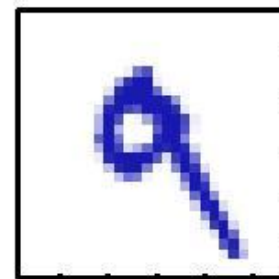
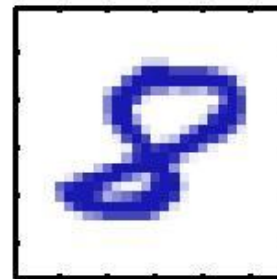
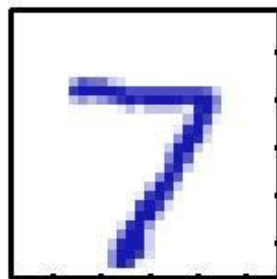
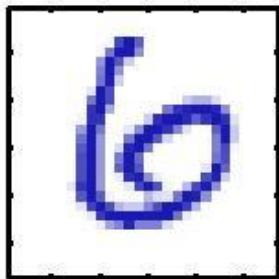
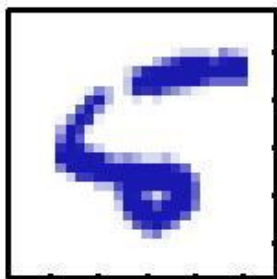
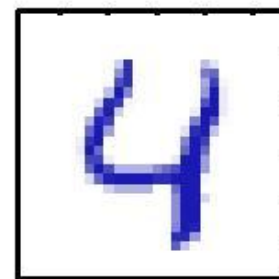
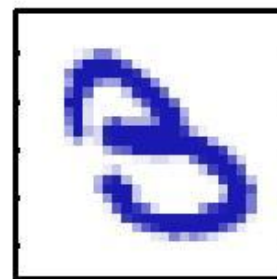
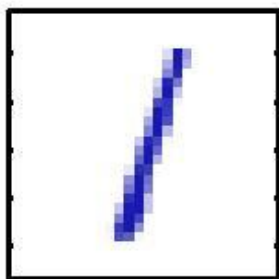
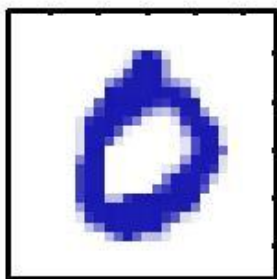

PATTERN RECOGNITION AND MACHINE LEARNING

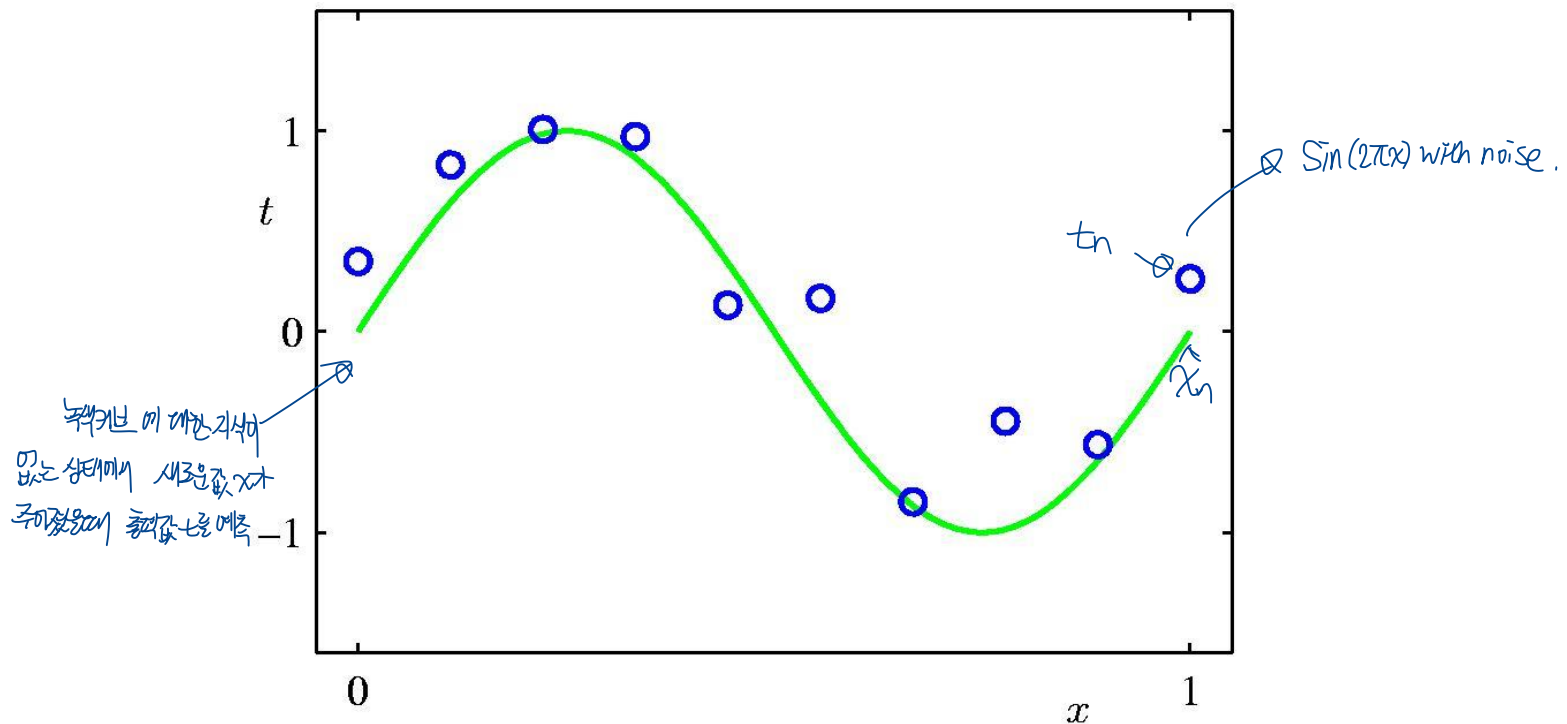
CHAPTER 1: INTRODUCTION

Example

Handwritten Digit Recognition

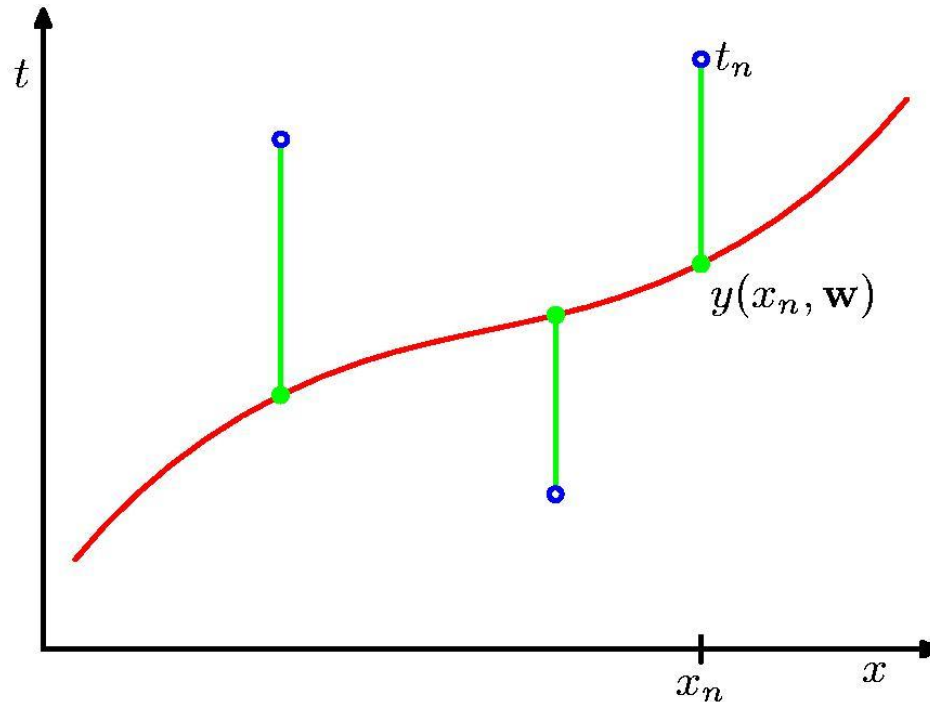


Polynomial Curve Fitting



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Sum-of-Squares Error Function



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

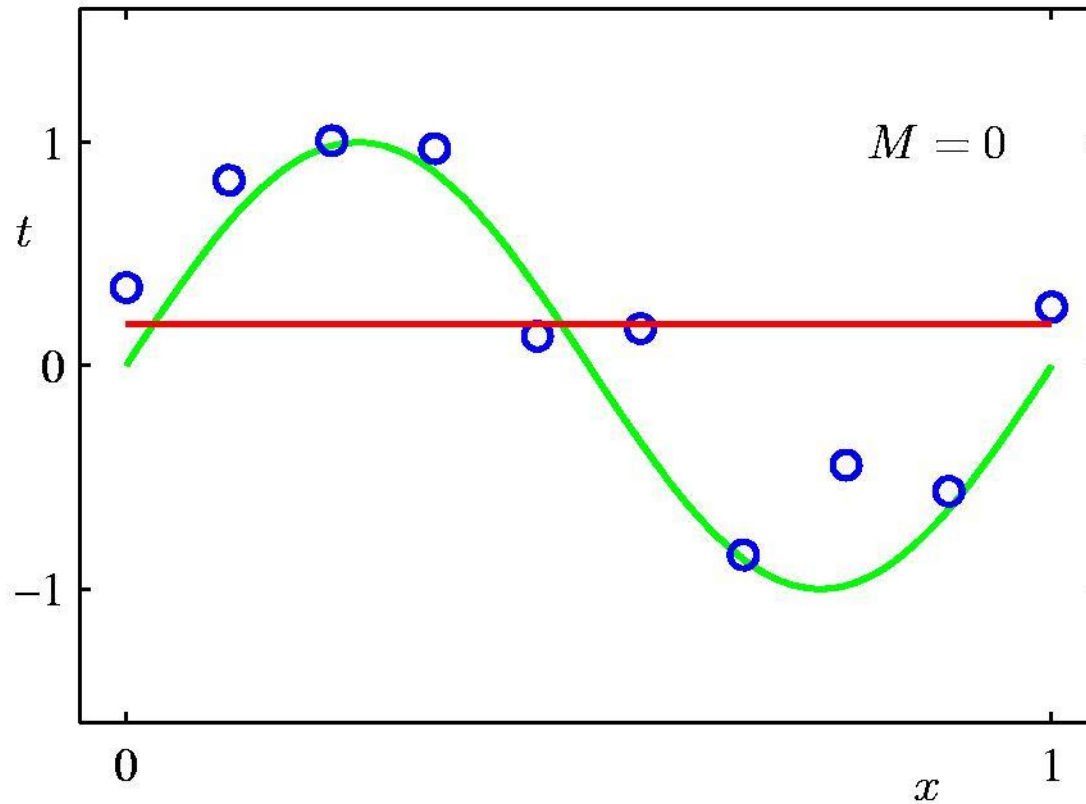
$E(\mathbf{w})$ 를 최소화하는 \mathbf{w} 선택, 미분 등 \mathbf{w} 에 대한 산점도에 $E(\mathbf{w})$ 를 최소화하는 \mathbf{w}^* 를 찾을 수 있다.

