# Chapter 2   Probability distributions

Introduction to useful probability distributions $<$ discrete / continuous

Discuss key statistical concepts such as Bayesian inference

Density estimation ( find $p(x)$ )          independent and identically

- iid assumption                    observations $x_1 ... x_N$

- unsupervised learning

- parametric vs non parametric

Conjugate distributions ( prior and posterior )

## 2.1 Discrete random variables

{ Bernoulli, binomial, beta } distribution

### Bernoulli distribution

R.v $x \in \{0, 1\}$, parameter $\mu$ denote the prob. of $x = 1$

$$p(x = 1 \mid \mu) = \mu \quad \text{— prob.} \qquad 0 \leq \mu \leq 1$$

parameter

$$\text{Bern}(x \mid \mu) = \mu^x (1 - \mu)^{1-x} \qquad (x \in \{0, 1\})$$

## Remark

- $p(x=0 \mid \mu) + p(x=1 \mid \mu) = (1-\mu) + \mu = 1$

- $\mathbb{E}[x] = 0 \cdot p(x=0 \mid \mu) + 1 \cdot p(x=1 \mid \mu) = \mu$

- $\mathrm{Var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \mu - \mu^2 = \mu(1-\mu)$

$$x_i \overset{D}{\sim} \mathrm{Bern}(x \mid \mu)$$

Assume $\{x_1, \dots x_N\}$ is drawn independently from some Bernoulli

Find a parameter $\mu$ $(= p(x=1))$ in frequentist setting

See likelihood function (of $\mu$) ($D = \{x_1, \ldots x_N\}$, $x_i \in \{0,1\}$)

$$P(D | \mu) = \prod_{n=1}^{N} P(x_n | \mu) = \prod_{n=1}^{N} \mu^{x_n} (1-\mu)^{(1-x_n)} \qquad (2.5)$$

Estimate (find) a parameter $\mu$ by maximizing (2.5)

$$\ln P(D | \mu) = \sum_{n=1}^{N} \{ x_n \ln \mu + (1-x_n) \ln (1-\mu) \}$$

$\ln P(D | \mu)$ only depends on $x_n$ ($\sum x_n$). Find a value for $\mu$ of

$$\frac{d}{d\mu} \ln P(D | \mu) = 0$$

$$\Rightarrow \quad \mu_{ML} = \frac{1}{N} \sum_{n=1}^{N} x_n = \frac{m}{N} \qquad \text{sample mean}$$

where $m = \#$ of $x = 1$

( likelihood 를 최대화 하는 $\mu$ )

The sample mean is an example of sufficient statistic.

E.g. Flip a coin 3 times, observe

$$\Rightarrow \quad \mu_{ML} = 1$$

# Binomial distribution  이항분포

Let $x = 0$ or $1$ and $N$ be ※ trials.

R.v $m \in \{0, 1, \ldots, N\}$ to be ※ of $x = 1$.

From (2.5)

$m = 0, 1, \ldots$ or $N$

binomial distribution $\propto \mu^m (1-\mu)^{N-m}$

$Ⓜ$

$m = $ ※ of $x = 1$, $\mu$ : the probability of $x = 1$

To normalize prob. dist. calculate all of possible ※
of obtaining m, x = 1. Denote by

$$Bin(m \mid N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$$

$$m = 0, 1, 2, \dots N$$

where $\binom{N}{m} := \dfrac{N!}{(N-m)!\, m!} = {}_N C_m$

We have

$$E[m] = \sum_{m=0}^{N} m \, Bin(m \mid N, \mu)$$

$$= \sum_{m=0}^{N} m \binom{N}{m} \mu^m (1-\mu)^{N-m} = N\mu$$

(평균 횟수)

$$\text{Var}[m] = \sum_{m=0}^{N} (m - \mathbb{E}[m])^2 \, \text{Bin}(m \mid N, \mu) = N\mu(1-\mu)$$

# 2.1.1 The beta distribution

Conjugate distribution

$p(\theta | x)$ and $p(\theta)$ have the same distributions.

$\boxed{p(\theta)}$ is called a conjugate prior for $p(x|\theta)$ likelihood

Goal
parameter $\theta$

※ $x$ 가 어떤 분포를 따를 지 가정

$\Downarrow$

posterior of $\theta$     likelihood function of $\theta$     prior of $\theta$

Recall the $p(D|\mu)$ of Bernoulli distribution

likelihood
$$p(D|\mu) = \prod_{n=1}^{N} \mu^{x_n} (1-\mu)^{(1-x_n)} \qquad x_n \in \{0, 1\}$$

$$p(x|\mu) = \mu^x (1-\mu)^{1-x}$$

To see Bayesian approach, we need to introduce $p(\mu)$.

beta distribution    normalization constant

(2.13)

$$\text{Beta}(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \qquad 0 \leq \mu \leq 1$$

where $\Gamma(x) := \int_0^\infty u^{x-1} e^{-u} \, du$ is the gamma function (1.414)

# Remark

- $$\int_0^1 \text{Beta}(\mu \mid a, b) \, d\mu = 1$$

- $$E[\mu] = \frac{a}{a+b}, \qquad \text{Var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)} \qquad \text{Excercise.}$$

