials in which $X = x_i$ (regardless Y)

$$P(X = x_i) = \frac{c_i}{N}$$

$$P(x, Y) \rightarrow P(x) = \sum Y P(x, Y)$$

Sum rule

$$P(X = x_i) = \sum_{j=1}^r P(X = x_i, Y = y_j) \quad \text{(marginal probability)}$$



Conditional probability

$$P(Y = y_j | X = x_i) = P(Y = y_j, X = x_i)$$

$$P(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i} = P(X = x_i) > 0 \quad \text{인 경우}$$

product rule

$$\begin{aligned} P(X=x_i, Y=y_j) &= \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} \\ &= P(Y=y_j | X=x_i) P(X=x_i) \end{aligned}$$

The rules of probability

X, Y random variables

sum rule $P(X) = \sum_Y P(X, Y)$

product rule $P(X, Y) = P(Y|X) P(X)$

$$= P(X|Y) P(Y)$$

$P(X, Y)$: the probability of X and Y
 $P(Y|X)$: the probability of Y given X

By symmetry property $P(X, Y) = P(Y, X)$,

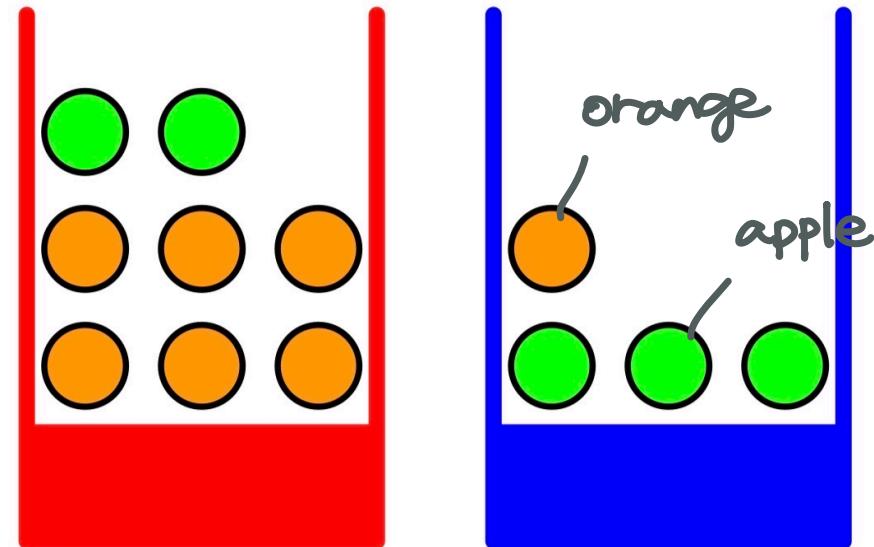
$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \quad \left(= \frac{\underline{P(X, Y)}}{P(X)} \right)$$

Bayes's theorem

$$P(X) = \sum_Y P(X|Y) P(Y)$$

Example

Figure 1.9 We use a simple example of two coloured boxes each containing fruit (apples shown in green and oranges shown in orange) to introduce the basic ideas of probability.



pick red box 40% of the time

" blue " 60%

$$P(B = \text{red}) = 4/10$$

$$P(B = \text{blue}) = 6/10$$

Suppose pick a box randomly and it turns out to be the blue box. Probability of ($F = \text{apple}$) given ($B = \text{blue}$)

$$P(F = \text{apple} \mid B = \text{blue}) = \frac{3}{4}$$

$$P(F = a \mid B = r) = \frac{1}{4}$$

$$P(F = o \mid B = r) = \frac{3}{4}$$

$$P(F = a \mid B = b) = \frac{3}{4}$$

$$P(F = o \mid B = b) = \frac{1}{4}$$

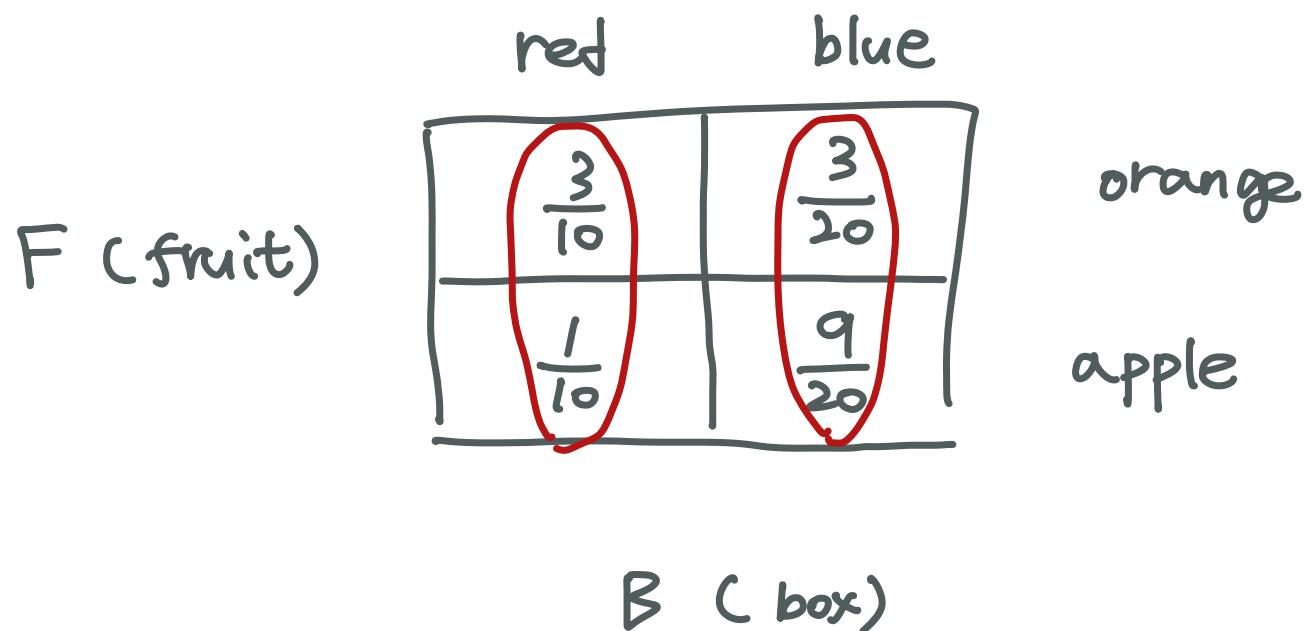
$$P(F=a) = P(F=a | B=r) P(B=r) + P(F=a | B=b) P(B=b)$$

//

$$= \frac{1}{4} \times \frac{4}{10} + \frac{3}{4} \times \frac{6}{10} = \frac{11}{20}$$

sum
rule

$$P(F=a, B=r) + P(F=a, B=b)$$



$$P(F=0) = \frac{q}{2^6}$$

$$P(B=r | F=0) = \frac{P(B=r, F=0)}{P(F=0)} = \frac{P(F=0 | B=r) P(B=r)}{P(F=0)}$$
$$= \frac{2}{3}$$

Interpretation of Bayes's theorem

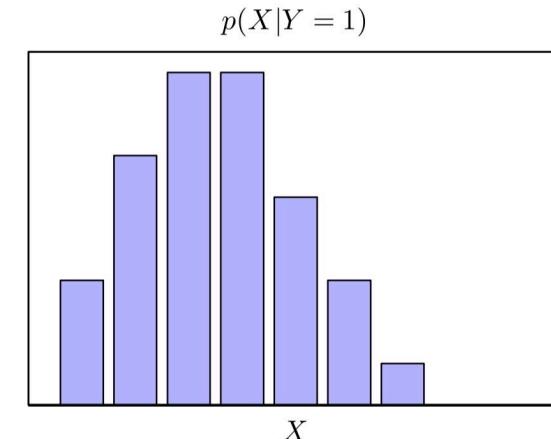
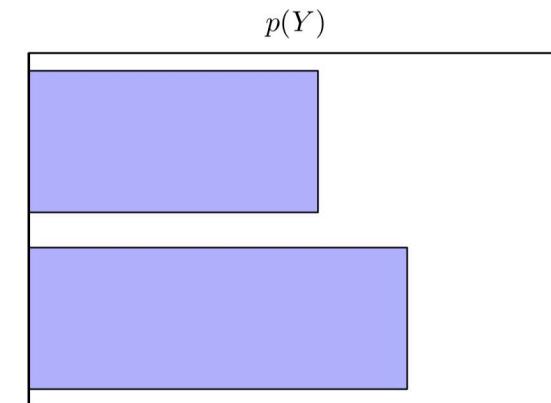
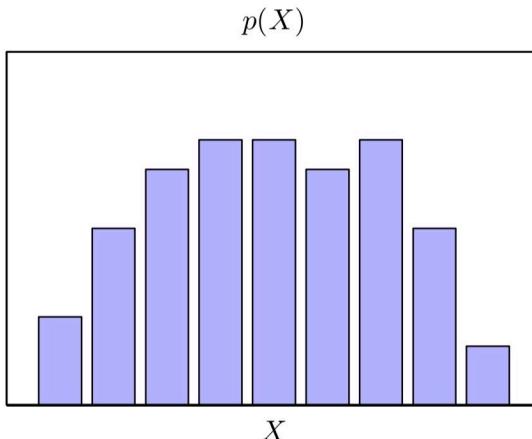
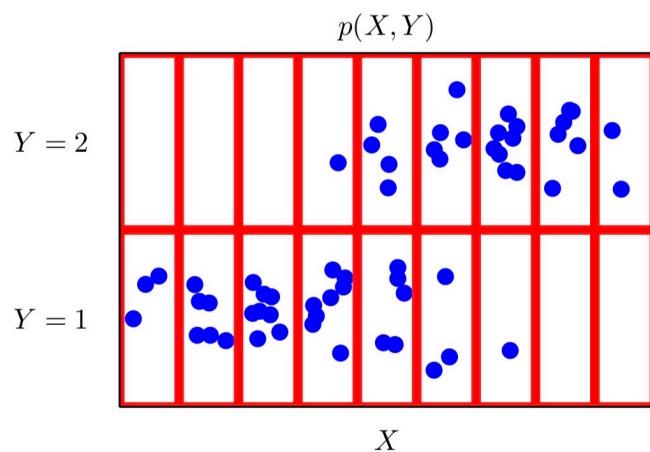
$P(B)$ prior probability (이전 box 2를 골랐는지)
(probability available before we observe the identity of the fruit)