0th Order Polynomial

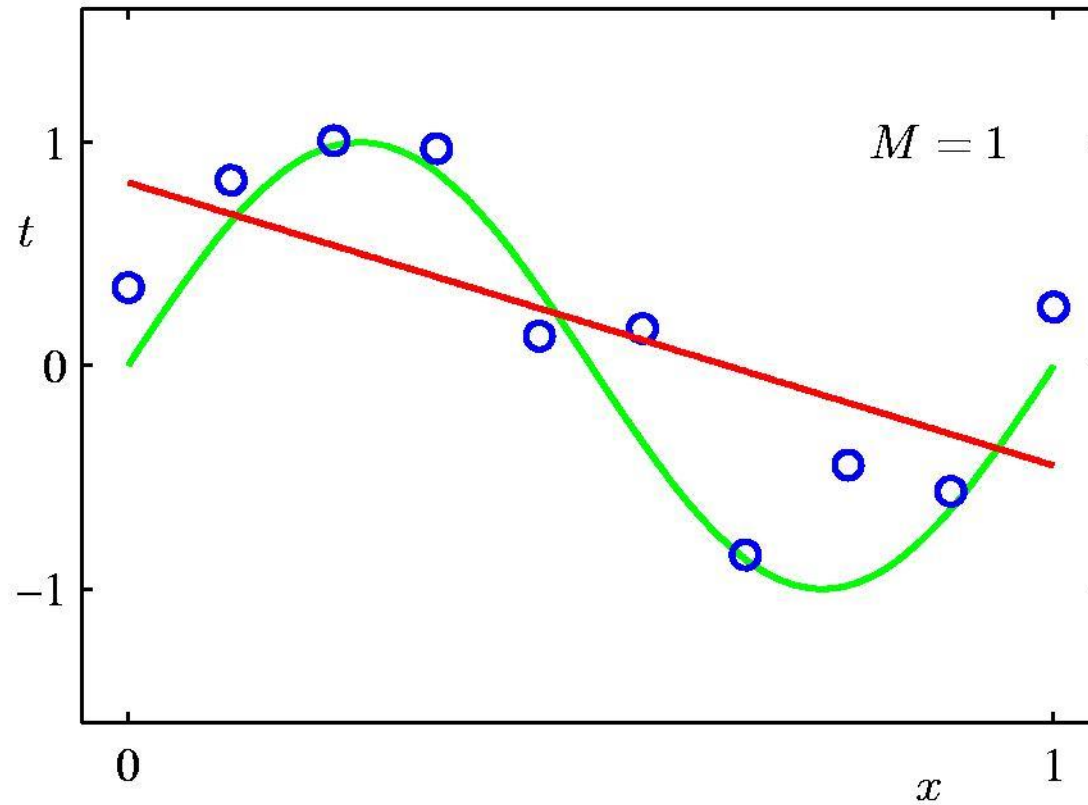
모델링

과적합(over-fitting)

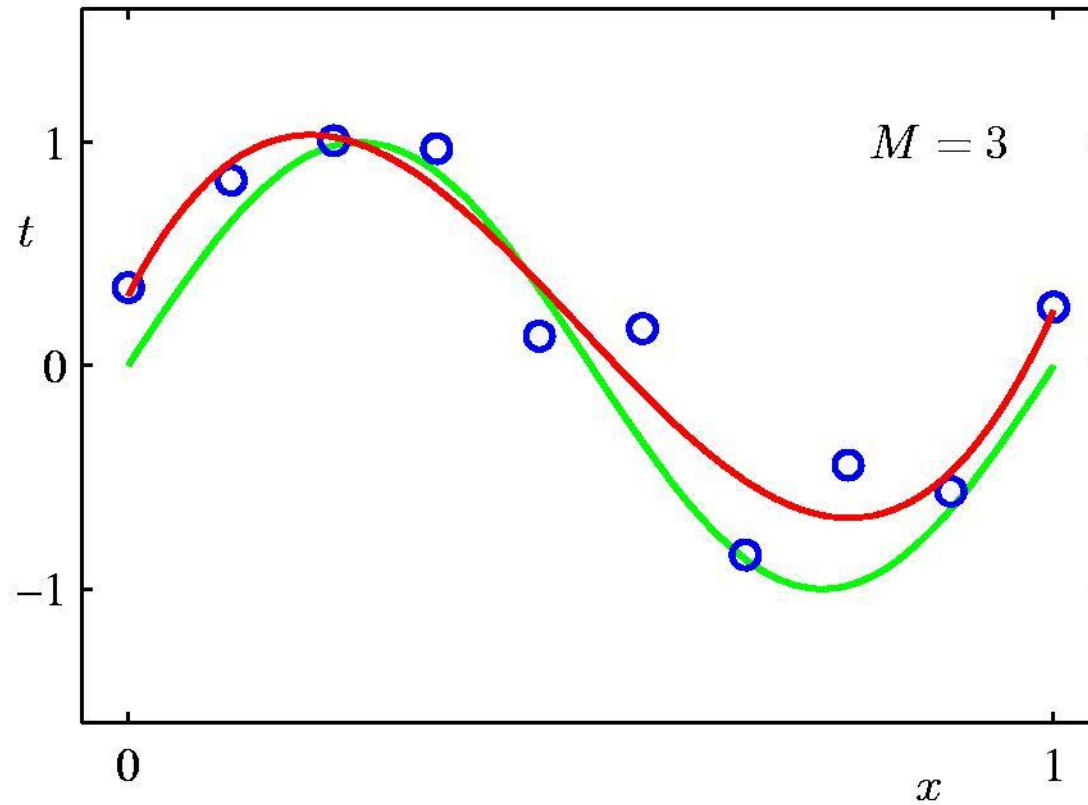
과소적합(under-fitting)