- $a, b$ are called hyperparameters

The posterior dist. of $\mu$ has the form as

$$P(\mu \mid m, \ell, a, b) \propto \mu^{m+a-1} (1-\mu)^{\ell+b-1}$$

where $\ell = N - m$.  $\quad m:$ # of $x_n = 1$

$$\text{posterior} \propto P(D \mid \mu) \quad P(\mu)$$
$$\mu^m (1-\mu)^\ell \quad \mu^{a-1} (1-\mu)^{b-1}$$

posterior

$$\Rightarrow \quad P(\mu \mid m, \ell, a, b) = \frac{\Gamma(m+a+\ell+b)}{\Gamma(m+a)\Gamma(\ell+b)} \mu^{m+a-1} (1-\mu)^{\ell+b-1} \qquad (2.18)$$
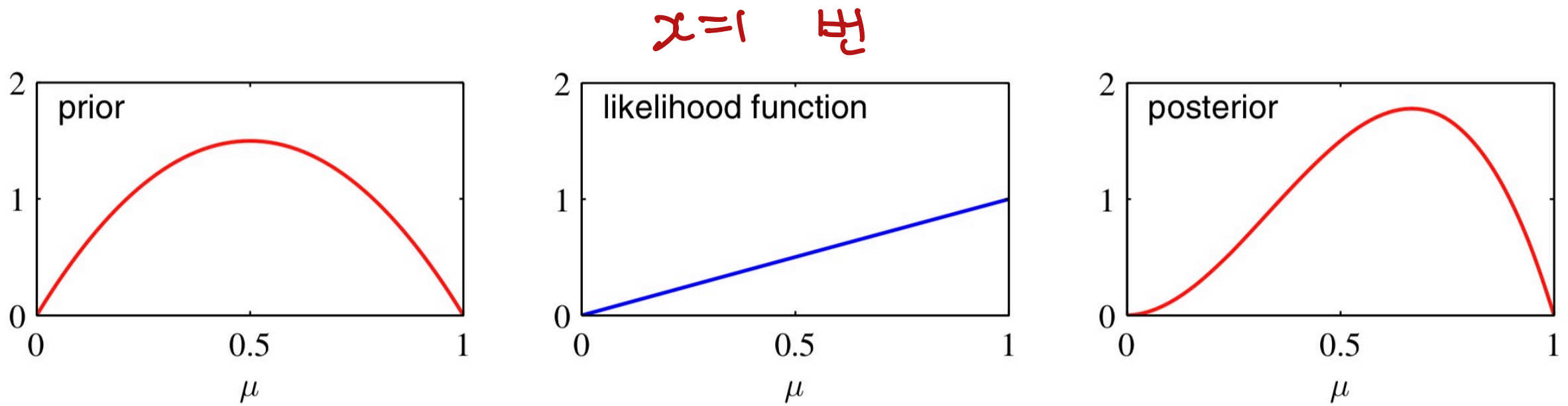
$a, b$ : parameters of prior

$m, \ell$ : result of observation
$\left. \right) \longrightarrow$ posterior

In view of (2.18) and def of beta dist,

$a$ , $b$ can be <u>interpreted</u> as effective # observations

of $x = 1$ and $x = 0$.

Sequential approach ( Bayesian view point)

$x=1$ 번

**Figure 2.3** Illustration of one step of sequential Bayesian inference. The prior is given by a beta distribution with parameters $a = 2$, $b = 2$, and the likelihood function, given by (2.9) with $N = m = 1$, corresponds to a single observation of $x = 1$, so that the posterior is given by a beta distribution with parameters $a = 3$, $b = 2$.

시행 횟수 1

Bern

$x_1 \dots x_N$

$\mu$

observation

Let us predict the outcome of the next trial

$$p(x=1 \mid D) = \int_0^1 p(x=1, \mu \mid D)\, d\mu = \int_0^1 \frac{p(x=1, \mu, D)}{p(D)}\, d\mu$$

$$= \int_0^1 \frac{p(x=1, \mu, D)}{p(\mu, D)} \cdot \frac{p(\mu, D)}{p(D)}\, d\mu$$

$$= p(x=1 \mid \mu, D) \cdot p(\mu \mid D)$$

conditionally independent

$$= \int_0^1 p(x=1 \mid \mu) \cdot p(\mu \mid D)\, d\mu$$

$A \perp\!\!\!\perp B \mid C$

$$\Leftrightarrow p(A \mid B, C)$$
$$= p(A \mid C)$$

$$= \int_0^1 \mu \, p(\mu \mid D)\, d\mu = \mathbb{E}_\mu[\mu \mid D]$$

posterior

$$= \frac{m + a}{m + a + \ell + b}$$

$$m + \ell = N$$

$$p(x=1 \mid D) = \frac{m+a}{m+a+l+b} = \frac{m+a}{N+a+b} \qquad (l = N - m)$$

$a, b$ are hyperparameters of prior.

$m, l$ are from the result of experiment ($m$ = # of $x=1$)

Remark

- $m, l \to \infty$ (huge observations), then ML $\approx$ Bayesian

- $a, b \to \infty$, then variance $\to 0$ (both prior and posterior)

## 2.2 Multinomial variables

Extend Bernoulli, binomial, beta distributions

Discrete variables that can take on one of K possible case

1-of-K scheme ( one hot encoding )

| ✗ | orange | apple | grape | |
|---|--------|-------|-------|---|
| orange | 1 | 0 | 0 | $(1,0,0)^T$ |
| apple | 0 | 1 | 0 | $(0,1,0)^T$ |
| apple | 0 | 1 | 0 | $(0,1,0)^T$ |
| grape | 0 | 0 | 1 | $(0,0,1)^T$ |
| orange | 1 | 0 | 0 | $(1,0,0)^T$ |
| | $M_1$ | $M_2$ | $M_3$ | |

$\Rightarrow$ $\Rightarrow$

$k$-dim. vector $x = (x_1, x_2, \ldots x_k)^T$ as $(0, 0, \ldots 1, 0, 0 \ldots 0)^T$

$$\sum_{k=1}^{K} x_k = 1$$

경우의 수 $= K$

$x \in \{0, 1\}$

Let $\mu_k$ be $P(x_k = 1)$.