$P(B|F)$ posterior probability
(과일 확인 후 이전 box)

Two random variables X and Y are independent
 P q

$$\underline{P(X, Y) = P(X) P(Y)}$$

$$\Rightarrow P(Y|X) = P(Y) \quad \text{and} \quad P(X|Y) = P(X)$$



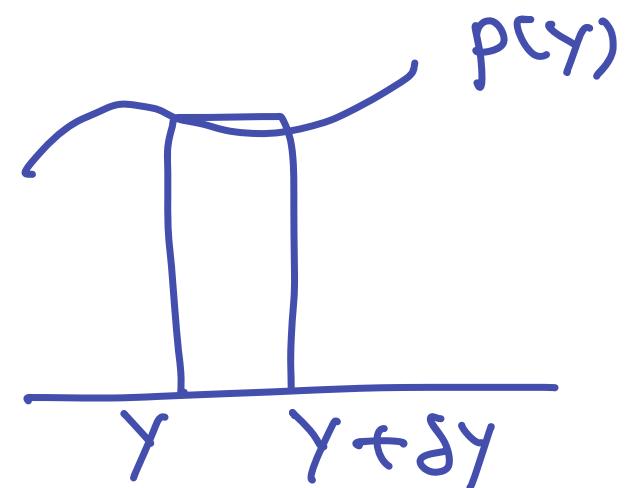
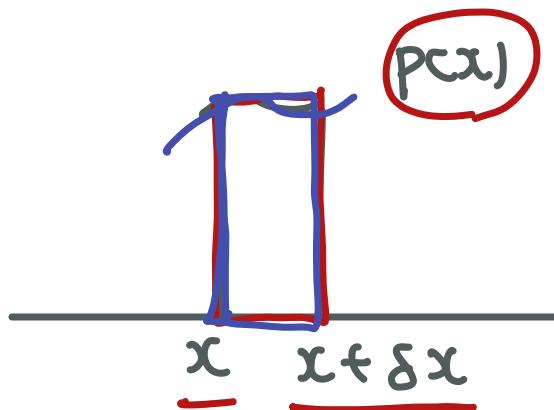
1. 2. 1 Probability density

continuous variable

open interval

$$P(\text{real valued variable falling in } (x, x + \delta x)) \approx \frac{p(x) \delta x}{p(x) \cdot \Delta x}$$

when δx is sufficiently small



$$\underline{x} + \Delta x$$

$$p(y, y + \delta y) \approx \underline{p(y)} \delta y$$

$p(x)$ is called the probability density over x

$$P(x \in (a,b)) = \int_a^b p(x) dx$$

$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

$$P_y(Y) dy = P_x(x) dx$$

$$\frac{dx}{dy}$$

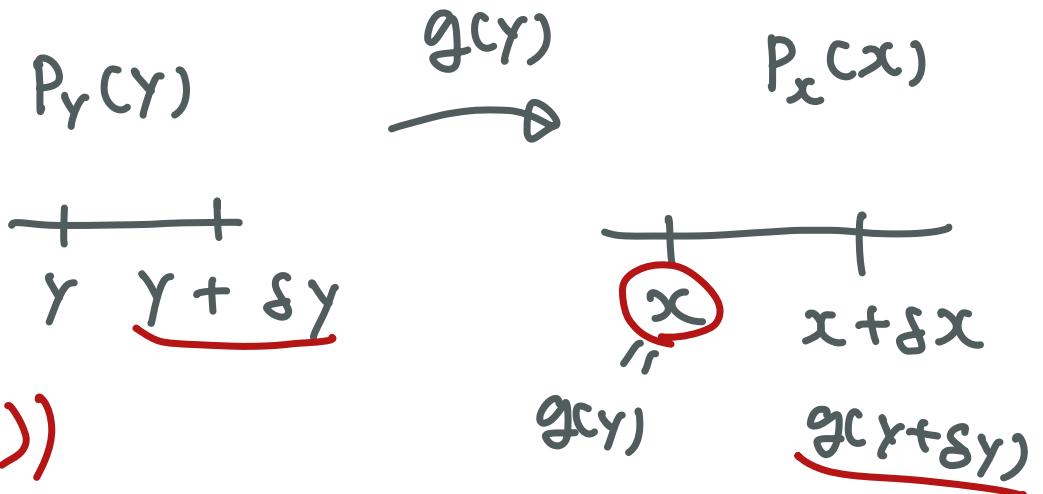
Consider a change of variables.

Bijective function

$$x = g(y)$$

$$P_y(y, y + \delta y) = P_x(g(y), g(y + \delta y))$$

$$P_y(y) dy$$



$$\approx P_x(x, x + \delta x) \text{ Chapter 1 Introduction}$$

$$P_x(x) dx$$

$$\tilde{f}(y) = f(g(y))$$

Find the density $P_Y(y)$

$(x + \delta x)$ is transformed into the range $(y, y + \delta y)$

$$P_X(x) \delta x \approx P_Y(y) \delta y$$

$$\Rightarrow P_Y(y) = P_X(x) \left| \frac{dx}{dy} \right|$$

$$= P_X(g(y)) |g'(y)|$$

Cumulative

large small

$$P(z) = \int_{-\infty}^z \textcircled{P(x)} dx$$

density

Multivariate continuous variables x_1, \dots, x_D , denoted by \mathbf{x} .

joint probability density $p(\mathbf{x}) = p(x_1, x_2, \dots, x_D)$ s.t

p (multi. variable falling in an infinitesimal

volume $\delta \mathbf{x}$ containing \mathbf{x})



$$\approx p(\mathbf{x}) \delta \mathbf{x}$$

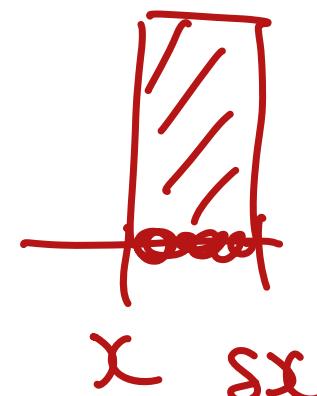
Multivariate probability density

measurable set $E \subset \mathbb{R}^D$

$$p(\mathbf{x} \in E) = \int_E p(\mathbf{x}) d\mathbf{x}$$

$$p(\mathbf{x}) \geq 0$$

$$\int_{\mathbb{R}^D} p(\mathbf{x}) d\mathbf{x} = 1$$



We can also consider joint probability distributions over a combination of discrete and continuous

In the discrete case

$p(x)$ is called probability mass function.

1.2.2 Expectations and covariances

Variable x under a probability distribution $p(x)$
(mass density)

Expectation of $f(x)$

In discrete case,

$$\mathbb{E}[f] := \sum_x p(x) f(x) \quad (\text{weighted mean})$$

In continuous case

$$\mathbb{E}[f] := \int p(x) f(x) dx$$

Random variable x .

Expectation of x is denoted by

discrete

$$\mathbb{E}[x] := \sum_x p(x) x \quad (f(x) = x)$$

continuous

$$= \int p(x) x dx$$

$E[f]$ can be approximated as a finite sum (N samples)

$$E[f] \approx \frac{1}{N} \sum_{n=1}^N f(x_n) \quad (\text{Monte Carlo Integration})$$

Average of the function $f(x,y)$ w.r.t distribution of x

$$E_x [f(x,y)] = \int f(x,y) p(x) dx \quad (\text{function of } y)$$

(or $\sum_x f(x,y) p(x)$)

Conditional expectation

$$E_x [f | y] = \int f(x) p(x|y) dx = \sum_x p(x|y) f(x)$$

(function of y)

Variance of $f(x)$

exp. of f

$$\text{var}[f] := \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

$$\Rightarrow = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

As $\mathbb{E}[x]$, $\text{Var}[x] := \mathbb{E}[(x - \mathbb{E}[x])^2]$ ($f(x) = x$)

Two variables x and y . covariance is defined by

$$\text{cov}[x, y] = \mathbb{E}_{x,y}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])]$$

x 와 y 가

$$p(x,y) = \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

같아나 같아
변동하는 x

If x and y are independent, $\text{cov}[x, y] = 0$

$$\iint x \cdot y p(x,y) dx dy = \iint xy p(x) p(y) dx dy$$

Two vectors of random variables \mathbf{x} and $\mathbf{y} \in \mathbb{R}^D$.