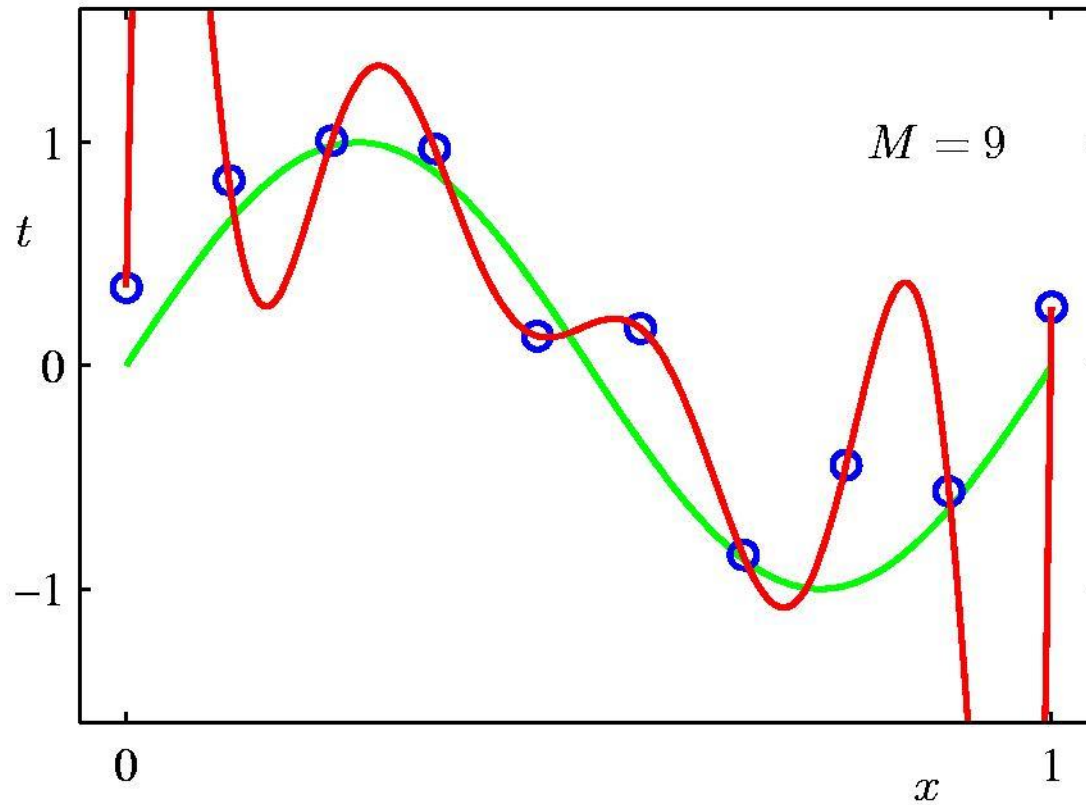
1st Order Polynomial



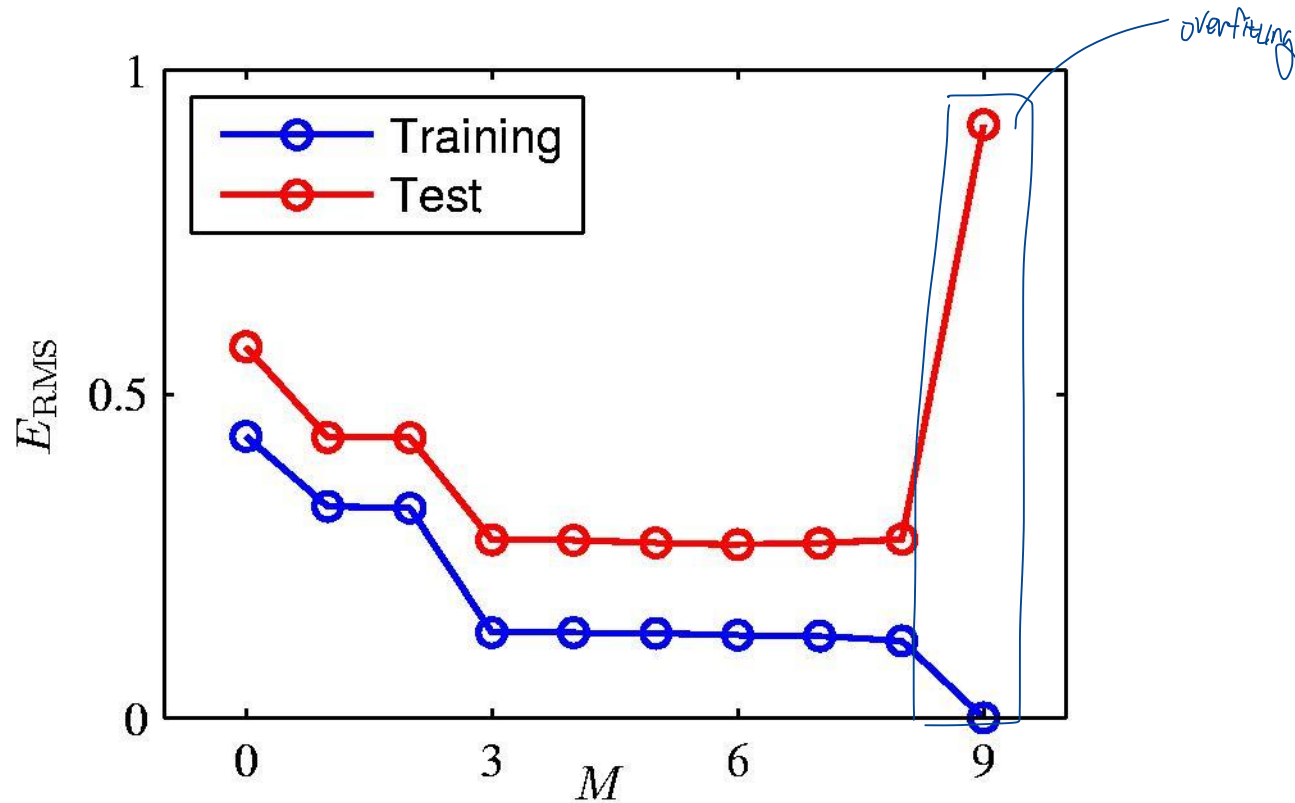
3rd Order Polynomial



9th Order Polynomial



Over-fitting



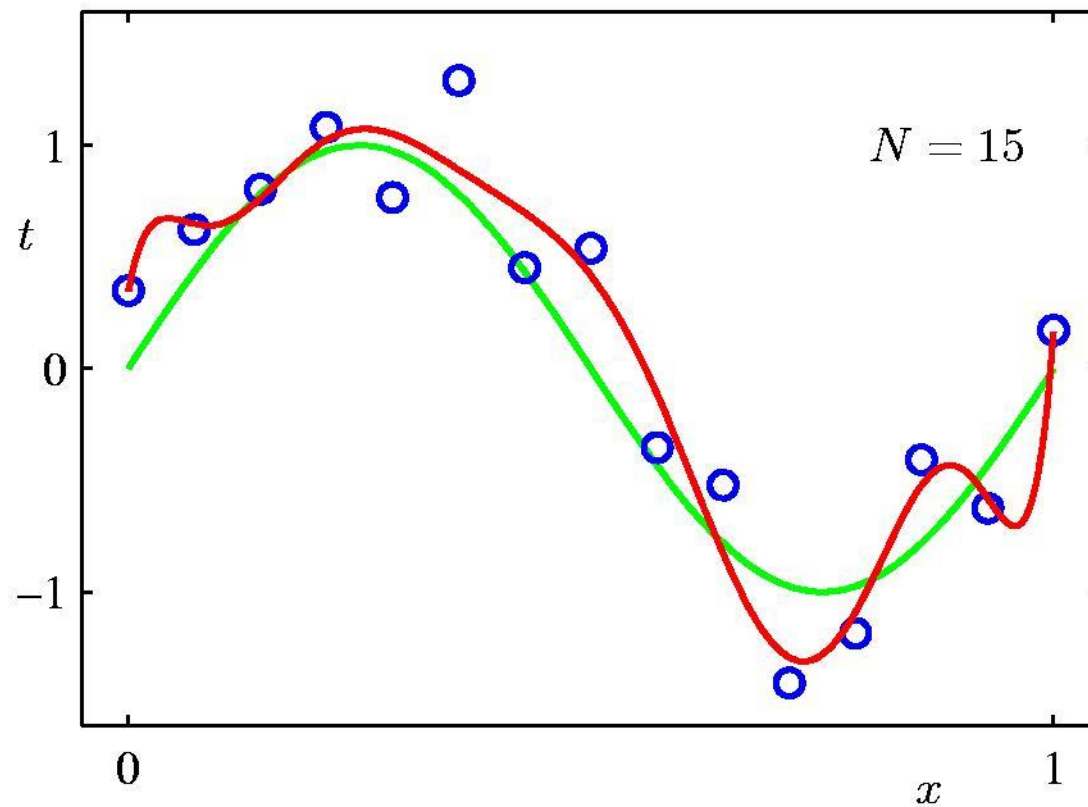
Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	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

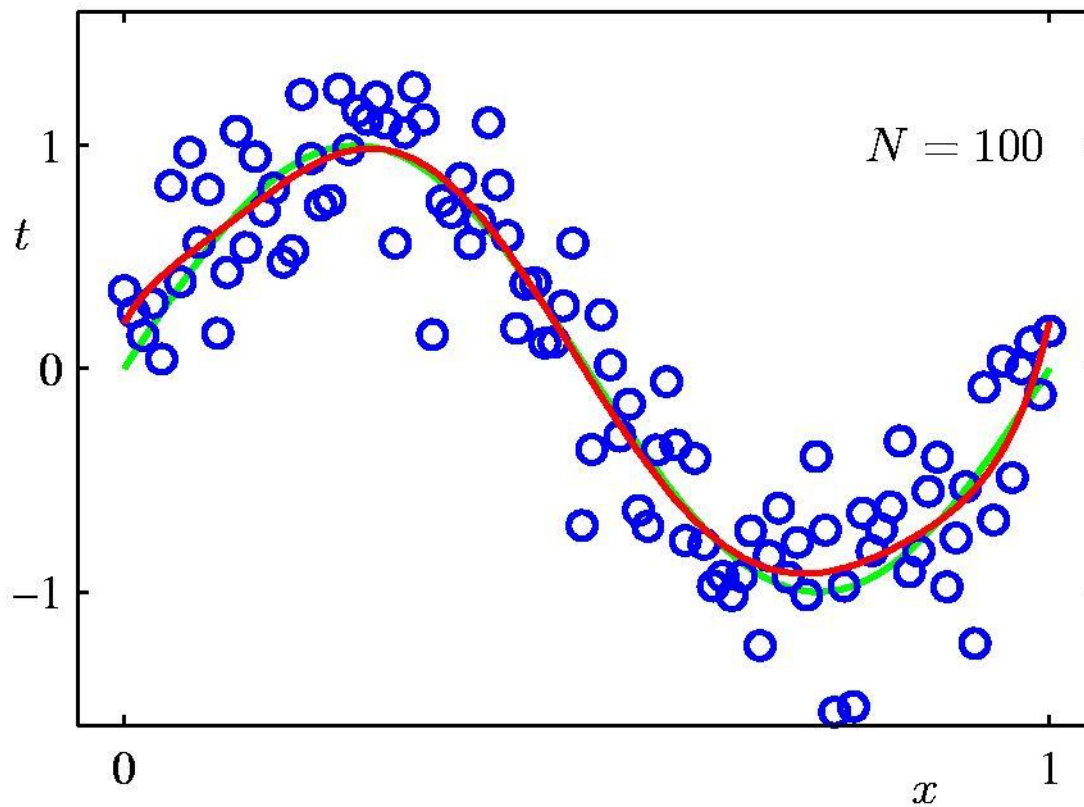
Data Set Size: $N = 15$

9th Order Polynomial



Data Set Size: $N = 100$

9th Order Polynomial *데이터를 높이기 과적합이 되어있다.*



Regularization 정규화

Penalize large coefficient values

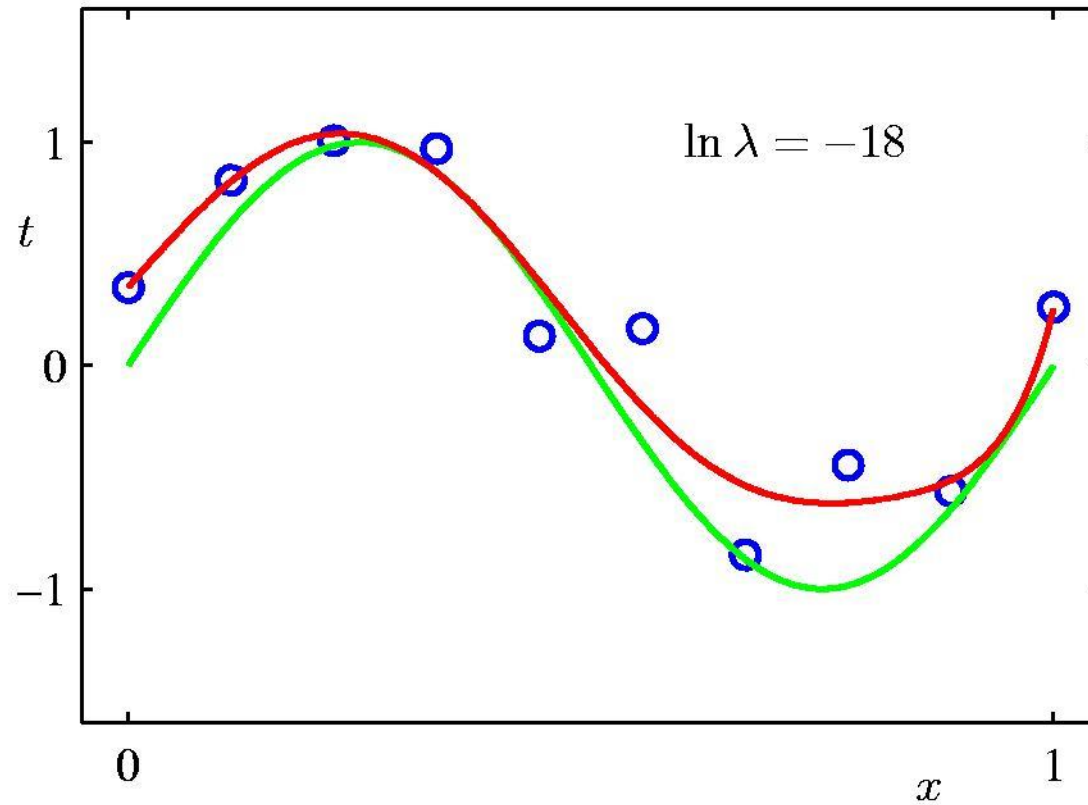
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{정규화}}$$

↗ also called, weight decay
가중치 감쇠.