Bernoulli

The distribution of $x$ ( Categorical distribution )

$$P(x \mid \mu) = \prod_{k=1}^{K} \mu_k^{x_k}$$

where $\mu := (\mu_1, \ldots \mu_k)^T$ with $\mu_k \geq 0$, $\sum_{k=1}^{K} \mu_k = 1$

# Remark

- $x$ can take $K$ possible cases.

- $\sum_x p(x \mid \mu) = \sum_{k=1}^{K} \mu_k = 1$

- $\mathbb{E}[x \mid \mu] = \sum_x p(x \mid \mu) x = (\mu_1, \ldots \mu_K)^T = \mu$
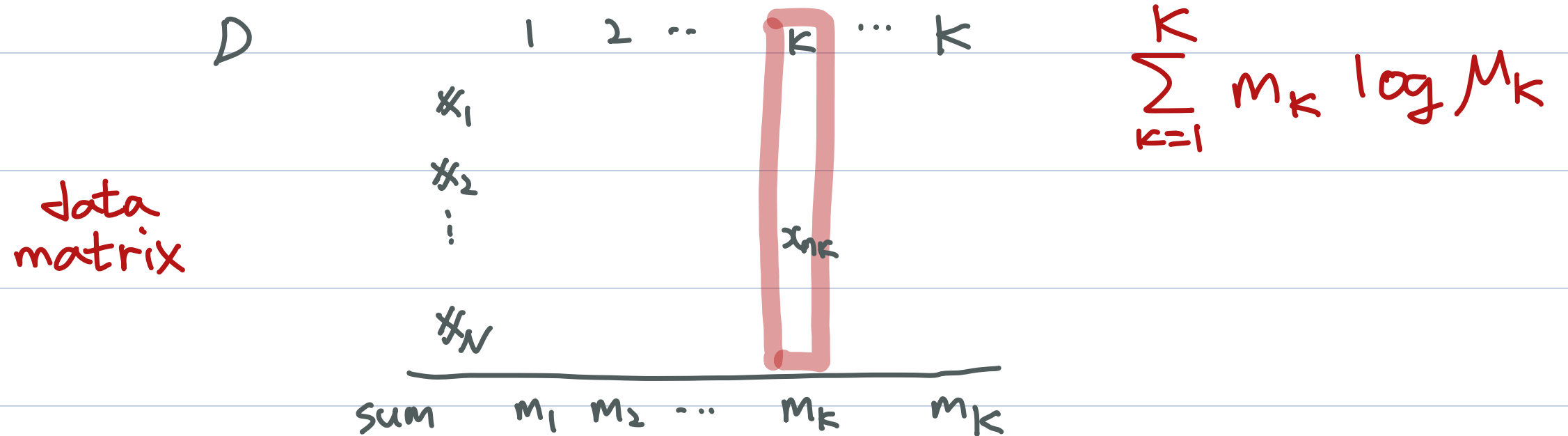
Consider $D$ of $N$ independent observations $x_1, \ldots x_N$.

Likelihood

$$P(D \mid \mu) = \prod_{n=1}^{N} \prod_{k=1}^{K} \mu_k^{x_{nk}} = \prod_{k=1}^{K} \mu_k^{\left(\sum_n x_{nk}\right)} = \prod_{k=1}^{K} \mu_k^{m_k}$$

where $\quad m_k = \sum_n x_{nk} \quad$ ( # of observations $x_k = 1$)



D      1   2   ..    k   ...   K

$x_1$

$x_2$

⋮

$x_N$

data matrix

$x_{nk}$

$$\sum_{k=1}^{K} m_k \log \mu_k$$

sum    $m_1$   $m_2$   ...   $m_k$    $m_K$

Maximize log likelihood $\ln p(D|\mu)$ for $\mu$ with constraint

$\sum \mu_k = 1$. Using Lagrange multiplier and maximizing

$$\sum_{k=1}^{K} m_k \ln \mu_k + \lambda \left( \sum_{k=1}^{K} \mu_k - 1 \right) \qquad (2.31)$$

$$\text{constraint}$$

Set the derivative of (2.31) wrt $\mu_k$ to zero

$$\Rightarrow \qquad \mu_k = \frac{m_k}{N}$$

$$\mu_k^{ML} = \frac{m_k}{N}$$

$$\mu_{ML} = (\mu_1^{ML}, \dots \mu_K^{ML})^T$$

Multinomial distribution $\quad m := (m_1, \ldots m_k)^T$

Joint distribution of the quantities $m_1, \ldots m_k$ conditioned

on $\mu$ and on ※ of $N$ total observations.

$$\text{Mult}(m_1, m_2, \ldots m_k \mid \mu, N) = \binom{N}{m_1\, m_2 \cdots m_k} \prod_{k=1}^{k} \mu_k^{m_k}$$

$$0 \le \mu_k \le 1$$

$$\text{where} \binom{N}{m_1\, m_2 \cdots m_k} := \frac{N!}{m_1!\, m_2! \cdots m_k!} \qquad \sum_{k=1}^{K} m_k = N$$