Let $f(\mathbf{x})$ be a multivariable function (can be vector function)

$$\mathbb{E}[f] := \sum p(x) f(x) = \int f(x) p(x) dx$$

$$\underline{f(x) = x}$$

If $f(x)$ is a vector function, $\mathbb{E}[f]$ is a vector

Covariance matrix is defined by

$1 \times D$

$$\begin{aligned} \text{cov}[\mathbf{x}, \mathbf{y}] &:= \mathbb{E}_{\substack{\mathbf{x}, \mathbf{y} \\ D \times D}} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}]) (\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]) \right] \\ &= \mathbb{E}_{\substack{\mathbf{x}, \mathbf{y}}} [\mathbf{x} \mathbf{y}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{y}^T] \end{aligned}$$

$p(x, y)$

$$\text{Var}[\mathbf{x}] = \text{cov}[\mathbf{x}, \mathbf{x}]$$

i,j component

$= x_i, y_j$ 의 variance

1. 2. 3 Bayesian probability

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} = P(X, Y)$$

비례

posterior

사후

$$P(Y|X) \propto P(X|Y) P(Y)$$

likelihood

가능도 함수

Prior

사전분포

Classical frequentist

vs

Bayesian

- frequencies of random
(repeatable events)

데이터 수가 적을 때
over-fitting

quantification of uncertainty
정량

데이터 수가 적은 경우 우리

직관적인 학습은 일반적으로 classic

빨간 상자를 선택했을 학습 vs 사람 가 나왔을 때 빨간 상자를
선택했을 학습

우리는 정보 (data) 를 활용함. 정보에 의해서 예측값 변화를 반영

Polynomial fitting curve

Aim to predict $w = (w_0, w_1, \dots, w_M)$

$$y(x, w) := w_0 + w_1 x + \dots + w_M x^M \quad (\text{M-th order})$$

Before observing data, assume w_j is in the form

(PC w/)

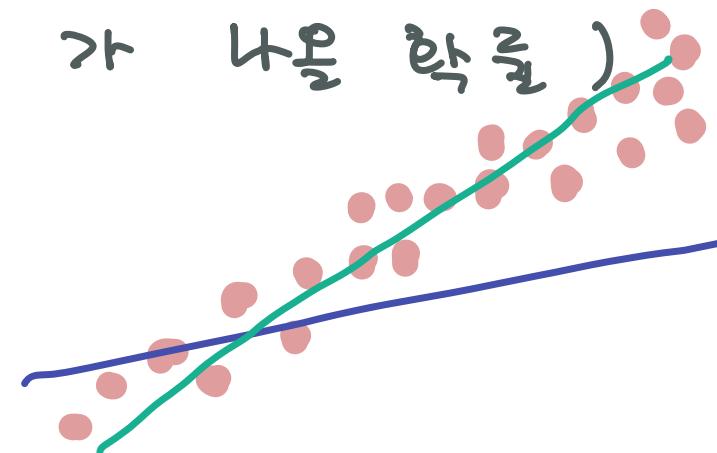
The effect of $D = \{t_1, \dots, t_N\}$ is expressed through

Prior

$PCD(w)$

C W 3 해당 데이터

$$t \sim N(x_1 | \gamma(x, w), \sigma^2)$$



Uncertainty in w after we have observed D

posterior likelihood prior

$$\underline{P(w|D)} = \frac{P(D|w) P(w)}{\underline{P(D)}}$$

Since data set D is given, $P(D|w)$ is a function
of w , called $P(D|w)$ likelihood function (가능도 함수)

Remark

- $P(D|w)$ is not a prob. distribution over w
- $\int P(D|w) dw \neq 1$

$$\underline{\int p(x|Y) dx = 1}$$

$$\underline{\int p(x|Y) dy \neq 1}$$

- $P(D)$ can be seen as a normalization constant
- Likelihood $P(D|w)$ has an important role both approaches

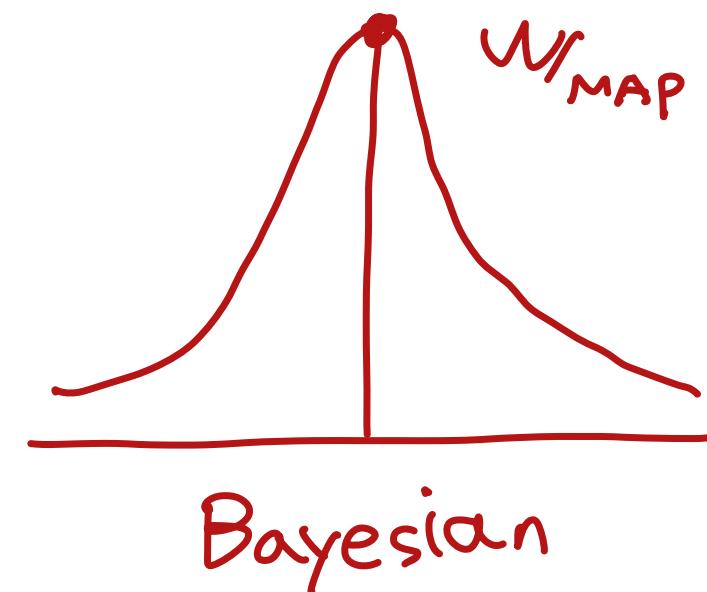
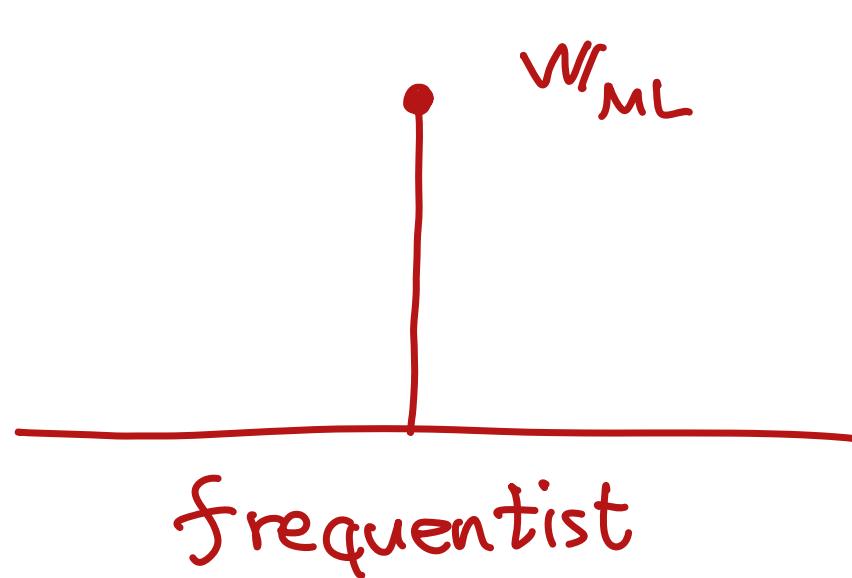
$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

Frequentist (point estimator)

- w is considered to be a fixed parameter whose value is by some form of estimator

Bayesian (distribution estimator)

- There is only a single data set D and uncertainty in the parameter is expressed through a probability distribution over w



MLE

MAP

빈도론에서 주로 사용하는 방법으로 maximum likelihood

error function = $-\text{likelihood}$ or $-\log(\text{likelihood})$

Remark (베이지안 특징)

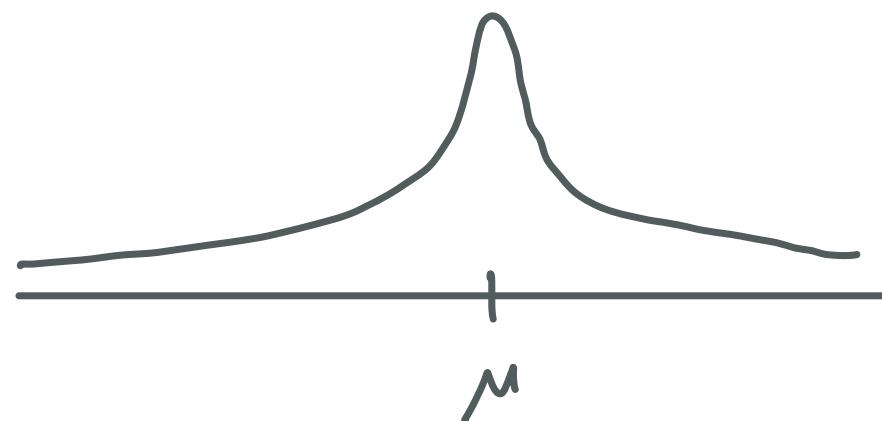
- 사전지식 (사전 분포)을 추론 과정에 반영할 수 있음
- But 일반적으로 사전 분포를 실제 사정의 믿음 보다는 수학적 편의성을 고려해서 선택 (conjugate)
- 주관적인 의견이 추론 과정에 반영
- Marginalization (주변화) over w 필요.

03-17

1.2.4 Gaussian distribution

Real-valued variable x , Gaussian distribution is defined by

$$N(x|\mu, \sigma^2) := \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} (x-\mu)^2 \right\} \quad (\text{pdf})$$



unimodal
continuous

μ : mean, σ^2 : variance

(2 parameters)

σ : standard deviation, $\beta = \frac{1}{\sigma^2}$: precision

Observe

$$N(x|\mu, \sigma^2) > 0$$

$$\int_{-\infty}^{\infty} N(x|\mu, \sigma^2) dx = 1 \quad (*)$$

$$E[x] = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x dx = \mu$$

second moment

$$E[x^2] = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2$$

$$\Rightarrow \text{Var}[x] = E[x^2] - E[x]^2 = \sigma^2$$

Let $I := \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2} x^2\right) dx$

$$I^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2} x^2\right) \exp\left(-\frac{1}{2\sigma^2} y^2\right) dx dy$$

$$= \int_0^{2\pi} \int_0^{\infty} \exp\left(-\frac{1}{2\sigma^2} r^2\right) r dr d\theta \quad \begin{aligned} x &= r \cos \theta \\ y &= r \sin \theta \end{aligned}$$

$$= 2\pi \int_0^{\infty} \exp\left(-\frac{1}{2\sigma^2} r^2\right) r dr$$

$$= 2\pi \sigma^2$$

$$\Rightarrow I = \sqrt{2\pi} \sigma$$

D-dimensional multivariate Gaussian distribution

$$N(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

determinant

$\boldsymbol{\mu}$: mean vector, Σ : covariance matrix ($D \times D$)

(see later!)

$\mathbf{x} \in \mathbb{R}^D$ random vector

N scalar variables x_1, \dots, x_N . Let $\mathbf{x} = (x_1, \dots, x_N)^T$

Assume x_i are drawn from a normal dist.