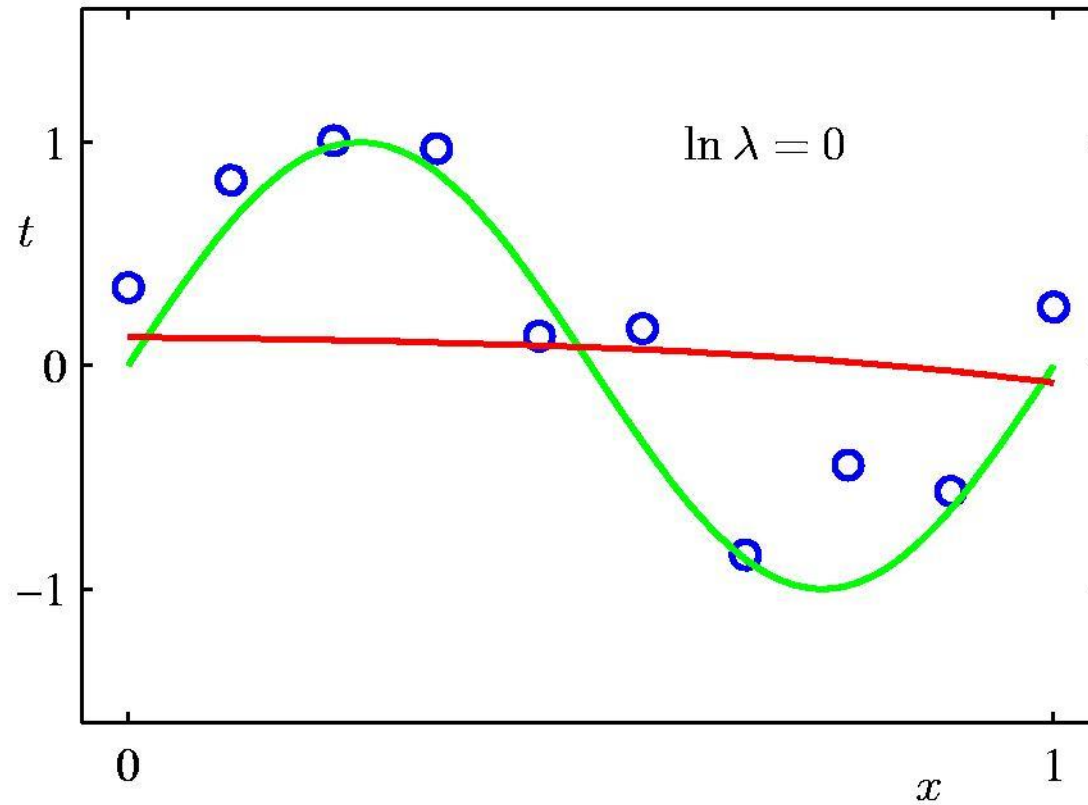
↳ ridge 회귀

↳. 일차함 → 2차함수

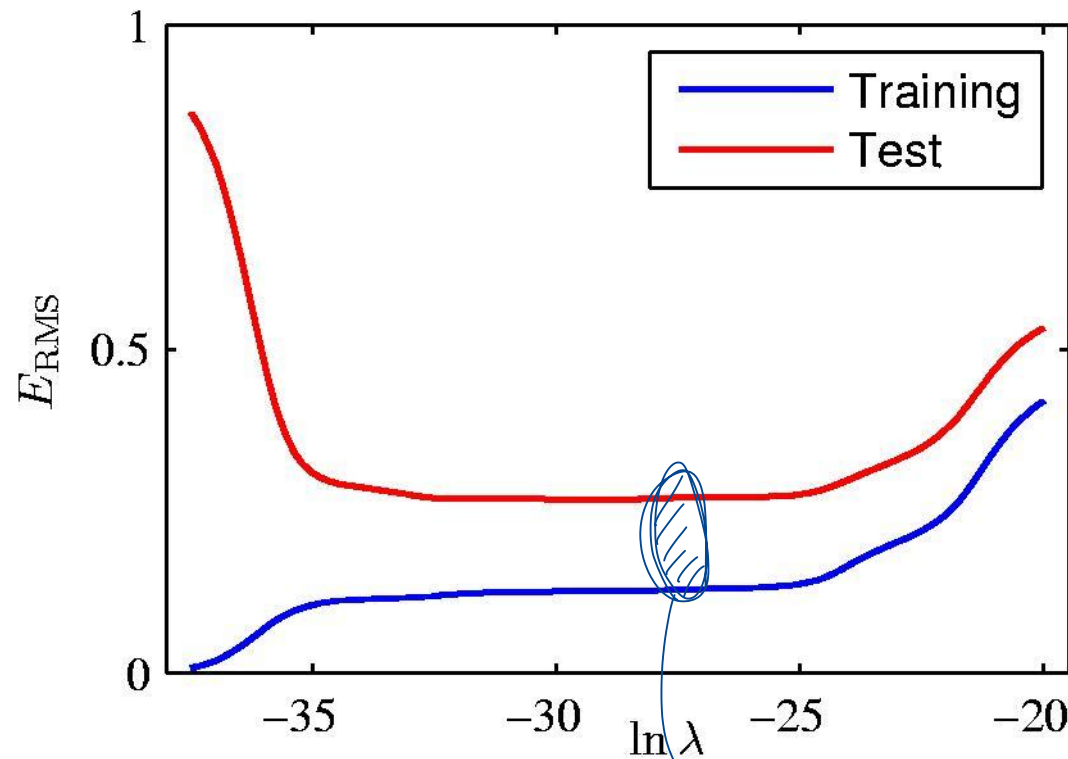
Regularization: $\ln \lambda = -18$



Regularization: $\ln \lambda = 0$



Regularization: E_{RMS} vs. $\ln \lambda$



올바른 hyperparameter를 찾아야 한다.

Polynomial Coefficients

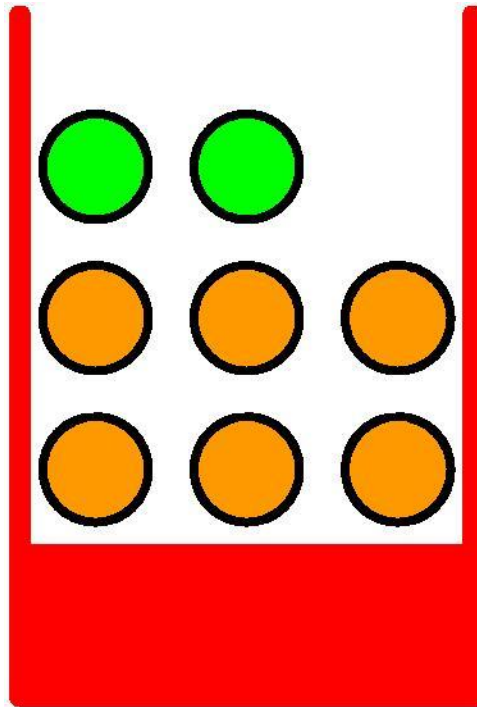
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01

Probability Theory

불확실성 \rightarrow 확률론 통한 해소.

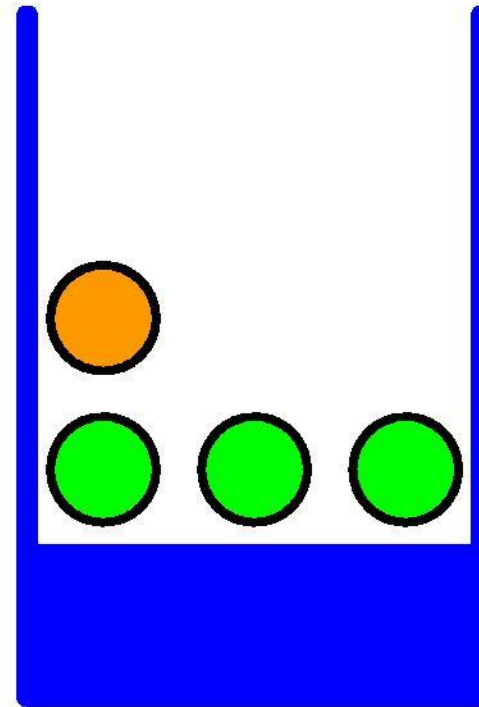
Apples and Oranges

- 불확실성
- 과일 색깔과 맛은 모두 주어진.



각 색깔을 고를 확률.

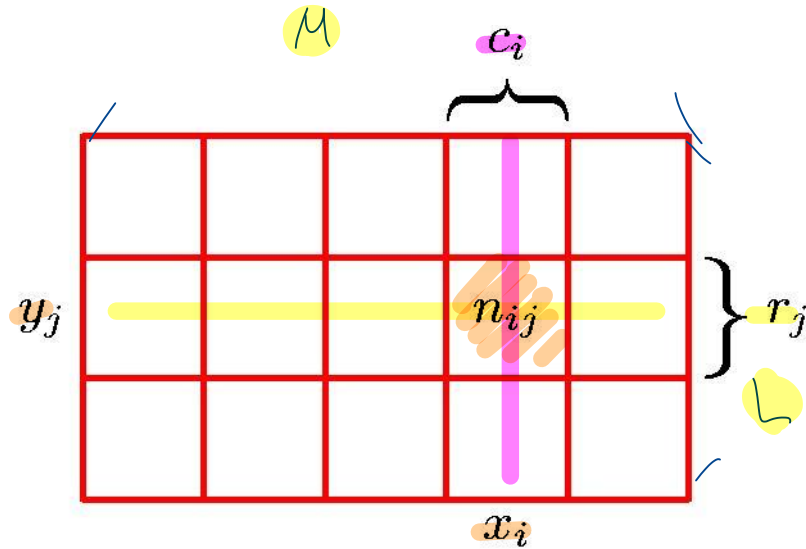
40%



60%

B - & 상자에 대한 확률 변수
F - & 과일 맛 고를 확률.

Probability Theory



Marginal Probability

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

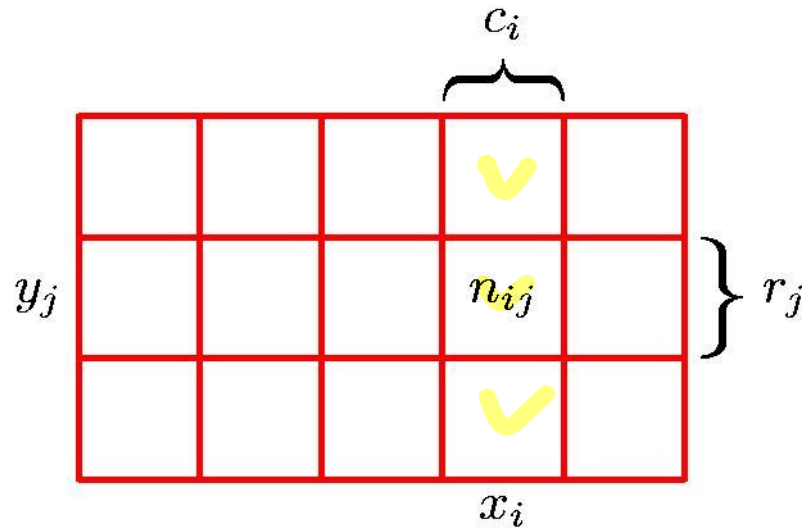
Joint Probability

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

Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Probability Theory



Sum Rule

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij}$$

$$= \sum_{j=1}^L p(X = x_i, Y = y_j)$$

Product Rule

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

$$= p(Y = y_j | X = x_i) p(X = x_i)$$

$\frac{P(Y=y_j | X=x_i)}{P(X=x_i)} \cdot P(X=x_i)$

The Rules of Probability

Sum Rule

주변확률 (marginal probability)

$$p(X) = \sum_Y p(X, Y)$$

Product Rule

합동확률 (joint probability)

$$p(X, Y) = \underbrace{p(Y|X)}_{\text{조건부 확률 (conditional probability)}} p(X)$$

Bayes' Theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

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

$$p(X) = \frac{p(X|Y)p(Y)}{p(Y|X)}$$

posterior \propto likelihood \times prior

다. 독립

$$P(X, Y) = P(X)P(Y)$$

각각의 샘플이 같은 수의 사고와 오작

→ 사고와 오작을 고를 확률은
이런 샘플을 학습과 관련 없다.

고려해야 할 것.

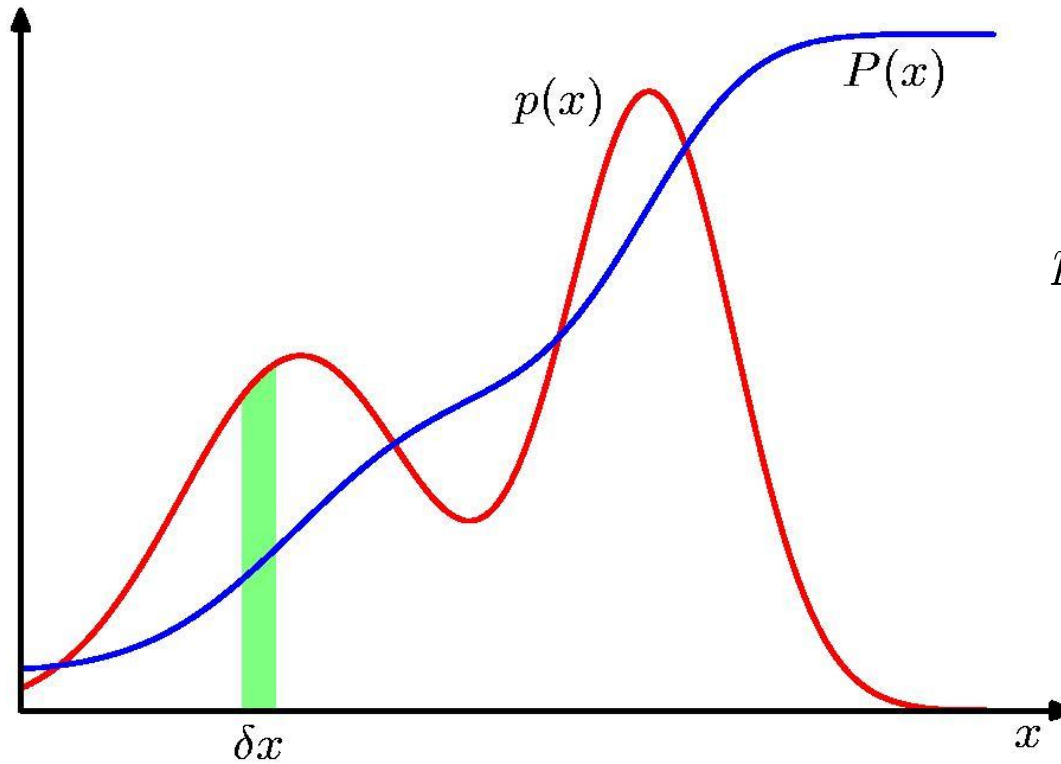
→ 어떤 과일이 선택되었는지 알기 위해 어떤 바스를 선택했는가? → $P(B)$

$\frac{4}{10}$ vs $\frac{6}{10}$

$$P(B|F) \propto P(F|B) \times P(B)$$

Probability Densities

연속적인 변수에 대한 확률



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

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

↑
cdf (누적분포함)

$$P'(x) = p(x)$$

e.g. 열역학.