$k$ 개 categorical 변수를 갖는 $N$ 개 자료에서 각 $k$-class 가

$m_k$ 씩 가질 확률 ( $\mu$ 가 주어졌을 때)

## 2.2.1 Dirichlet distribution ( multi-dim version of beta)

Consider the prior distributions for the parameters $\{\mu_k\}$ of multinomial distribution. (or <span style="color:red">categorical distribution)</span>

Recall

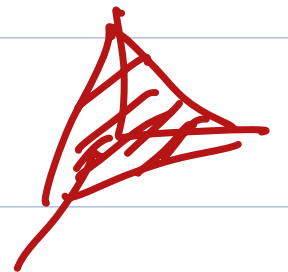$$\text{Mult} (m_1, \ldots m_K \mid \mu, N) \propto \prod_{k=1}^{K} \mu_k^{m_k}$$

$$\Rightarrow \quad p(\mu \mid \alpha) \propto \prod_{k=1}^{K} \mu_k^{\alpha_k - 1} \qquad \underset{\mu_1}{\mu}^{a-1} \; \underset{\mu_2}{(1-\mu)}^{b-1}$$

where $0 \le \mu_k \le 1$, $\sum_{k} \mu_k = 1$. Here $\alpha := (\alpha_1, \ldots \alpha_K)^T$ is

the parameter

$\{\mu_k\}$ is confined to a simplex of dim $K-1$

# Dirichlet distribution

$$\text{Dir}(\mu \mid \alpha) := \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^{K} \mu_k^{\alpha_k - 1} \qquad 0 \leq \mu_k \leq 1$$

where $\Gamma(x)$ is the Gamma function, $\alpha_0 := \sum \alpha_k$

So the posterior distribution for the parameters $\{\mu_k\}$

$$\underset{\text{posterior}}{p(\mu \mid D, \alpha)} \propto \underset{\text{likelihood}}{p(D \mid \mu)} \underset{\text{prior}}{p(\mu \mid \alpha)} \propto \prod_{k=1}^{K} \mu_k^{\alpha_k + m_k - 1}$$

Posterior dist. again follow Dirichlet dist.

$$\prod \mu_k^{m_k} \times \prod \mu_k^{\alpha_k - 1}$$

$$\Rightarrow \quad P(\mu \mid D, \alpha) = Dir(\mu \mid \underbrace{\alpha}_{\text{Prior}} + \underbrace{m}_{D})$$

$$= \frac{\Gamma(\alpha_k + m_k)}{\Gamma(\alpha_1 + m_1) \cdots \Gamma(\alpha_k + m_k)} \prod_{k=1}^{K} \mu_k^{\overset{\text{Prior}}{\alpha_k} + \overset{\text{observation}}{m_k} - 1}$$

where $\quad m := (m_1, \dots m_K)^T$.

As binomial dist. with beta prior, we can interprete

$\alpha_k$ of Dirichlet prior as an effective $\#$ of $x_k = 1$.

## 2.3 The Gaussian distribution (a.k.a. normal dist.)

Single real value $x \in \mathbb{R}$

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

where $\mu$: mean, $\sigma^2$: variance

$D-\dim$ vector $x \in \mathbb{R}^D$

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

where $\mu$: $D-\dim$ mean vector, $\Sigma$: $D \times D$ covariance matrix

$|\Sigma|$ : determinant of $\Sigma$

# Remark: Why Gaussian is important?

- Fits many natural phenomena

- Maximum entropy in continuous r.v

- Central limit theorem

Fix $N$. iid random samples of vector $X_1, X_2 \ldots X_N$ are drawn from a population with $M$, $\Sigma$.

$\Rightarrow$ R.V. $\overline{X}_N := \frac{1}{N} \sum_{n=1}^{N} X_n \approx N(M, \frac{1}{N}\Sigma)$

sample mean



**Figure 2.6** Histogram plots of the mean of $N$ uniformly distributed numbers for various values of $N$. We observe that as $N$ increases, the distribution tends towards a Gaussian.

Let us see geometrical form of Gaussian.

$$\Delta^2 := (x - \mu)^T \Sigma^{-1} (x - \mu) \qquad \mathcal{N}(x | \mu, \Sigma)$$

<span style="color:red">1×D      D×D      D×1</span>

$\Delta$ is called Mahalanobis distance from $\mu$ to $x$

WLOG. assume $\Sigma$ is symmetric (and real).

<span style="color:red">$\Delta$ = constant 인</span>    <span style="color:red">(등고선)</span>
<span style="color:red">$x$ 의 모임</span>

Consider the eigenvector equation of $\Sigma$

$$\Sigma u_i = \lambda_i u_i \qquad i = 1, \ldots D$$

$$(\lambda_i, u_i)$$

Eigenvalues $\lambda_1 \ldots \lambda_D$ are real and its eigenvectors can be chosen to form an <u>orthonormal set</u>, so that

$$u_i^T u_j = I_{ij}$$

where $I_{ij}$ is the $i,j$ element of the identity matrix.
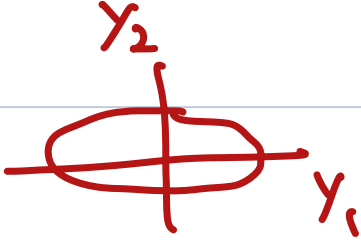
By eigenvector decomposition, $\Sigma$ can be expressed as

$$\Sigma = \sum_{i=1}^{D} \overset{\text{value}}{\lambda_i} \, u_i \, u_i^T \qquad D\times 1 \quad 1\times D$$