How to find the parameters μ and σ^2

(parametric inference)

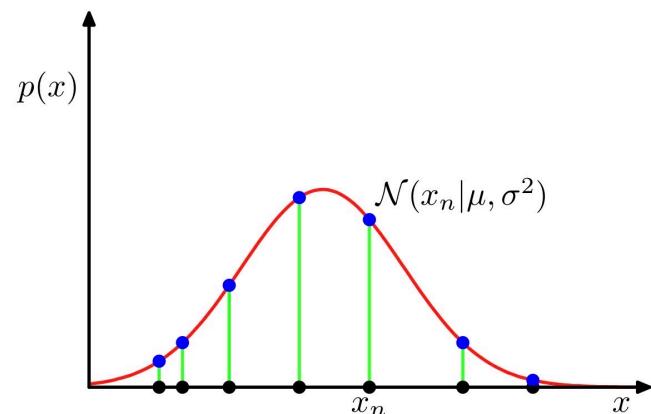
x_i are samples the same distribution independently (iid)

$$p(\mathbf{x} | \mu, \sigma^2) = \prod_{n=1}^N N(x_n | \mu, \sigma^2)$$

likelihood

Maximize likelihood

Figure 1.14 Illustration of the likelihood function for a Gaussian distribution, shown by the red curve. Here the black points denote a data set of values $\{x_n\}$, and the likelihood function given by (1.53) corresponds to the product of the blue values. Maximizing the likelihood involves adjusting the mean and variance of the Gaussian so as to maximize this product.



maximize likelihood \Leftrightarrow maximize $\log(\text{likelihood})$

$$\ln p(X|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

constant

Maximizing wrt μ ,

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \quad (\text{sample mean})$$

Maximizing wrt σ^2

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2 \quad (\text{sample variance})$$

$$N(x|\mu, \sigma^2)$$

Fix N and some dist. μ_{ML} and σ_{ML}^2 are functions of the choice x_1, \dots, x_N

$$\mathbb{E}[\mu_{ML}] = \mu \quad \text{unbiased}$$

$$\mathbb{E}[\sigma_{ML}^2] = \left(\frac{N-1}{N}\right) \sigma^2$$

$$\mathbb{E}[\mu_{ML}] = \mathbb{E}\left[\frac{1}{N} \sum x_n\right] = \frac{1}{N} \sum \mathbb{E}[x_n] = \mu$$

1.2.5 Curve fitting (re-visited)

Probabilistic perspective

input value x , target value t

N input values $\mathbf{x} = (x_1, \dots, x_N)^T$

"target" " $\mathbf{t} = (t_1, \dots, t_N)^T$

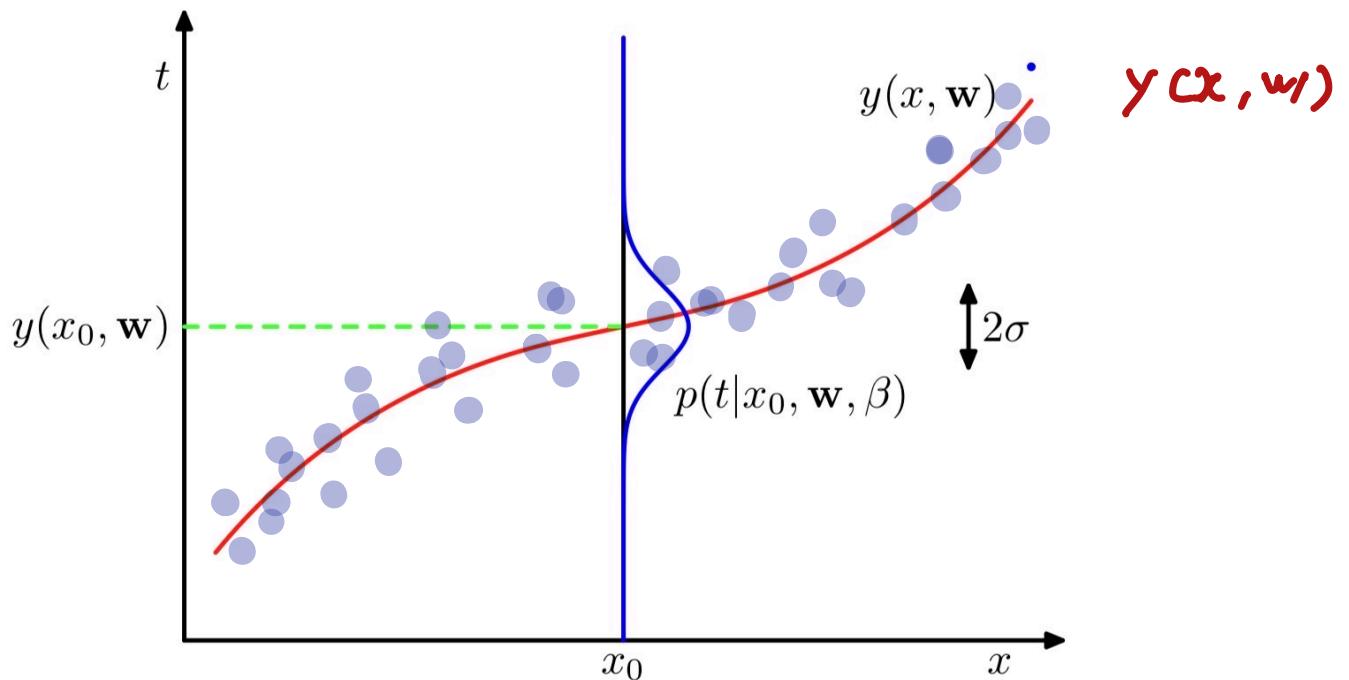
Predict target value t for some new input x

Express uncertainty over target t using prob. dist

Assume

$$p(t|x, w, \beta) = N(t | y(x, w), \frac{1}{\beta}) \quad (1.60)$$

Figure 1.16 Schematic illustration of a Gaussian conditional distribution for t given x given by (1.60), in which the mean is given by the polynomial function $y(x, w)$, and the precision is given by the parameter β , which is related to the variance by $\beta^{-1} = \sigma^2$.



(t : 평균이 $y(x, w)$ 이고 분산이 $\frac{1}{\beta}$ 인 Gaussian)

Determine w and β by MLE.

Assume data x and t are drawn independently from (1.60)

$$p(t|x, w, \beta) = \prod_{n=1}^N N(t_n | y(x_n, w), \frac{1}{\beta})$$

by def

$$\text{assumed to be fixed} = \prod \left(\frac{\beta}{2\pi} \right)^{1/2} \exp \left\{ -\frac{\beta}{2} (y(x_n, w) - t_n)^2 \right\}$$

Maximize log likelihood wrt w and β

$$\ln p(t|x, w, \beta) = -\frac{\beta}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln (2\pi)$$

solutions for w_{ML} , β_{ML} !

\downarrow
 $y(x, w_{ML})$, \hat{x}

Point prediction (not use β_{ML})

New input x , $t := \gamma(x, w_{ML})$

Distribution prediction (not Bayesian)

$$p(t|x, w_{ML}, \beta_{ML}) = N(t | \gamma(x, w_{ML}), \beta_{ML}^{-1})$$

Take a step towards a more Bayesian approach.

Prior distribution over w as the form

$$\begin{pmatrix} \alpha & & \\ & \ddots & \\ & & \alpha \end{pmatrix}$$

$$p(w|\alpha) = N(w|0, \alpha^{-1}I) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M+1}{2}} \exp\left\{-\frac{\alpha}{2} w^T w\right\}$$

where α is a precision and $M+1$ is # of coeff. w

($y|x, w$) m th order polynomial, w_0, w_1, \dots, w_M)

α is called hyperparameter (parameter of parameter)

$$p(w|\alpha, \beta, \alpha, \beta) \propto p(t|x, w, \beta) p(w|\alpha)$$

likelihood

데이터 $x, \#$ 에 대해 가능성이 높은 w 를 찾기.

Maximize posterior distribution (MAP, maximum posterior)

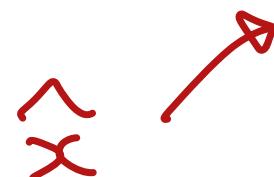
$$\Leftrightarrow \text{Minimize} \quad \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{\alpha}{2} w^T w$$

see later!

w_{MC}

w_{MAP}

$y(x, w_{MAP})$



$\hat{t}_{\phi(\Sigma)}$

1. 2. 6 Bayesian curve fitting

We still make a point estimate of w .

In a fully Bayesian approach, we should apply sum and product rules of probability. So we need to integrate over all values of w . (Marginalization)

Compared with MAP (maximize posterior), we will use

posterior $p(w | \mathcal{X}, \mathcal{T})$ (\mathcal{X}, \mathcal{T} training data set)

all information

Predict the value of t for a new test point x

\Rightarrow Evaluate $p(t|x, \alpha, \beta)$

(새로운 x 값에 대해서 예측값 t 의 분포, 기존 데이터 정보 활용)

Now α and β are assumed to be given and fixed.

$$p(t|x, \alpha, \beta) = \int \underbrace{p(t|w)}_{\text{posterior}} p(w|\alpha, \beta) dw$$