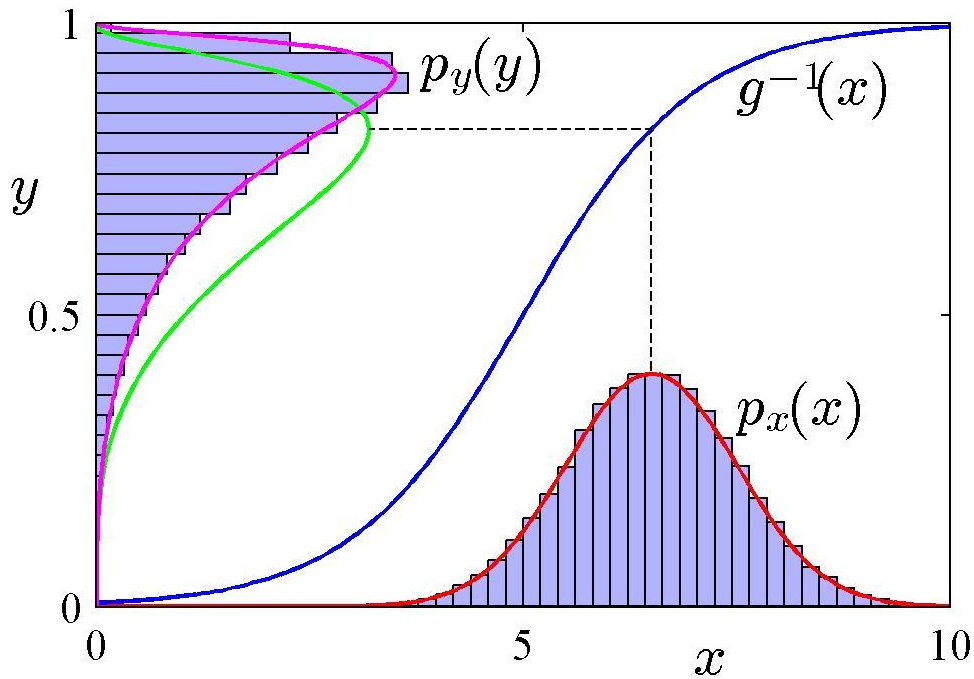
$$P(x) = \int p(x, y) dy$$

$$p(x, y) = p(y|x)p(x)$$

$$p(x) \geq 0$$

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

Transformed Densities



$$x = g(y)$$

$$f(x) = f(g(y))$$

$$(x, x + \delta x)$$

$$(y, y + \delta y)$$

$$p_x(x) \delta x \approx p_y(y) \delta y$$

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right|$$

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

→ 확률 밀도의 변환은 미분변수를 생각하기에 따라서 달라질 수 있다.

<연속>

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

Conditional Expectation (discrete)

Approximate Expectation (discrete and continuous)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Variances and Covariances

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

제곱의 평균 - 평균의 제곱.

예) $\text{Var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$

$$\begin{aligned} \text{cov}[x, y] &= \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}] \\ &= \mathbb{E}_{x,y} [xy] - \mathbb{E}[x]\mathbb{E}[y] \end{aligned}$$

$$\begin{aligned} \text{cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\} \{\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]\}] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\mathbf{x} \mathbf{y}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{y}^T] \end{aligned}$$

→ 두 확률 변수 x 와 y 가 비대칭적 공분산 형성된다.

- $\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$ ^[1]

[증명]

$$(X - \mu_X)(Y - \mu_Y) = XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y$$
를 이용하면

$$\begin{aligned}\text{Cov}(X, Y) &= \mathbb{E}(XY) - \mathbb{E}(\mu_X Y) - \mathbb{E}(\mu_Y X) + \mathbb{E}(\mu_X \mu_Y) \\ &= \mathbb{E}(XY) - \mu_X \mathbb{E}(Y) - \mu_Y \mathbb{E}(X) + \mu_X \mu_Y \\ &= \mathbb{E}(XY) - \mu_X \mu_Y\end{aligned}$$

- $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$

[증명]

◦ 일반화: $\text{Var}\left(\sum_{k=1}^n X_k\right) = \sum_{k=1}^n \text{Var}(X_k) + 2 \sum_{i < j} \text{Cov}(X_i, X_j)$

- $|\text{Cov}(X, Y)| \leq \sqrt{\text{Var}(X) \cdot \text{Var}(Y)}$

2.1. 모공분산

[편집]

모공분산은 **모집단**의 공분산이다. $\text{Cov}(X, Y)$ 또는 σ_{XY} 로 쓴다. X 와 Y 는 확률 변수, N 은 모집단의 표본의 개수, X_i 와 Y_i 는 각 확률 변수의 **도수**, μ 는 **모평균**을 뜻한다.

$$\begin{aligned}\text{Cov}(X, Y) &= \sigma_{XY} \\ &= \frac{1}{N} \sum_{i=1}^n (X_i - \mu_X)(Y_i - \mu_Y) \\ &= \mathbb{E}\{(X - \mu_X)(Y - \mu_Y)\}\end{aligned}$$

곧, 모공분산이란 X 의 편차와 Y 의 편차의 곱의 평균이다.

2.2. 표본공분산

[편집]

표본공분산은 표본집단의 공분산이다. S_{XY} 로 쓴다. X 와 Y 는 확률 변수, n 은 표본집단의 표본의 개수, X_i 와 Y_i 는 각 확률 변수의 **도수**, \bar{X} 와 \bar{Y} 는 표본평균을 뜻한다.

$$\begin{aligned} S_{XY} &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\ &= \mathbb{E}\{(X - \bar{X})(Y - \bar{Y})\} \end{aligned}$$

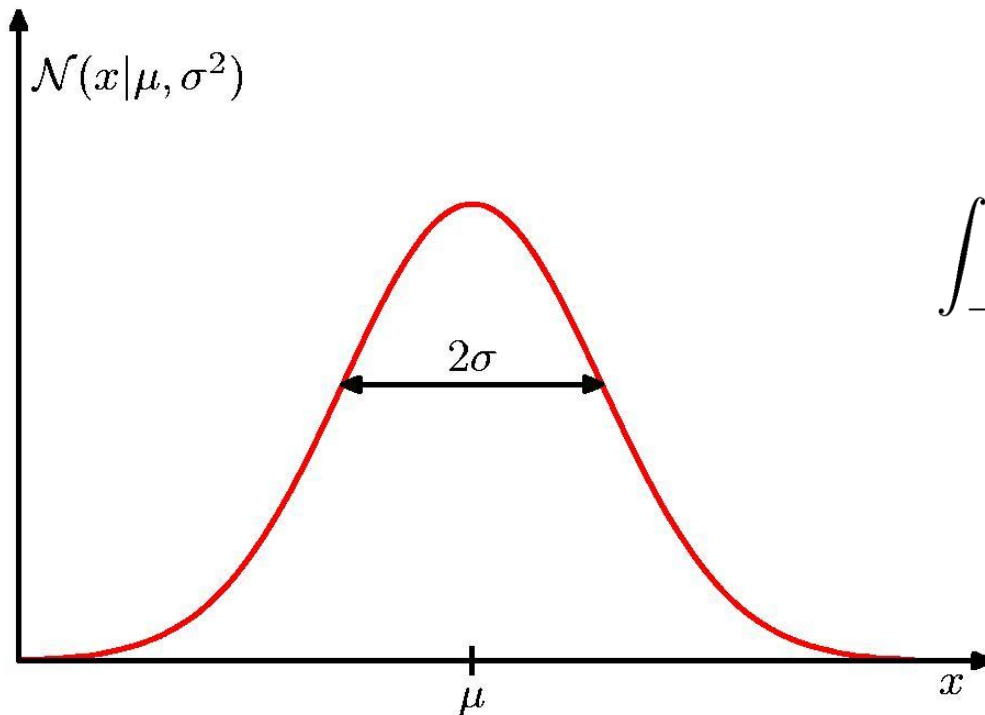
곧, 표본공분산이란 X 의 편차와 Y 의 편차의 곱의 평균이다. 주의할 점은 (표본의 개수) -1 로 나누는 것이다. n 이 아니라 $n - 1$ 로 나누는 것은 오차를 줄이기 위함으로, 일반적인 표본 분산의 계산법과 같다.

The Gaussian Distribution

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

↑ ↑
평균 분산

e.g. σ → 표준편차
 σ^2 → 분산



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

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

Gaussian Mean and Variance

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

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

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

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

From the definition of the Gaussian distribution, X has probability density function:

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

From the definition of the expected value of a continuous random variable:

$$E(X) = \int_{-\infty}^{\infty} x f_X(x) \, dx$$

So:

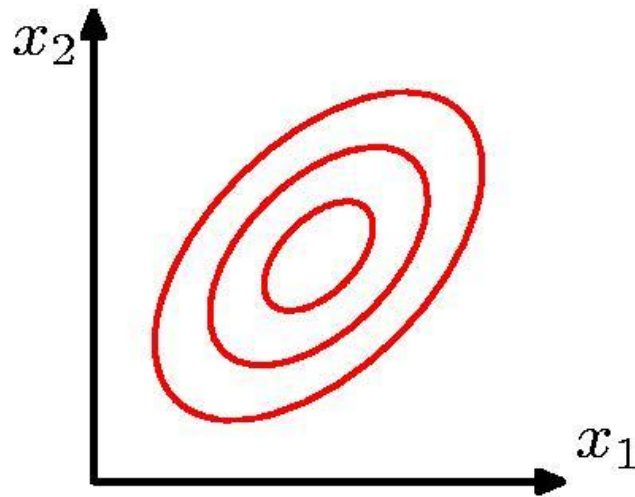
$$\begin{aligned} E(X) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} x \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \\ &= \frac{\sqrt{2}\sigma}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} (\sqrt{2}\sigma t + \mu) \exp(-t^2) dt && \text{substituting } t = \frac{x-\mu}{\sqrt{2}\sigma} \\ &= \frac{1}{\sqrt{\pi}} \left(\sqrt{2}\sigma \int_{-\infty}^{\infty} t \exp(-t^2) dt + \mu \int_{-\infty}^{\infty} \exp(-t^2) dt \right) \\ &= \frac{1}{\sqrt{\pi}} \left(\sqrt{2}\sigma \left[-\frac{1}{2} \exp(-t^2) \right]_{-\infty}^{\infty} + \mu \sqrt{\pi} \right) && \text{Fundamental Theorem of Calculus, Gaussian Integral} \\ &= \frac{\mu\sqrt{\pi}}{\sqrt{\pi}} && \text{Exponential Tends to Zero and Infinity} \\ &= \mu \end{aligned}$$

$$\begin{aligned} Var[X] &= \int_{-\infty}^{\infty} (x-\mu)^2 f_X(x) \, dx \\ &= \int_{-\infty}^{\infty} (x-\mu)^2 \frac{1}{\sqrt{2\pi}\sigma^2} e^{-(x-\mu)^2/2\sigma^2} \, dx \\ &= \frac{1}{\sqrt{2\pi}\sigma^2} \int_{-\infty}^{\infty} z^2 e^{-z^2/2\sigma^2} \, dz && (\because \text{let } z = x - \mu, \, dz = dx) \\ &= \left[-\frac{1}{\sqrt{2\pi}\sigma^2} z \sigma^2 e^{-z^2/2\sigma^2} \right]_{-\infty}^{\infty} + \frac{1}{\sqrt{2\pi}\sigma^2} \sigma^2 \int_{-\infty}^{\infty} e^{-z^2/2\sigma^2} \, dz \\ &= 0 + \sigma^2 && (\because \int_{-\infty}^{\infty} e^{-z^2/2\sigma^2} \, dz = \sqrt{2\pi\sigma^2}) \end{aligned}$$