$D\times D$        vector

and

$$\Sigma^{-1} = \sum_{i=1}^{D} \frac{1}{\lambda_i} u_i \, u_i^T \qquad \text{(inverse)}$$

$$\Rightarrow \quad \Delta^2 = (\textbf{x} - \mu I)^T \Sigma^{-1} (\textbf{x} - \mu I) = (\textbf{x} - \mu I)^T \left( \sum_{\lambda=1}^{D} \frac{1}{\eta_\lambda} u_\lambda u_\lambda^T \right) (\textbf{x} - \mu I)$$

Mahalanobis
distance



$$= \sum_{\lambda=1}^{D} \frac{y_\lambda^2}{\eta_\lambda} = C \quad \text{의}$$

$$\textbf{x} \text{의 모임}$$

where $\quad y_\lambda := u_\lambda^T (\textbf{x} - \mu I) \quad$ ( inner product of $u_\lambda$ , $(\textbf{x} - \mu I)$ )

$$1 \times D \qquad D \times 1 \qquad \qquad \textbf{x}$$
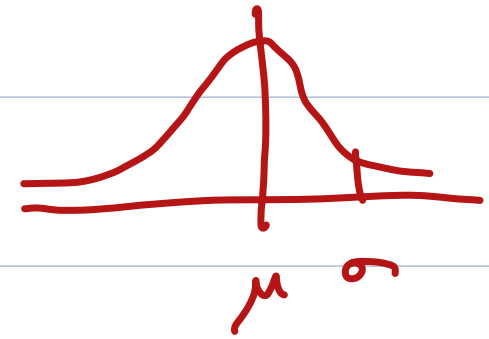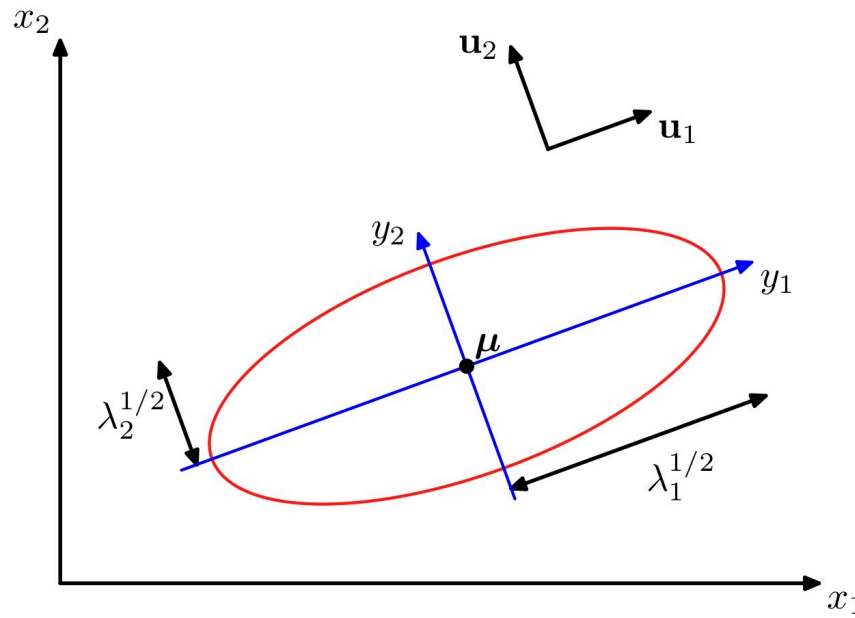
$y_\lambda$ : new coordinate system defined by $u_\lambda$, shifted and rotated

Let $\quad y := (y_1, \dots, y_D)^T$. Then $\quad y = U (\textbf{x} - \mu I) \quad$ where

$$U = \begin{pmatrix} u_1^T \\ \vdots \\ u_D^T \end{pmatrix} \quad \text{whose rows are} \quad u_\lambda^T$$

( orthogonal matrix )

**Figure 2.7** The red curve shows the elliptical surface of constant probability density for a Gaussian in a two-dimensional space $\mathbf{x} = (x_1, x_2)$ on which the density is $\exp(-1/2)$ of its value at $\mathbf{x} = \boldsymbol{\mu}$. The major axes of the ellipse are defined by the eigenvectors $\mathbf{u}_i$ of the covariance matrix, with corresponding eigenvalues $\lambda_i$.



## Remark

— If $\lambda_i > 0 \;^{\forall} i = 1 \ldots D$, contour surface of $\triangle$ is ellipsoid.

— center $\mu$, axes oriented along $u_i$ and scailing factors are given by $\lambda_i^{\frac{1}{2}}$

WLOG, assume all eigenvalues of $\Sigma$ are strictly positive
Otherwise the distribution cannot be normalized (see ch 12)
i.e. $\Sigma$ is assume to be positive definite.

$$y := U(x - \mu)$$

Now consider the Gaussian dist. in $y_i$ coordinate system.