(1.60)

$$\sim N(t|y(x, w), \beta^{-1})$$

assumption

We will see and perform analytically later. The result is given by a Gaussian

$$p(t|x, \Phi, \#) = N(t | m(x), s^2(x))$$

where $m(x) = \beta \Phi(x)^T S \sum_{n=1}^N \Phi(x_n) t_n$

$$s^2(x) = \beta^{-1} + \Phi(x)^T S \Phi(x)$$

Here matrix S is given by

$$S^{-1} = \alpha I + \beta \sum_{n=1}^N \Phi(x_n) \Phi(x_n)^T$$

I : unit matrix ($M+1$ dim), $\Phi(x) = (\phi_0(x), \dots \phi_M(x))^T$, $\phi_i(x) = x^i$

we will generalize later

Recall assumptions

- Prior distribution over w : $p(w|\alpha) = N(w|0, \alpha^{-1}I)$
- likelihood : $p(t|x, w, \beta) = \prod_{n=1}^N N(t_n | y(x_n, w), \beta^{-1})$
- posterior $p(w|x, t, \alpha, \beta) \propto p(t|x, w, \beta) p(w|\alpha)$

$$p(t|x, \alpha, \beta) = \int p(t|x, w) p(w|\alpha) dw$$

/

likelihood of single point

dependent of x

Curve fitting

$$y(x, w) := \sum_{j=1}^M x^j w_j \quad (M \text{ fixed})$$

model assumption

Set any error function

$$\text{Assume } t \sim N(t | y(x, w), \frac{1}{\beta})$$

Find w^* s.t

$$\text{likelihood } \prod_{n=1}^N N(t_n | y(x_n, w), \frac{1}{\beta})$$

minimize error function

$$\text{prior } w \sim N(w_0 | 0, \alpha^{-1} I)$$

$$y(x, w^*)$$

$$\Rightarrow \text{posterior } p(w | x, t, \alpha, \beta)$$

$$\propto p(t | x, w, \beta) p(w | \alpha)$$

$$\Rightarrow p(t | x, \hat{x}, \hat{t})$$

$$= \int p(t | x, w) p(w | \alpha) d w$$

1.3 Model selection

of parameter (coefficients) \approx model complexity

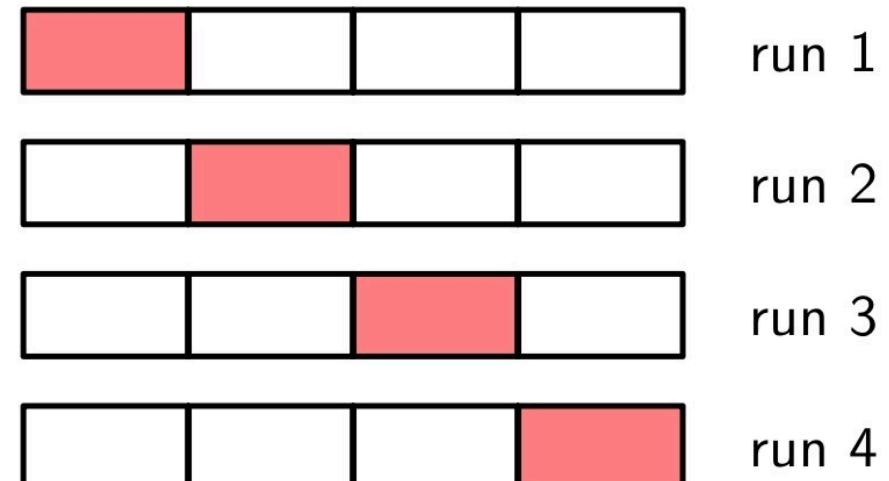
The performance on the training set is not a good indicator of predictive performance (because of over-fit)

S - fold cross-validation



Figure 1.18

The technique of S -fold cross-validation, illustrated here for the case of $S = 4$, involves taking the available data and partitioning it into S groups (in the simplest case these are of equal size). Then $S - 1$ of the groups are used to train a set of models that are then evaluated on the remaining group. This procedure is then repeated for all S possible choices for the held-out group, indicated here by the red blocks, and the performance scores from the S runs are then averaged.



when S increases, # of training runs is increases

Another method is to add penalty term

E.g. Akaike information criterion

$$\ln P(D|W_{ML}) - M$$

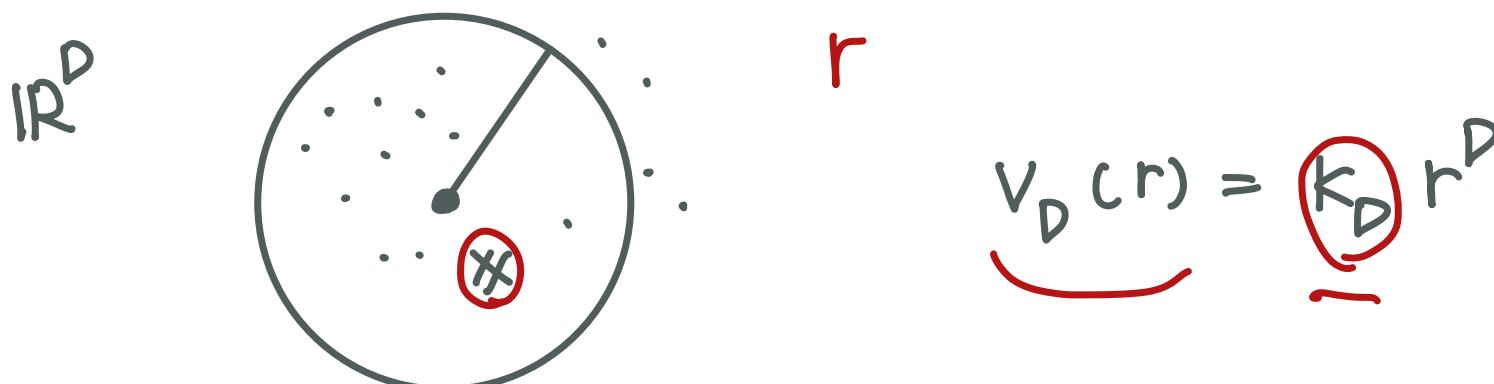
where M is # of parameters

1.4 The curse of Dimensionality

The dimension of input vector \mathbf{x}

결론 : 차원이 높다고 (feature의 수가 많다고) 예측은 좋은건 X

- 높은 차원 만큼 많은 데이터가 필요함
- 차원수 (또는 D^M)에 비례한 coefficients를 찾아야 함
- 고차원에서의 직관은 저차원의 직관과 매우 다름



$$\frac{V_D(1) - V_D(1-\epsilon)}{V_D(1)} = 1 - (1-\epsilon)^D$$

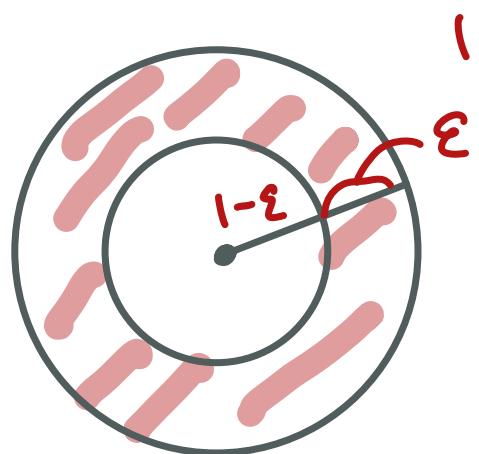
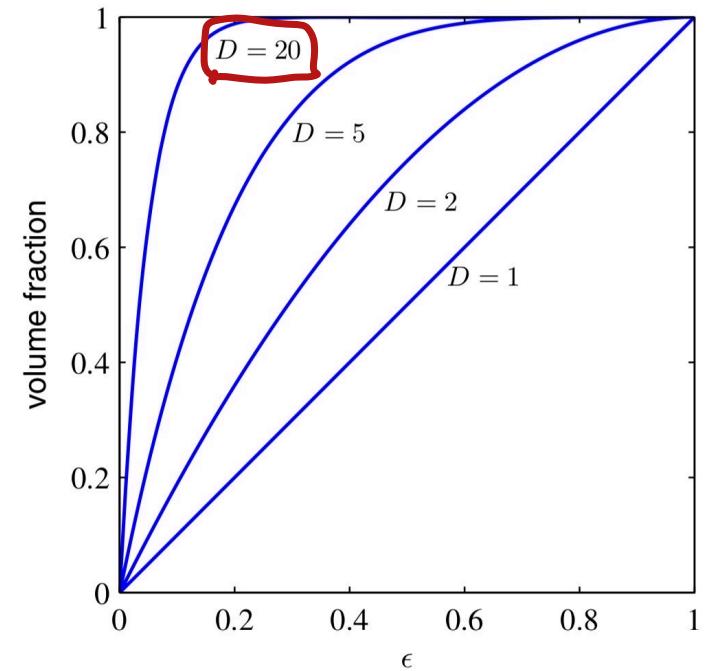


Figure 1.22 Plot of the fraction of the volume of a sphere lying in the range $r = 1 - \epsilon$ to $r = 1$ for various values of the dimensionality D .



Generally, real data is confined to a region of the space having lower dimensionality (manifold)

1.5 Decision Theory

Input vector \mathbf{x} , target vector \mathbf{t}

$$\begin{matrix} p(\mathbf{x}, \mathbf{t}) \\ p(\mathbf{t}|\mathbf{x}) \end{matrix} ?$$

For regression, \mathbf{t} comprise continuous value 2131

For classification, \mathbf{t} represent class label

Joint prob. $p(\mathbf{x}, \mathbf{t})$ provides a 'complete' summary of the uncertainty associated with \mathbf{x}, \mathbf{t} .

Determination of $p(\mathbf{x}, \mathbf{t})$ is an typical example. (difficult)

Input image \times (x-ray) binary classification

Class label C_1, C_2

Inference: modeling $P(C|x, C_k)$

Decision: for given x , decide which of two classes x is.

We are interested in $P(C_k | x)$ ✓

By Bayes's theorem,

$$P(C_k | \mathbf{x}) = \frac{P(x|C_k) P(C_k)}{P(x)}$$

class density

prior

posterior

P(x|C_k) is circled in red.

The entire fraction $P(x|C_k) P(C_k) / P(x)$ is underlined in red.

Any quantities appearing in theorem can be obtained from $p(x, C_k)$ by either marginalizing or conditioning

Minimize the chance to assigning \mathbf{x} to wrong class

\Rightarrow Choose the class having the higher posterior prob.