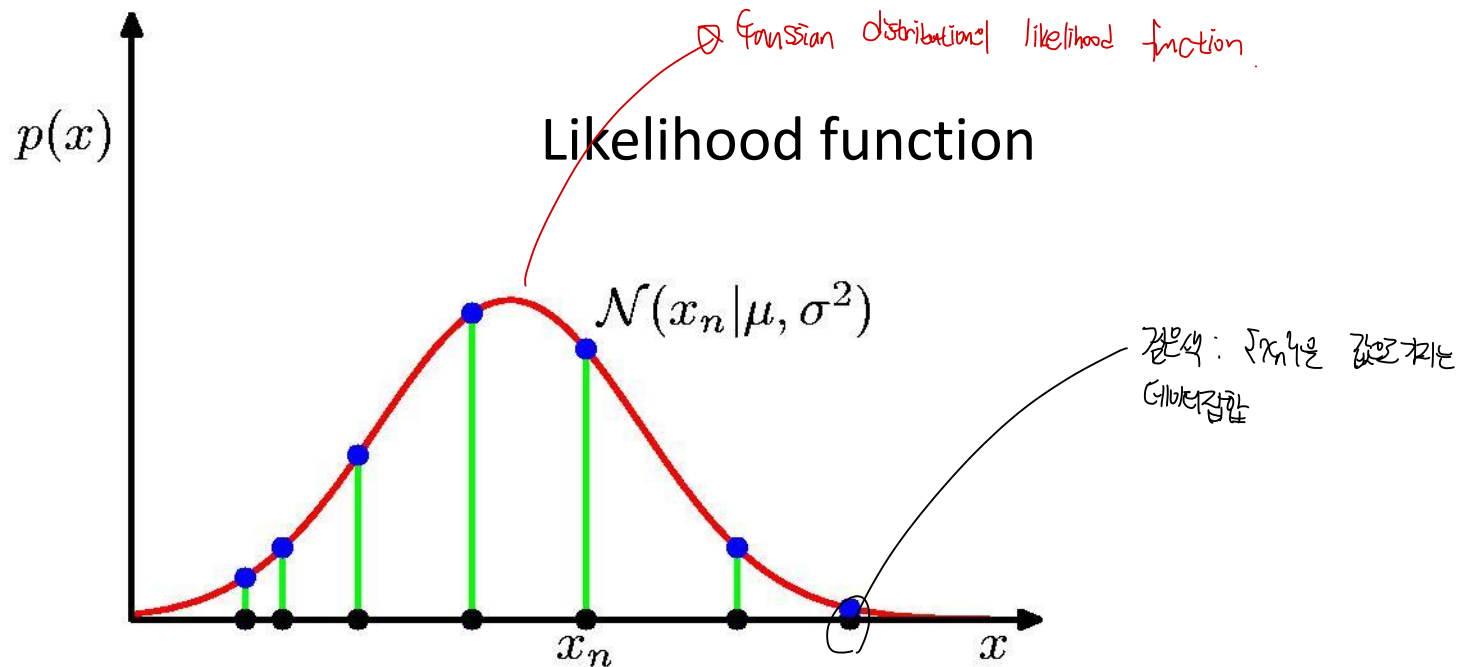
The Multivariate Gaussian

D차원 벡터 \mathbf{x} 는 \mathbb{R}^D 공간, $\mathbf{\mu}$ 는 평균, $\mathbf{\Sigma}$ 는 공분산

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



Gaussian Parameter Estimation



$$p(\mathbf{x} | \mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \sigma^2) \quad \text{집합 } x \text{가 i.i.d이다.}$$

파라미터 집합의 곱은 의미한다.

μ, σ^2 주어져서 해당 곱을 최대화할 때
자료를 최대화 함수 있다.

Maximum (Log) Likelihood

log는 변수에게 의존하는 함수 \therefore 서로 하면. 값을 찾는다는 의미가 없다.

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

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\frac{\partial L}{\partial \mu} = 0$$

표본평균 (Sample mean)

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

$$\frac{\partial L}{\partial \sigma^2} = 0$$

표본분산 (Sample variance)

Properties of μ_{ML} and σ_{ML}^2

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

분포의 중심을 $\frac{N-1}{N}$ 만큼 과소 평가해 있다.

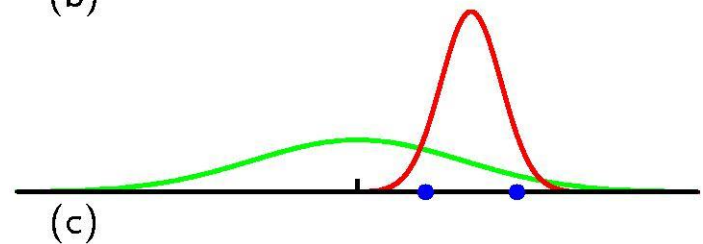
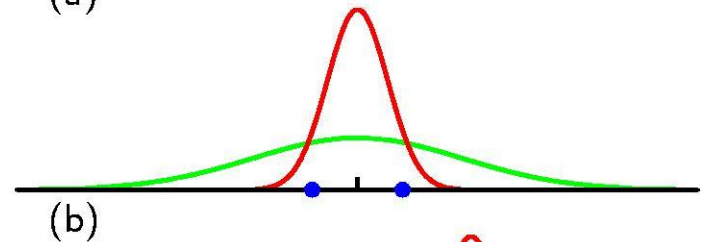
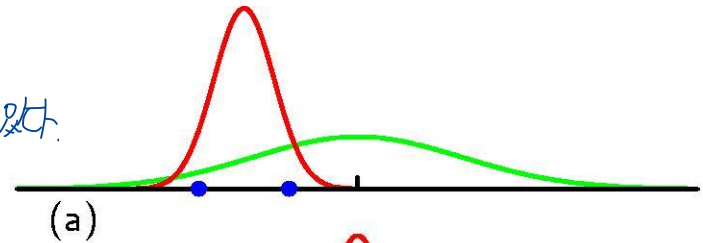
$$\mathbb{E}[\sigma_{\text{ML}}^2] = \left(\frac{N-1}{N} \right) \sigma^2$$

$\hat{\sigma}^2$ 분산량

$$= \frac{N}{N-1} \sigma_{\text{ML}}^2$$

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

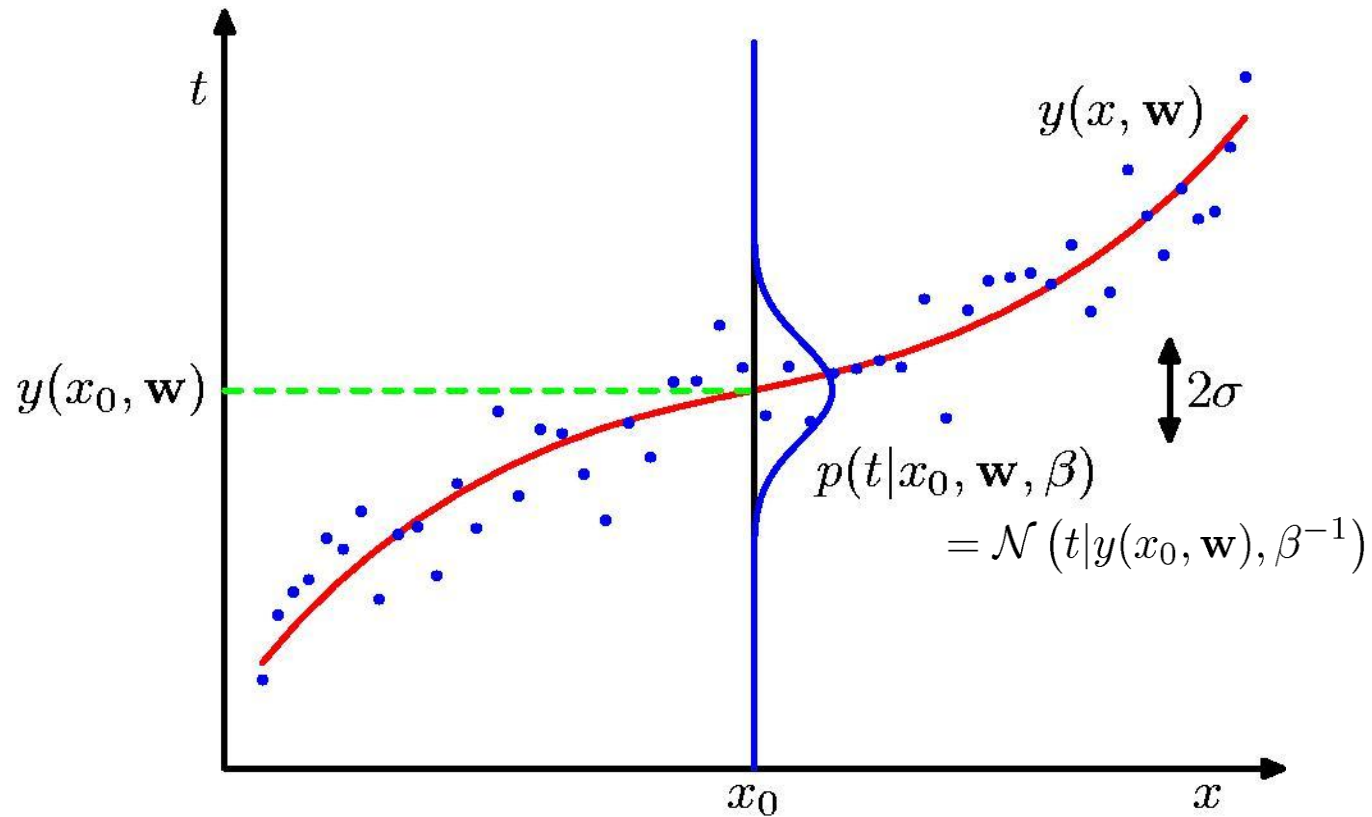
비편향



분산이 다름에 과소평가.

※ 두 추정값의 표준값이 실제 분산과 같아지면 unbiased, 다른 biased이다.

Curve Fitting Re-visited



Maximum Likelihood

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

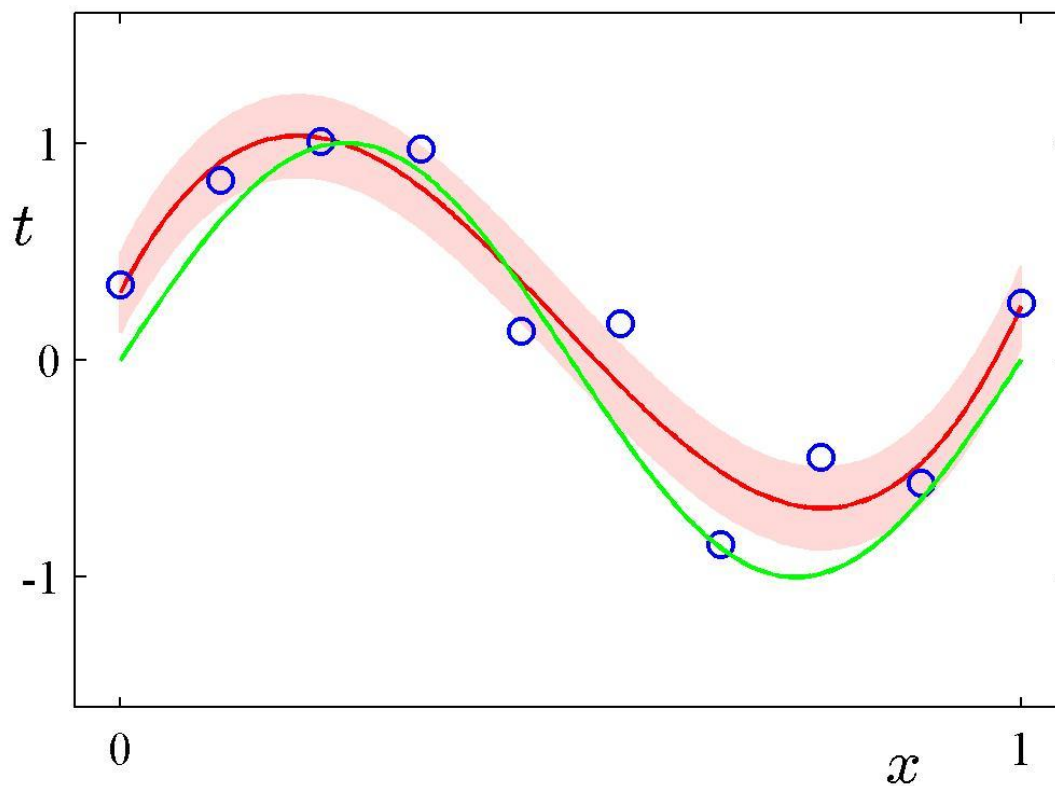
$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = - \underbrace{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

Determine \mathbf{w}_{ML} by minimizing sum-of-squares error, $E(\mathbf{w})$.