Jacobian matrix $J$ with

(1.27) $\underline{x = g(y)}$

$$P_y(y) = P_x(g(y))|g'(y)|$$

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji}$$

$$U^T y + \mu = x$$

where $U_{ji}$ are element of $U^T$

So $\quad |J|^2 = |U^T|^2 = |U^T||U| = |U^T U| = |I| = 1$

and hence $|J| = \pm 1$. Also $|\Sigma|$ can be written as

$$|\Sigma| = \prod_{j=1}^{D} \lambda_j \qquad \qquad \mathbf{x} \qquad \qquad y = U(\mathbf{x} - \mu)$$

Thus in the $y_j$ coordinate system,

$$P(y) = P(\mathbf{x}) \underset{\pm 1}{|\det(J)|} = \prod_{j=1}^{D} \frac{1}{(2\pi \lambda_j)^{1/2}} \exp\left\{-\frac{y_j^2}{2\lambda_j}\right\}$$

$$\mathcal{N}(y \mid 0, \text{diag}(\sqrt{\lambda_1}...\sqrt{\lambda_D})) \qquad = \prod_{j=1}^{D} P_{y_i}(y_i)$$

$P(y)$ is the product of $D$ independent univariate Gaussian.

In (1.49), (1.51), we found univariate Gaussian dist has

$$\mathbb{E}[x] = \mu, \quad \text{var}[x] = \sigma^2$$

Now we will intepret parameters $\mu$ (D-dim), $\Sigma$ (D×D).

$$\mathbb{E}[x] = \frac{1}{(2\pi)^{\frac{1}{2}D}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\} x \, dx$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}D}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2} z^T \Sigma^{-1} z\right\} (z+\mu) \, dz$$

where we have changed variables using $z = x - \mu$.

Note that the exponent is even. So the term $z$

in the factor $(z + \mu)$ will vanish.

$$\Rightarrow \qquad \mathbb{E}[x] = \mu.$$

$x \in \mathbb{R}^D$

Now consider the second order moments of multivariate

Gaussian. In univariate case, the second order moment

is given by $\mathbb{E}[x^2]$. In multivariate Gaussian, there are

$D^2$ second order moments given by $\mathbb{E}[x_i x_j]$.

$$\mathbb{E}[x_i x_j] \qquad\qquad x = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix}$$

$$\mathbb{E}[\pmb{x}\,\pmb{x}^T] = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left\{-\frac{1}{2}(\pmb{x}-\pmb{\mu})^T \Sigma^{-1}(\pmb{x}-\pmb{\mu})\right\} \pmb{x}\,\pmb{x}^T \, d\pmb{x}$$

$\underset{D\times 1}{} \underset{1\times D}{}$

$$= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left\{-\frac{1}{2}\pmb{z}^T \Sigma^{-1}\pmb{z}\right\} \underbrace{(\pmb{z}+\pmb{\mu})(\pmb{z}+\pmb{\mu})^T}_{(*)} \, d\pmb{z}$$

$\pmb{z}\pmb{z}^T$

$$\pmb{z} := \pmb{x} - \pmb{\mu}.$$

$(*) \qquad (\pmb{z}+\pmb{\mu})(\pmb{z}^T+\pmb{\mu}^T) = \pmb{z}\pmb{z}^T + \overbrace{\pmb{z}\pmb{\mu}^T + \pmb{\mu}\pmb{z}^T}^{\text{vanish by symmetry}} + \pmb{\mu}\pmb{\mu}^T$

constant

transpose of

Recall $\quad y := U(\pmb{x}-\pmb{\mu}),\quad$ rows of $U$ are eigenvectors of $\Sigma$

$\underset{D\times D}{} \quad \underset{1\times 1}{}$

$$\Rightarrow \quad \pmb{z} = U^T y = (u_1, \ldots u_{1_D})\begin{pmatrix} Y_1 \\ \vdots \\ Y_D \end{pmatrix} = \sum_{j=1}^{D} Y_j \, u_{1_j}$$

$\underset{D\times D}{} \qquad \underset{D\times 1}{} \qquad\qquad \underset{D\times 1}{}$

$$\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left\{-\frac{1}{2} \not{z}^T \Sigma^{-1} \not{z}\right\} \not{z} \not{z}^T d\not{z}$$

$$= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \sum_{i=1}^{D} \sum_{j=1}^{D} u_i u_j^T \int \exp\left\{-\sum_{k=1}^{D} \frac{y_k^2}{2\lambda_k}\right\} y_i y_j \, dy$$

$$= \sum_{i=1}^{D} u_i u_i^T \lambda_i = \Sigma$$

we have used $|\Sigma| = \prod_{i=1}^{D} \lambda_i$ and $\frac{1}{(2\pi)^{1/2}} \frac{1}{\lambda^{1/2}} \exp\left\{-\frac{y^2}{2\lambda}\right\} \sim N(0,\lambda)$

e.g. <span style="color:red">D=2</span>

$$\sum_{i} \sum_{j} \iint \exp\left(-\frac{y_1^2}{2\lambda_1}\right) \exp\left(-\frac{y_2^2}{2\lambda_2}\right) y_i y_j \, dy_1 \, dy_2$$

will vanish when $i \neq j$

Thus we obtain $\mathbb{E}[xx^T] = \mu_1\mu_1^T + \Sigma$ (DxD matrix)