1.5.1 Minimizing the misclassification rate

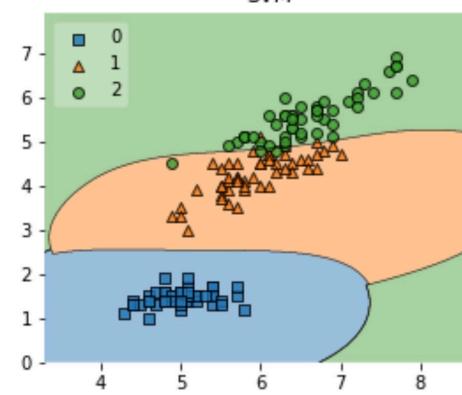
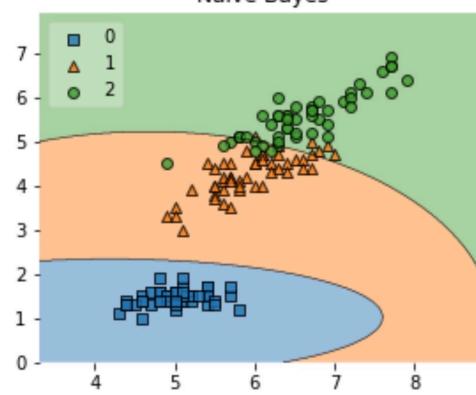
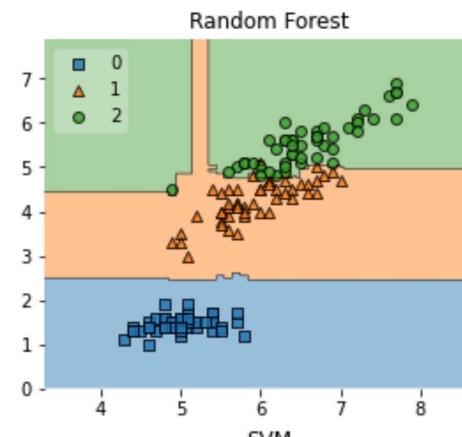
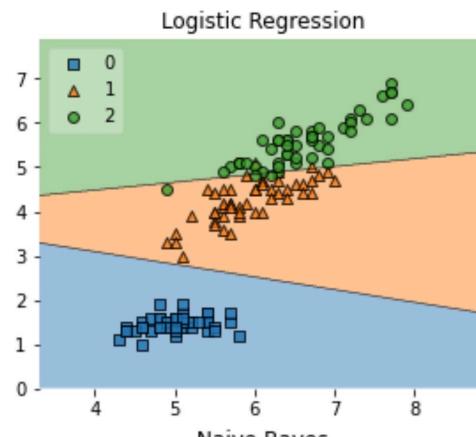
Decision regions R_k

if $x \in R_i$, then x

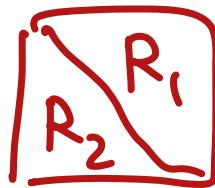
\times

- Regions in domain space assigning points to class C_k

Decision boundary : boundaries between R_k



Suppose our goal is to make as few misclassification as possible.



$$P(\text{mistake}) = P(x \in R_1, C_2) + P(x \in R_2, C_1)$$

error

$$= \int_{R_1} P(x, C_2) dx + \int_{R_2} P(x, C_1) dx$$

Minimize $P(\text{mistake}) !!$

If $\underline{P(x, C_1)} > \underline{P(x, C_2)}$, assign that x to C_1 ($x \in R_1$)

if $P(x, C_2) > P(x, C_1)$, " x to C_2 ($x \in R_2$)

$$P(x, C_1) = P(C_1 | x) P(x), \quad P(x, C_2) = P(C_2 | x) P(x)$$

In order to minimize $P(\text{mistake})$,

assign x to the class for which $P(C_k | x)$ is largest.



General case of K classes

$$\begin{aligned} P(\text{correct}) &= \sum_{k=1}^K P(x \in R_k, C_k) \\ &= \sum_{k=1}^K \int_{R_k} P(x, C_k) dx \end{aligned}$$

$P(x, C_k) (\propto P(C_k | x))$ 가 충분히 R_k 를 대표

1.5.2 Minimizing the expected loss

cost function (loss)

The optimal solution is the one which minimizes the loss function

		예측		Confusion matrix
		양	음	
실제	양	0	100	
	음	1	0	

E.g. loss matrix

Let L_{kj} be the loss in case that true class is C_k

but it is assigned to class C_j (generally $L_{kk} = 0$)

$$\mathbb{E}[L] = \sum_k \sum_j \int_{R_j} L_{kj} P(x, c_k) dx$$

/
 c_k 이 있는데 c_j 를 판정한 경우

Minimize $\mathbb{E}[L]$ iff minimize $\sum_k L_{kj} P(c_k | x)$



1.5.3 The reject option

$$P(x, c_k)$$

$$\propto P(c_k | x)$$

1.5.4 Inference and decision

Inference stage : model for $P(C_k | \mathbf{x})$ target input $P(\mathbf{x}, C_k)$

Decision stage : function mapping \mathbf{x} into decision
(decision or discriminant)

Identify three distinct approaches

(a) Find $\underline{P(C_k | \mathbf{x})}$

- class conditional densities $P(\mathbf{x} | C_k)$

- prior class probabilities $P(C_k)$ (simple)

-

$$P(C_k | \mathbf{x}) = \frac{P(\mathbf{x} | C_k) P(C_k)}{P(\mathbf{x})}$$

-

$$P(\mathbf{x}) = \sum_k P(\mathbf{x} | C_k) P(C_k)$$

- Equivalently, model the joint distribution $P(\mathbf{x}, C_k)$ directly

and then obtain $P(C_k | \mathbf{x})$

- **Generative model**: modeling the distributions of \mathbf{x}, C_k

(b) Find $P(C_k | \mathbf{x})$ directly, apply decision theory

- Call it discriminative model
- Almost deep learning models

$S : S$ یعنی

\mathbf{x} test

(c) Find discriminant function $P(C_k | \mathbf{x})$??

- mapping each \mathbf{x} directly onto a class label
- probabilities play no role

Compensating for class prior

- Imbalanced data classification (difficult to generalize)

- $$PCC_E(x) = \frac{P(x|C_E) P(C_E)}{P(x)}$$

Combining models

E.g. Medical diagnosis

- blood test, x-ray image

- 두 정보를 통합 하기 vs >Hyp 전단 (사후 확률) 모델 만들기

Conditionally independent

A and B are conditionally independent given C ($A \perp\!\!\!\perp B | C$)

$$P(A, B | C) = P(A | C) \cdot P(B | C)$$

prop $A \perp\!\!\!\perp B | C \Leftrightarrow P(A | B, C) = P(A | C)$

양 or 정상

Suppose $x_I \perp\!\!\!\perp x_B | C_k$ (assumptions of Naive Bayes classifier)

$$\begin{aligned} \Rightarrow P(C_k | x_I, x_B) &\propto \underline{P(x_I, x_B | C_k)} \\ &= \underline{P(x_I | C_k)} P(x_B | C_k) P(C_k) \\ &\propto \underline{P(C_k | x_I)} P(C_k | x_B) / P(C_k) \end{aligned}$$

1.5.5 Loss functions for regression

Prediction $y(x)$ of the value t for each x

Let $L(t, y(x))$ be the loss function. (depends on $y(x)$)

$$\mathbb{E}[L] = \iint L(t, y(x)) p(x, t) dx dt \approx \frac{1}{n} \sum L(t_n, y(x_n))$$

In case $L(t, y(x)) := \{y(x) - t\}^2$

$$\mathbb{E}[L] = \iint \{y(x) - t\}^2 p(x, t) dx dt \approx \frac{1}{n} \sum (t_n - y(x_n))^2$$

$y(x)$

MSE

Goal: choose $y(x)$ as so minimize $E[L]$

Assume $y(x)$ is completely flexible

$y \rightarrow L \rightarrow \text{expected loss}$

$$\frac{\delta E[L]}{\delta y(x)} = 2 \int \{y(x) - t\} p(x,t) dt = 0$$

Functional F , $F[y] = \int G(y(x), y'(x), x) dx$

CD.5)

$$\Rightarrow \frac{\delta F}{\delta y} = \frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right)$$

Since $E[L] = \int \underbrace{\{(y(x) - t)^2 p(x,t) dt\}}_{G(x,y,\dot{y})} dx$

$$\frac{\delta \mathbb{E}[L]}{\delta y(x)} = \frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial \dot{y}} \right) = 2 \int \{y(x) - t\} p(x,t) dt$$

1

Continue to minimize $\mathbb{E}[L]$ for $y(x)$

In case
 $\{y(x) - t\}^2$

$$y(x) = \frac{\int t p(x,t) dt}{p(x)} = \int t \underline{p(t|x)} dt =: \mathbb{E}_t[t|x] \quad \checkmark$$

function of x

We used $\int p(x,t) dt = p(x)$.