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$

Predictive Distribution

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$



MAP: A Step towards Bayes

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

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

$$\beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}$$

Determine \mathbf{w}_{MAP} by minimizing regularized sum-of-squares error, $\tilde{E}(\mathbf{w})$.

Bayesian Curve Fitting

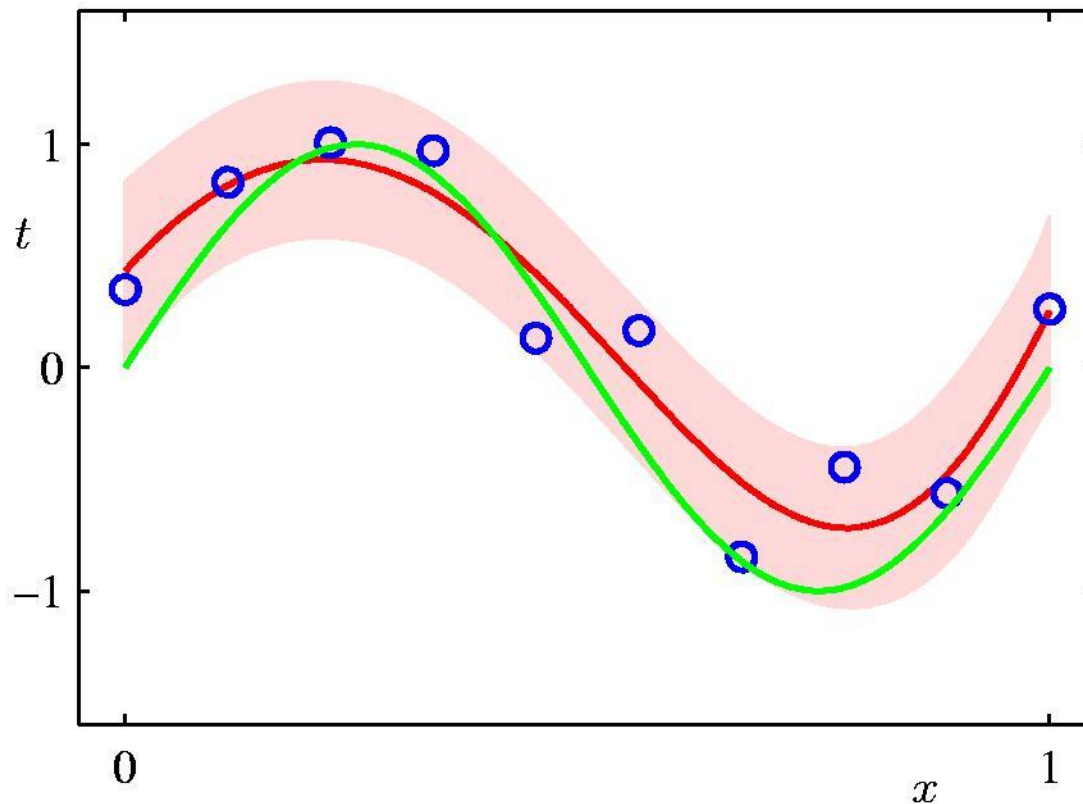
$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}(t|m(x), s^2(x))$$

$$m(x) = \beta \phi(x)^T \mathbf{S} \sum_{n=1}^N \phi(x_n) t_n \qquad s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x)$$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \qquad \phi(x_n) = (x_n^0, \dots, x_n^M)^T$$

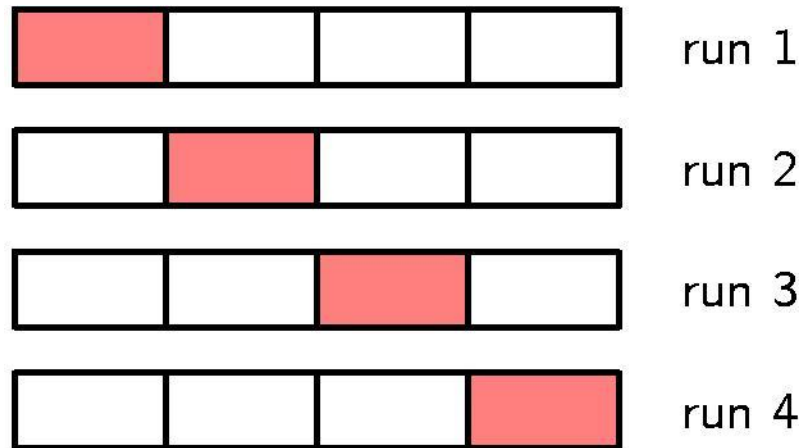
Bayesian Predictive Distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

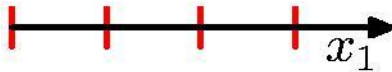


Model Selection

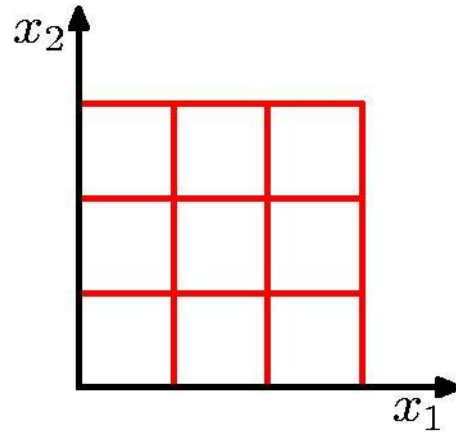
Cross-Validation



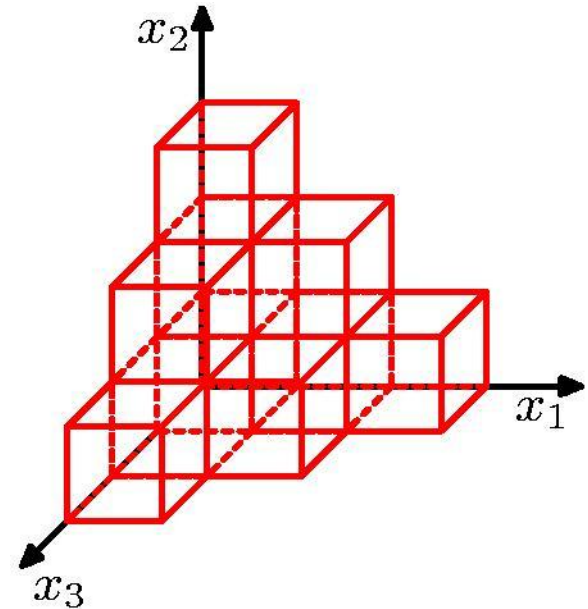
Curse of Dimensionality



$D = 1$



$D = 2$



$D = 3$

↓
形

Curse of Dimensionality

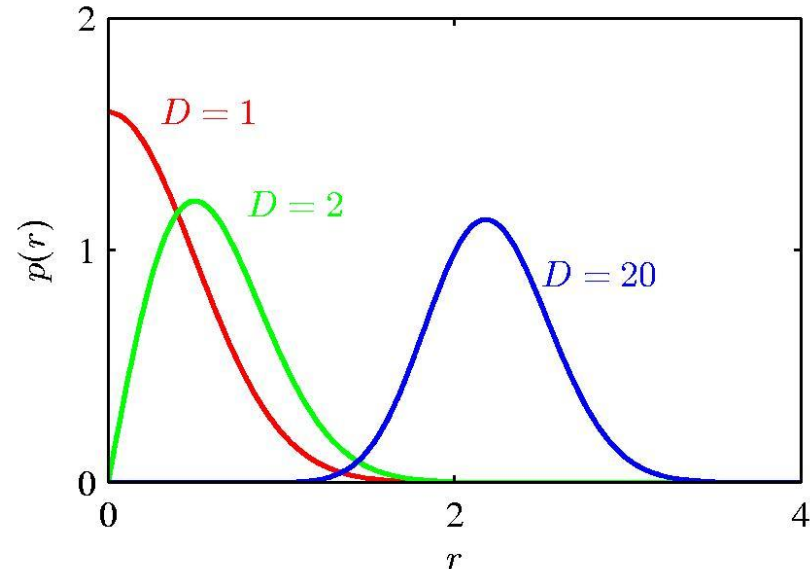
Polynomial curve fitting, $M = 3$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

Gaussian Densities in higher dimensions


→ 확률밀도가 거의 모든 근처의 값은 0에 가까워 있음을 확인할 수 있다.

→ Manifold 문제.



Decision Theory

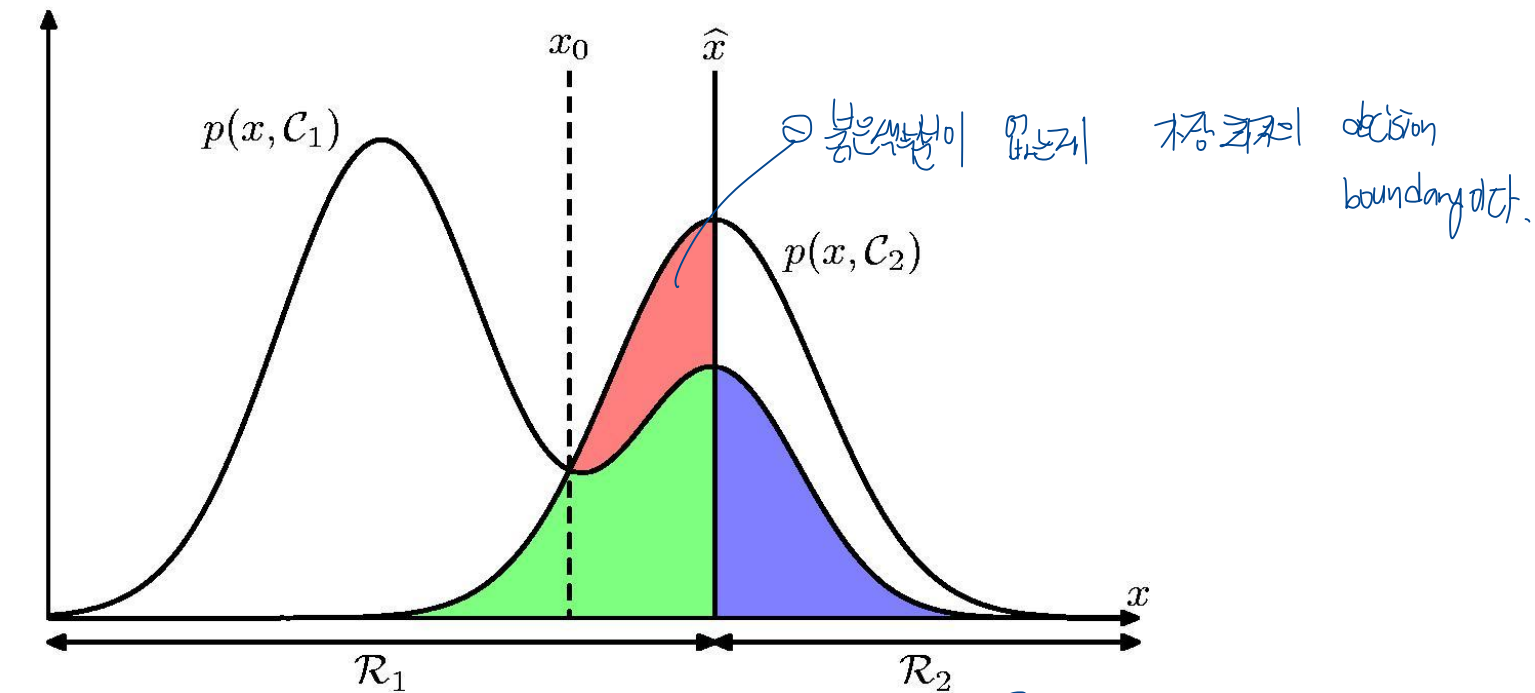
Inference step

Determine either $p(t|\mathbf{x})$ or $p(\mathbf{x}, t)$. 