and covariance of $x$ can be obtained by

$$\text{cov}[x] = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^T]$$

DxD

$$= \Sigma$$

# Remark

$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \Sigma)$   sym. real
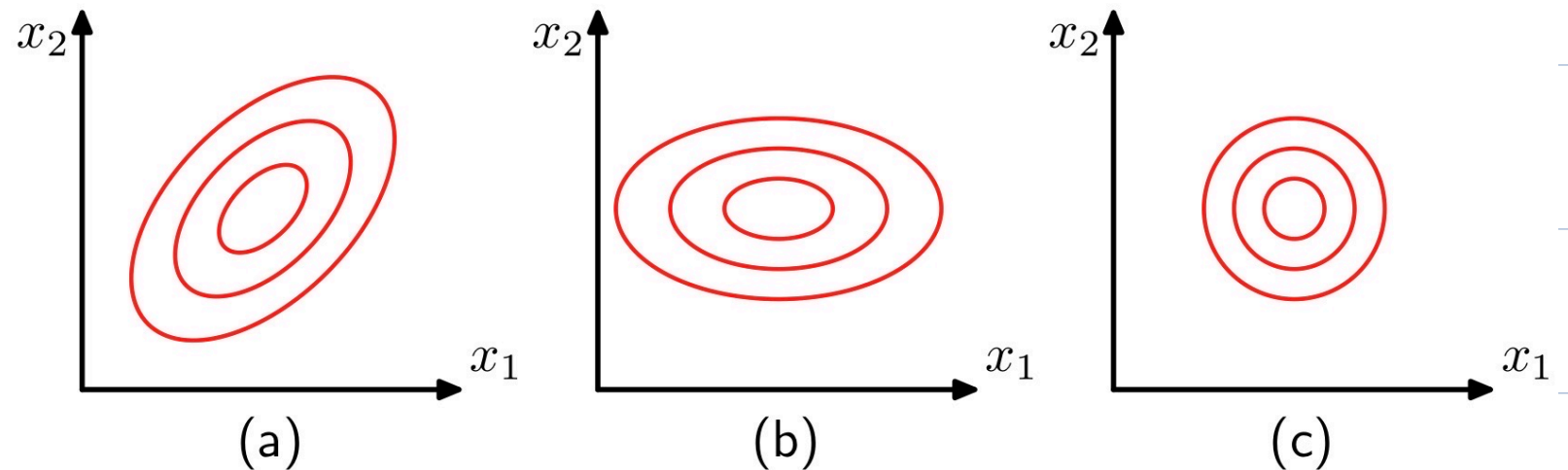
$\boldsymbol{\mu}$: $D$   $\Sigma$: $D \times D$

- \# of parameters : $\dfrac{D(D+3)}{2}$   quadratic

- $\Sigma = \text{diag}(\sigma_i^2)$   or   $\Sigma = \sigma^2 I$

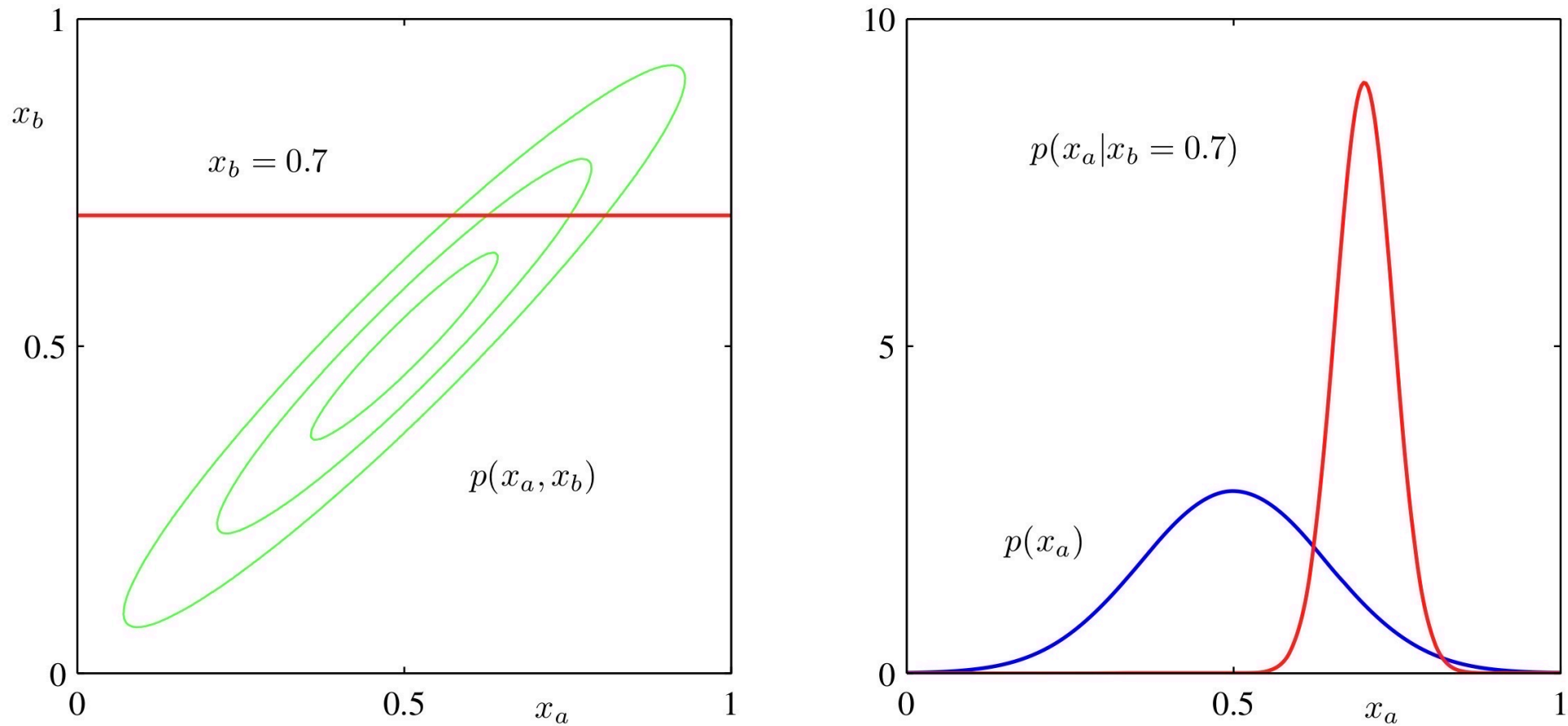  $2D$                         $D+1$

  deep dive
  into Gaussian

- Unimodal (single maximum)