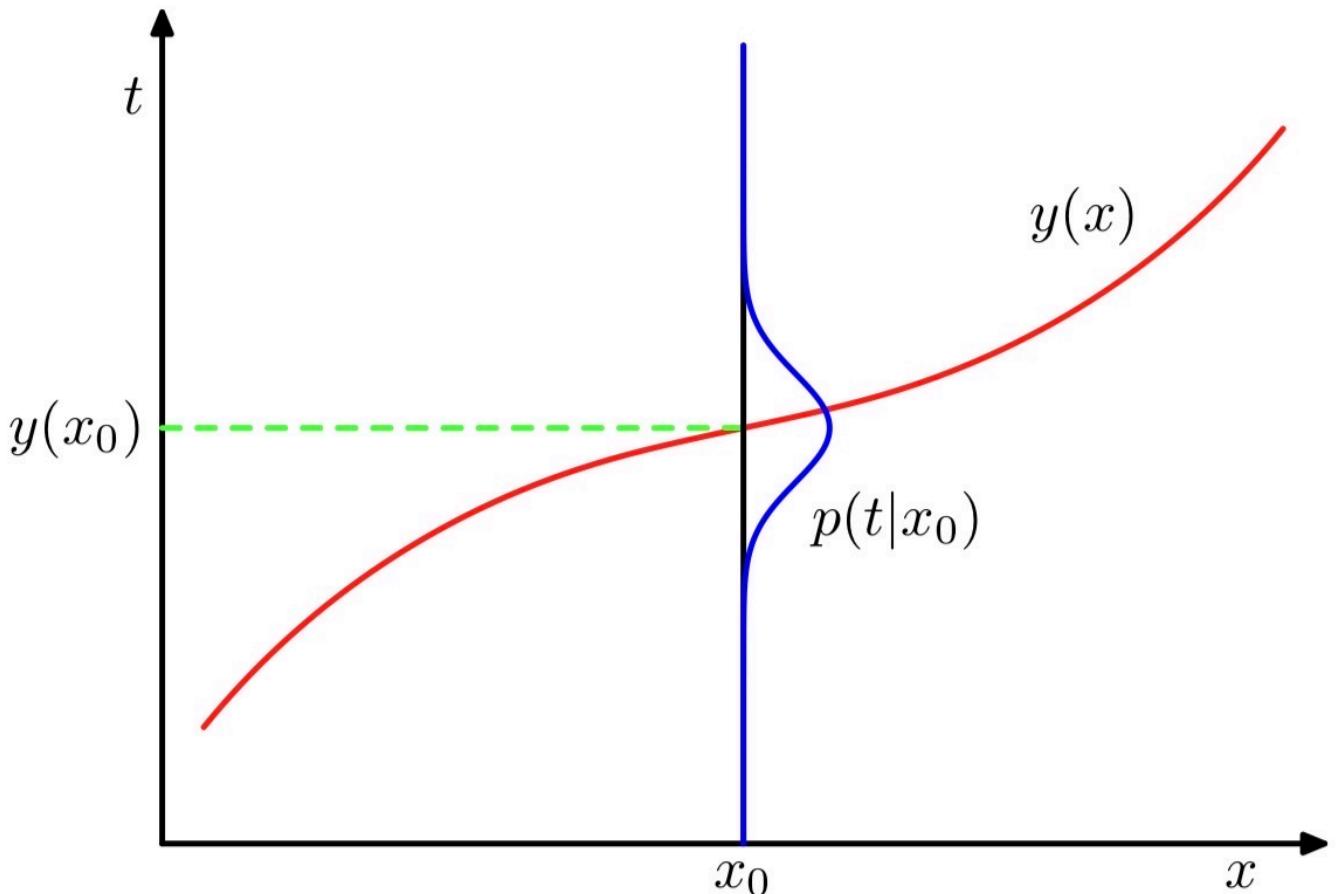
For multivariate target $\#$, the optimal solution can be expressed as

$$\mathbb{E}_{\#}[\#|x]$$

$L = \{y(x) - t\}^2$ 인 경우 (MSE) $E[L]$ (expected loss) \equiv

최소화 하는 $y(x) = E_t[t|x] = \int t p(t|x) dt$ (conditional exp.)
x의 대상 함수

Figure 1.28 The regression function $y(x)$, which minimizes the expected squared loss, is given by the mean of the conditional distribution $p(t|x)$.



Note that

$$\begin{aligned} \{y(x) - t\}^2 &= \{y(x) \pm E[t|x] - t\}^2 \\ &= \{y(x) - E[t|x]\}^2 + 2\{y(x) - E[t|x]\}\{E[t|x] - t\} \\ &\quad + \{E[t|x] - t\}^2 \end{aligned}$$

$$\textcircled{1} \quad \mathbb{E}[y(x)] = \mathbb{E}[t|x] \int^2 p(x,t) dt \quad p(x,t) = p(t|x) \cdot p(x)$$

$$= \int y(x) - \mathbb{E}[t|x] \int^2 p(x) dx \cdot \underbrace{\int p(t|x) dt}_{=1}$$

$$\textcircled{2} \quad \mathbb{E}[y(x)] = \mathbb{E}[t|x] \int (\mathbb{E}[t|x] - t) p(x,t) dt$$

$$= \int y(x) - \mathbb{E}[t|x] \int \left(\int (\mathbb{E}[t|x] - t) p(t|x) dt \right) p(x) dx$$

$$= \mathbb{E}[t|x] \int p(t|x) dt - \int t p(t|x) dt$$

$$= 0$$

$$\begin{aligned}
 ③ \iint \{ \mathbb{E}[t(x)] - t^2 \} p(x, t) dx dt &= \iint \{ \mathbb{E}[t|x] - t^2 \} p(t|x) p(x) dx dt \\
 &= \underbrace{\left(\iint \{ \mathbb{E}[t|x] - t^2 \} p(t|x) dt \right)}_{\text{mean}} \underbrace{p(x)}_{\text{ }} dx \\
 &= \underbrace{\int \text{var}[t|x] p(x) dx}_{\text{ }} \\
 &\quad \boxed{\text{typo}}
 \end{aligned}$$

Thus,

$$\mathbb{E}[L] = \iint \{ y(x) - \mathbb{E}[t(x)] \}^2 p(x) dx + \underbrace{\int \text{var}[t|x] p(x) dx}_{\text{without } y(x)}$$

minimum of $\mathbb{E}[L]$
 / noise

(optimal least squares predictor is given by the conditional mean)

Three distinct approaches to solving regression problems

$$y(x) = \mathbb{E}_t[t|x] = \int t p(t|x) dt \quad (1.89)$$

(a) Determining $p(x, t)$, normalizing to find $p(t|x)$.

Calculating the conditional mean given by (1.89)

(b) Solving the inference of determining $p(t|x)$ directly.

Calculating the conditional mean given by (1.89)

(c) Find $y(x)$ directly from the training data.

classic
✓

Bayesian
✓

Minkowski loss

$$\mathbb{E}[L_q] = \iint |y(x) - t|^q p(x, t) dx dt$$

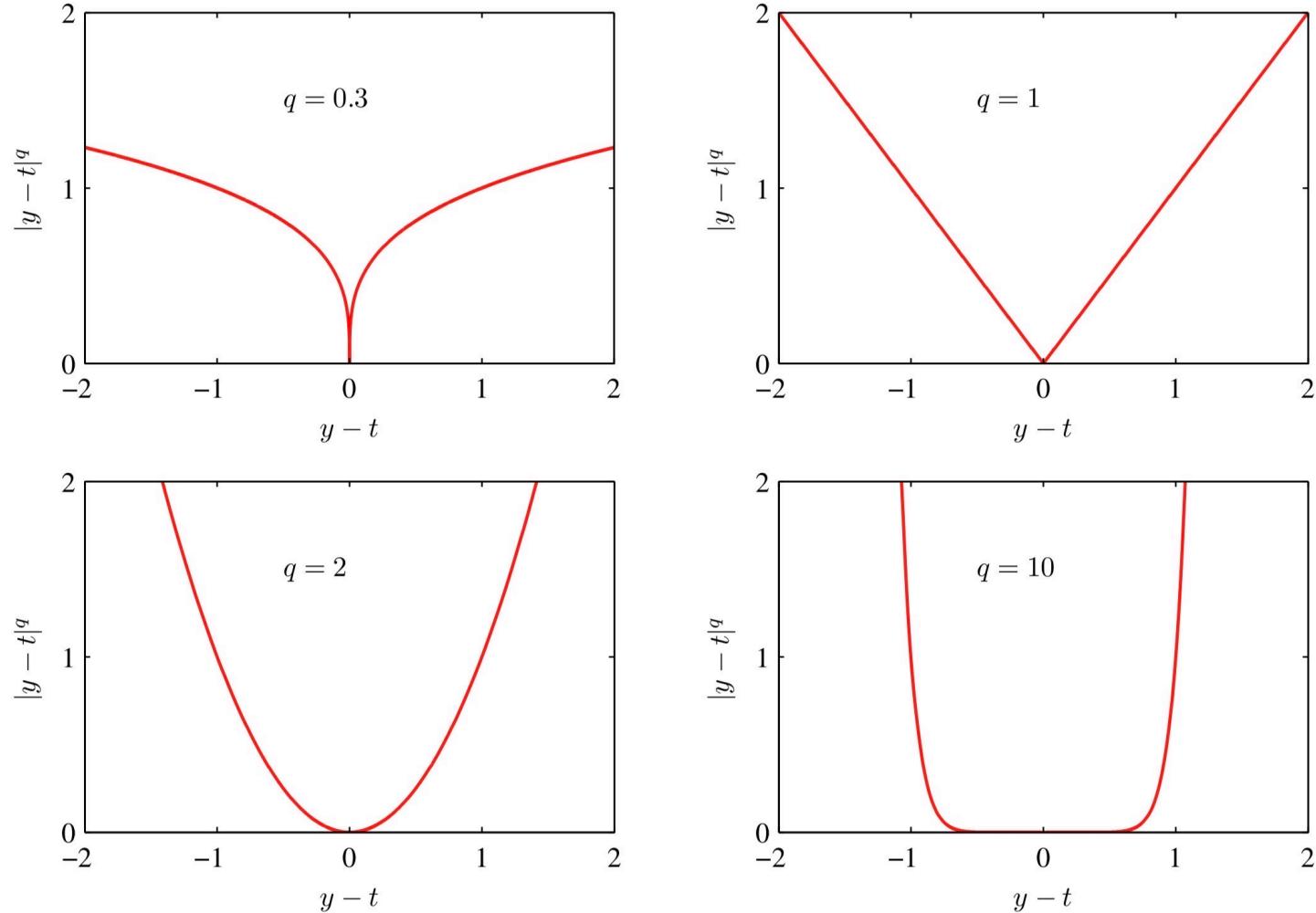


Figure 1.29 Plots of the quantity $L_q = |y - t|^q$ for various values of q .