Decision step

For given \mathbf{x} , determine optimal t .

Minimum Misclassification Rate



$$\begin{aligned}
 p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, C_2) + p(\mathbf{x} \in \mathcal{R}_2, C_1) \\
 &= \int_{\mathcal{R}_1} p(\mathbf{x}, C_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, C_1) d\mathbf{x}.
 \end{aligned}$$

Minimum Expected Loss

Example: classify medical images as 'cancer' or 'normal'

		Decision		
		cancer	normal	
Truth	cancer	0	1000	<i>손상</i>
	normal	1	0	

Minimum Expected Loss

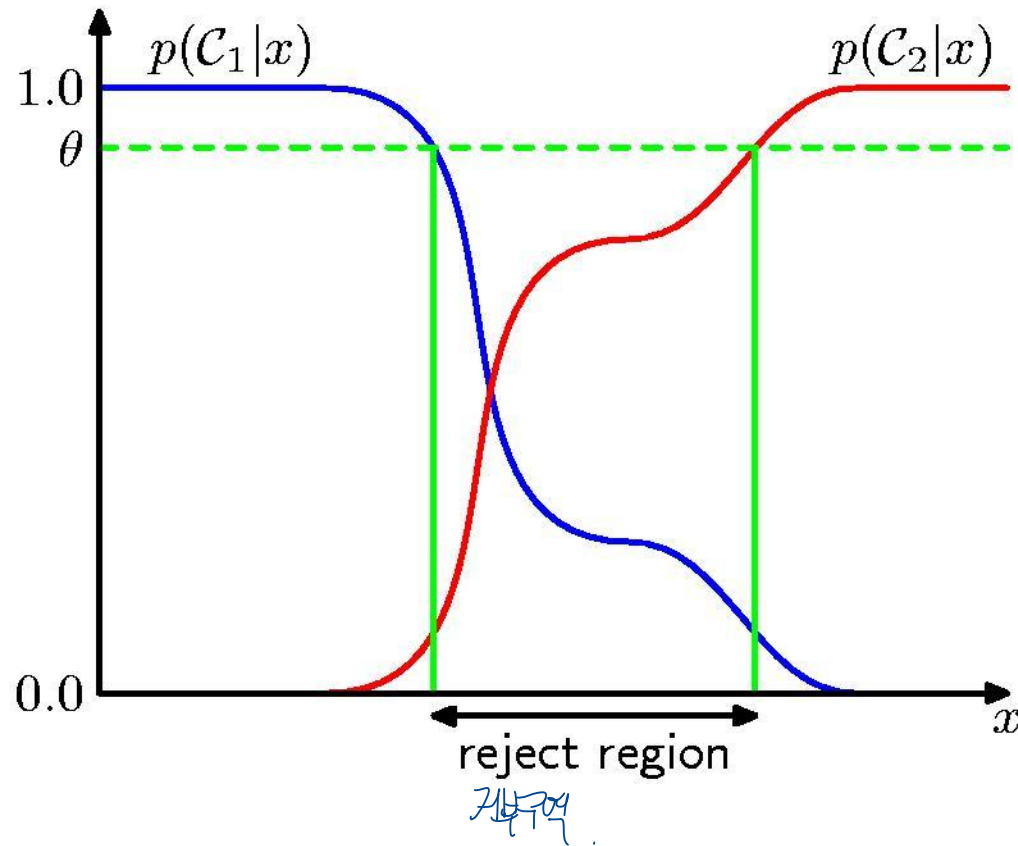
$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

Regions \mathcal{R}_j are chosen to minimize

$$\mathbb{E}[L] = \sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

$$p(\mathbf{x}, \mathcal{C}_k) = p(\mathcal{C}_k | \mathbf{x}) p(\mathbf{x})$$

Reject Option



Why Separate Inference and Decision?

- Minimizing risk (loss matrix may change over time)
- Reject option
- Unbalanced class priors
- Combining models

Decision Theory for Regression

Inference step

Determine $p(\mathbf{x}, t)$.

Decision step

For given \mathbf{x} , make optimal prediction, $y(\mathbf{x})$, for t .

Loss function: $\mathbb{E}[L] = \iint \underbrace{L(t, y(\mathbf{x}))}_{\text{loss}} p(\mathbf{x}, t) \, d\mathbf{x} \, dt$

Handwritten notes:
- $y(\mathbf{x})$ is the prediction
- $L(t, y(\mathbf{x}))$ is the loss

The Squared Loss Function

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

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

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) \, d\mathbf{x} + \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) \, d\mathbf{x}$$

$$y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$$

$$\text{Q. } \frac{\partial \mathbb{E}[L]}{\partial y(\mathbf{x})} = 2 \int \{y(\mathbf{x}) - t\} p(\mathbf{x}, t) \, dt = 0, \quad y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) \, dt}{p(\mathbf{x})} = \int t p(t|\mathbf{x}) \, dt = \mathbb{E}_t[t|\mathbf{x}]$$

Generative vs Discriminative

Generative approach:

Model $p(t, \mathbf{x}) = p(\mathbf{x}|t)p(t)$

Use Bayes' theorem $p(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)p(t)}{p(\mathbf{x})}$

Discriminative approach:

Model $p(t|\mathbf{x})$ directly

Entropy : 불확실성, 상태를 결정하는 정보의 양

$$h(x) = -\log_2 p(x)$$

$$h(x, y) = h(x) + h(y)$$

$$p(x, y) = p(x)p(y)$$

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

\nearrow
 $E[h(x)]$

Important quantity in

- coding theory
 - statistical physics
 - machine learning
-

Entropy

Coding theory: x discrete with 8 possible states; how many bits to transmit the state of x ?

All states equally likely

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

Entropy

x	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$
code	0	10	110	1110	111100	111101	111110	111111

$$\begin{aligned} H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\ &= 2 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{average code length} &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 \\ &= 2 \text{ bits} \end{aligned}$$

Entropy

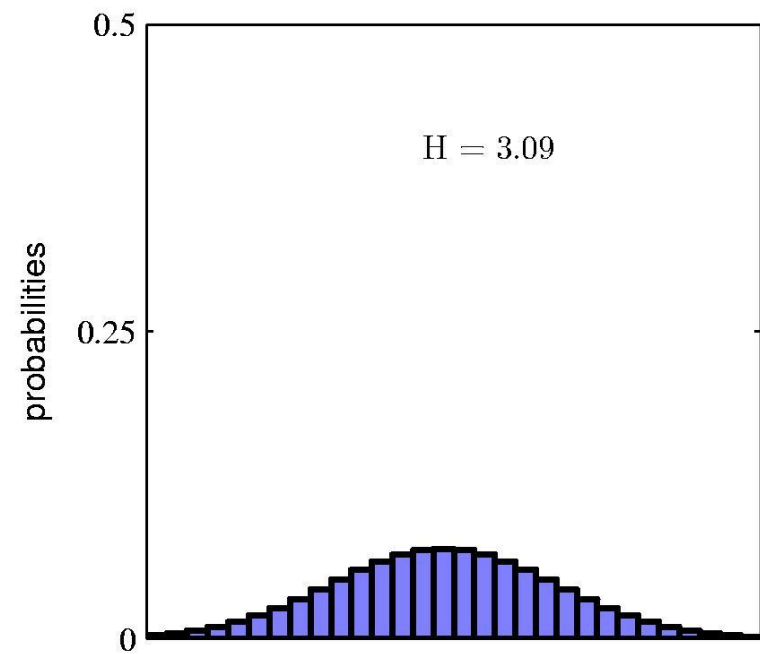
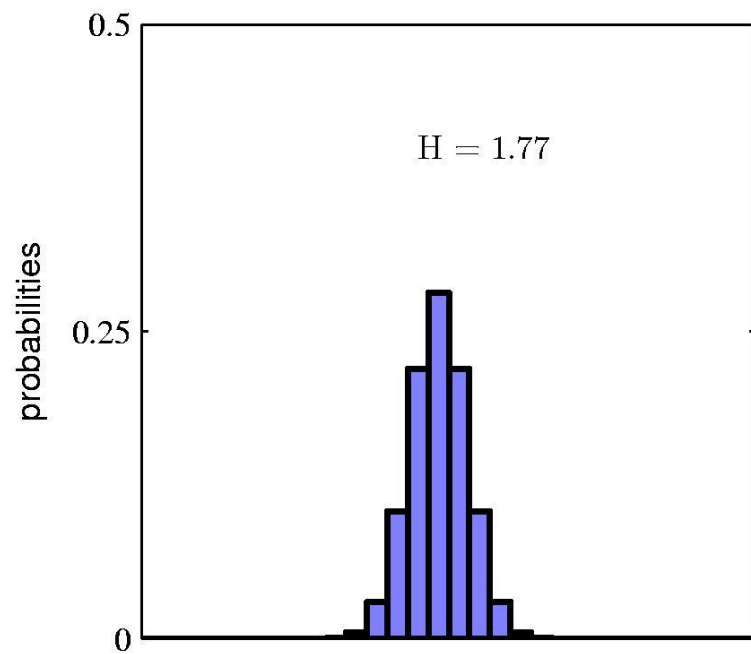
In how many ways can N identical objects be allocated M bins?

$$W = \frac{N!}{\prod_i n_i!}$$

$$H = \frac{1}{N} \ln W \simeq - \lim_{N \rightarrow \infty} \sum_i \left(\frac{n_i}{N} \right) \ln \left(\frac{n_i}{N} \right) = - \sum_i p_i \ln p_i$$

Entropy maximized when $\forall i : p_i = \frac{1}{M}$

Entropy



Differential Entropy

Put bins of width Δ along the real line

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

Differential entropy maximized (for fixed σ^2) when

$$p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

in which case

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

σ^2 가 클수록 엔트로피가 증가한다는 것을 다시 확인하였다.

Conditional Entropy

$$H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x}$$

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}]$$

-✓. \mathbf{x} 과 \mathbf{y} 를 독립짓기 위해 필요한 정보량은 \mathbf{x} 만 따로 독립짓기 위해 필요한 정보량의
가 주어졌을 때 \mathbf{y} 로 독립짓기 위해 필요한 정보량을 합친 것과 같다.

The Kullback-Leibler Divergence

$$\begin{aligned} \text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \end{aligned}$$

↑
한글자입은분포

↑
modeling한분포

$$\text{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^N \{ -\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n) \}$$

$$\text{KL}(p\|q) \geq 0$$

$$\text{KL}(p\|q) \neq \text{KL}(q\|p)$$

$$\text{KL}(p\|q)=0 \text{ 이면 } p(x)=q(x)$$

Mutual Information

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$