1.6 Information Theory

Discrete random variable x

The amount of information can be view as 'degree of surprise'

on learning x . (가능성 낮은 사건에 크게 놀람)

우리 나라가 월드컵 우승할 확률 vs 브라질이 우승할 확률

사건이 발생했을 때 놀라는 정도

Measure of information depends on $p(x)$

$h(x)$ express the information content

Consider the unrelated (independent) events x, y

$$h(x,y) = h(x) + h(y) \quad \text{각각 따로 일어났을 때 정보량의 합}$$

함께 일어났을 때 얻는 정보량

(When $x \perp\!\!\!\perp y$, $p(x,y) = p(x) \cdot p(y)$)

Let

$$\underline{h(x)} := -\log_2 p(x). \quad (\text{amount of information})$$

$$-\int p(x) \ln p(x) dx$$

Entropy

$$H[x] = -\sum_i p(x_i) \log_2 p(x_i)$$

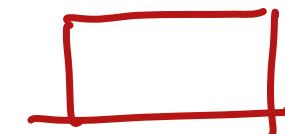
$E[h(x)]$

정보 전송에 필요한
평균 정보량

Take $p(x) \log_2 p(x) := 0$ when $p(x) = 0$.

Eg. discrete random variable x having 8 possible states
(length 3 bits)

① uniform prob $(\frac{1}{8}, \frac{1}{8}, \dots, \frac{1}{8})$ case



$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = \underline{3}$$

② prob $(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64})$



$$H[x] = \underline{2}$$

We see that non uniform distribution has a smaller entropy than uniform one.

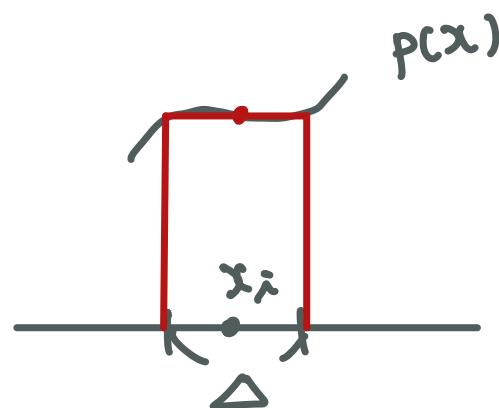
When x is discrete rv, uniform distribution has the maximum entropy.

Introduction to entropy analysis



Let us define entropy of conti. rv x with pdf $p(x)$.

Divide x into bins of width Δ and assume $p(x)$ is continuous. By mean value theorem, $\exists x_i \in [i\Delta, (i+1)\Delta]$ s.t



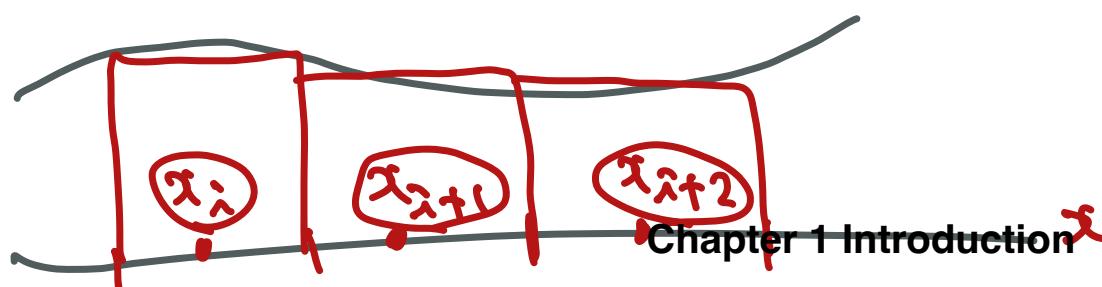
$$\int_{x_0}^{(i+1)\Delta} p(x) dx = p(x_i) \Delta$$

i th

(convert to step function)

Assigning any value $x \in i$ th bins to x_i

\Rightarrow the prob. of observing the values in λ^{th} bins = $P(x_\lambda) \Delta$



Entropy of such discrete distribution

$\log_2 \rightarrow \ln$

Probability

$$H_\Delta = - \sum_i p(x_i) \Delta \ln(p(x_i) \Delta) = - \sum_i p(x_i) \Delta \ln p(x_i) - \underline{\ln \Delta}$$

where we used $\sum_i p(x_i) \Delta = \int p(x) dx = 1$.

Now we omit $-\underline{\ln \Delta}$ and consider $\Delta \rightarrow 0$

$$\lim_{\Delta \rightarrow 0} \left\{ - \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx$$

$\mathbb{E}[h(x)]$

called 'differential entropy' (continuous form of entropy)

Multivariate continuous variable

$$H[x] := - \int p(x) \ln p(x) dx$$

$x \rightarrow \infty$

In the case of discrete distributions, uniform distribution has the maximum entropy.

Let us consider the maximum entropy configuration for a continuous variable.

$$(-\infty, \infty) \quad r.v \quad x$$

Constraints

$$p(x) = ?$$

$$\int p(x) dx = 1$$

$$\int x p(x) dx = \mu$$

$$\int (x - \mu)^2 p(x) dx = \sigma^2$$

fix

Appendix D

Using Lagrange multiplier and calculus of variations,
we set the derivative of this functional to zero giving

$$(*) \quad p(x) = \exp \{ -1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 \}$$

Apply $p(x)$ to constraints

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}$$

(*)

We need to maximize the following functional w.r.t $p(x)$

$$\begin{aligned} & \boxed{- \int p(x) \ln p(x) dx} + \pi_1 \left(\int p(x) dx - 1 \right) + \pi_2 \left(\int x p(x) dx - \mu \right) \\ & + \pi_3 \left(\int (x - \mu)^2 p(x) dx - \sigma^2 \right) \end{aligned}$$

Let

$$\begin{aligned} F(p(x)) &:= \int \{ -p(x) \ln p(x) + \pi_1 p(x) + \pi_2 x p(x) + \pi_3 (x - \mu)^2 p(x) \} dx \\ &+ (-\pi_1 - \pi_2 \mu - \pi_3 \sigma^2) \\ \text{functional} &= \int G(p(x)) dx + c \end{aligned}$$

$$\frac{\delta F(P(x))}{\delta P(x)} = \frac{\partial G(P)}{\partial P} = -\ln P(x) - P(x) \cdot \frac{1}{P(x)} + \pi_1 + \pi_2 x + \pi_3 (x - \mu)^2$$

$$= 0$$

$$\ln P(x) = \pi_1 + \pi_2 x + \pi_3 (x - \mu)^2 - 1$$

(appendix D)

변별법

Differential entropy of the Gaussian

$$\underline{H[x] = \frac{1}{2} \{ 1 + \ln(2\pi\sigma^2) \}}$$

$$\Sigma - \underbrace{P(x) \ln P(x)}_{\geq 0}$$

X, Y two random variables with joint prob $P(X, Y)$

- $\ln P(Y|X)$: additional information needed to specify Y , when X is given.

$P(Y|X)$

Conditional entropy of Y given X

$$H[Y|X] := -\iint p(y, x) \ln p(y|x) dy dx$$

We have

$$H[X, Y] = H[Y|X] + H[X]$$

(X, Y) 의 정보량
= X 정보량 +
 X 가 주어졌을 때
 Y 를 듣는 정보량

1. 6. 1 Relative entropy and mutual information

Unknown prob. dist. $p(x)$ (true dist.)

Approximating prob. dist. $q(x)$

Amount of additional information needed to send x using $q(x)$

$$KL(P \parallel q) := - \int p(x) \ln q(x) dx - \underbrace{\left(- \int p(x) \ln p(x) dx \right)}_{P \text{ 정보량}}$$

unknown

approximate

$$= - \int p(x) \ln \left\{ \frac{q(x)}{p(x)} \right\} dx$$

Kullback - Leibler divergence (relative entropy)

Remark

- $KL(p \parallel q) \neq KL(q \parallel p)$ not symmetric
- $KL(p \parallel q) \geq 0$ (*)
- $KL(p \parallel q) = 0$ iff $p = q$
- measurement of difference between two prob. dists.

(*) φ : convex function, f : a function of x
prob. measure

$$\varphi\left(\int f(x) p(x) dx\right) \leq \int \varphi(f(x)) p(x) dx$$

Jensen's
inequality

put $p(x) := -\ln(x)$. $f(x) = \frac{q(x)}{p(x)}$

Now suppose \mathbf{x} is generated from unknown $p(\mathbf{x})$.

We try to approximate $p(\mathbf{x})$ using $q(\mathbf{x}|\theta)$ (e.g. parametric dist.)

and observed

N points $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$\int f(x) p(x) dx \approx \frac{1}{N} \sum f(x_n)$$

Monte Carlo Integration

$$KL(P||Q) = - \int p(x) \ln \left\{ \frac{q(x)}{p(x)} \right\} dx \approx \sum_{n=1}^N \left\{ -\ln q(\mathbf{x}_n | \theta) + \underbrace{\ln p(\mathbf{x}_n)}_{\text{indep. of } q} \right\}$$

To minimize $KL(P||Q)$, we need to maximize

$$\sum \ln q(\mathbf{x}_n | \theta) \quad (\text{log likelihood})$$

θ 를 매개변수로 q 를 풀어서 $\mathbf{x}_1, \dots, \mathbf{x}_N$ 이 나올 확률

$$\ln \left(\prod q(\mathbf{x}_n | \theta) \right) \quad \text{likelihood}$$

Degree of Independence

$$I[X, Y] := \underbrace{KL(P(X, Y) \parallel P(X)P(Y))}_{\text{mutual information between } X, Y}$$

$$0 \leq I[X, Y] = H[X] - H[X|Y]$$

$$X \perp\!\!\!\perp Y \Rightarrow H[X]$$

$$= H[Y] - H[Y|X]$$

$$= H[X|Y]$$

The reduction in uncertainty about X as a consequence of the new Y .