**Figure 2.8** Contours of constant probability density for a Gaussian distribution in two dimensions in which the covariance matrix is (a) of general form, (b) diagonal, in which the elliptical contours are aligned with the coordinate axes, and (c) proportional to the identity matrix, in which the contours are concentric circles.

**Figure 2.9** The plot on the left shows the contours of a Gaussian distribution $p(x_a, x_b)$ over two variables, and the plot on the right shows the marginal distribution $p(x_a)$ (blue curve) and the conditional distribution $p(x_a|x_b)$ for $x_b = 0.7$ (red curve).

## 2.3.1 Conditional Gaussian distributions

$x$ : $D$-dimensional vector with $N(x \mid \mu, \Sigma)$ which is partitioned into $x_a$, $x_b$ with $\underline{x_a \in \mathbb{R}^M}$, $\underline{x_b \in \mathbb{R}^{D-M}}$

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$

corresponding partitions of mean vector, covariance matrix

$$\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \qquad \Sigma = \begin{pmatrix} \overset{M \times M}{\Sigma_{aa}} & \overset{M \times D-M}{\Sigma_{ab}} \\ \underset{D-M \times M}{\Sigma_{ba}} & \underset{D-M \times D-M}{\Sigma_{bb}} \end{pmatrix}$$

Note that $\Sigma^T = \Sigma$ implies $\Sigma_{aa}$, $\Sigma_{bb}$ are symmetric, $\Sigma_{ba} = \Sigma_{ab}^T$

Let $\Lambda := \Sigma^{-1}$. Inverse of covariance matrix, precision matrix

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$M \times M$

$D \times D$

$D-M \times D-M$

Note that $\Lambda_{aa}$, $\Lambda_{bb}$ are symmetric, $\Lambda_{ba} = \Lambda_{ab}^{T}$.

$\Lambda_{aa} \neq \Sigma_{aa}^{-1}$          excercise

Find conditional distribution $p(x_a | x_b)$. Fix $x_b$.

Consider the quadratic form in the exponent.

$$-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu) \qquad \underline{1 \times D} \quad \underline{D \times D} \quad \underline{D \times 1} \qquad (x_b \text{ fixed})$$

$$x_a, x_b$$

$$= -\frac{1}{2}(x_a - \mu_a)^T \Lambda_{aa} (x_a - \mu_a) - \frac{1}{2}(x_a - \mu_a)^T \Lambda_{ab} (x_b - \mu_b)$$

$$(2.70)$$

$$-\frac{1}{2}(x_b - \mu_b)^T \Lambda_{ba} (x_a - \mu_a) - \frac{1}{2}(x_b - \mu_b)^T \Lambda_{bb} (x_b - \mu_b)$$

First, $p(x_a | x_b)$ will be M-dim Gaussian, because density

function is a quadratic form of exponent.

Now we are going to find its mean vector and covariance

by "completeting the square".

$$(M \mid D-M) \overset{\frown}{\underset{\smile}{}}$$

$$(M \mid D-M) \begin{pmatrix} \bigcirc & \bigcirc \\ \bigcirc & \bigcirc \end{pmatrix} \begin{pmatrix} M \\ \overline{D-M} \end{pmatrix}$$

$$\underline{1 \times D} \qquad \underline{D \times D} \quad \underline{D \times 1}$$

E.g. in case $p(x) \propto \exp( ax^2 + bx + c)$ → $x \overset{i.i.d}{\sim}$ gaussian

① ⟹

$$p(x) \propto \exp\left\{ a\left(x^2 + \frac{b}{a}x + \frac{b^2}{4a^2}\right) - \frac{b^2}{4a} + c\right\}$$

$$\propto \exp\left\{ a\left(x + \frac{b}{2a}\right)^2\right\}$$

② $\quad p(x) = N(x \mid \mu, \sigma^2) \quad \mu = -\frac{b}{2a}, \quad \sigma^2 = -\frac{1}{2a}$

Since $\quad N(x \mid \mu, \sigma^2) \propto \exp\left\{ -\frac{1}{2\sigma^2}\left(x^2 - 2\mu x + \mu^2\right)\right\}$

$$\propto \exp\left\{ -\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x\right\}$$

$$\Rightarrow \quad a = -\frac{1}{2\sigma^2}, \quad b = \frac{\mu}{\sigma^2} \quad \left( \mu = -\frac{b}{2a}, \quad \sigma^2 = -\frac{1}{2a}\right)$$

Likewise, the exponent in D-dim Gaussian can be written

$$-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) = \boxed{-\frac{1}{2}} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu + \text{constant}$$

In view of (2.70), (second order in $x_a$)

$$-\frac{1}{2} x_a^T \Lambda_{aa} x_a$$

So we obtain covariance of $p(x_a | x_b)$ is given by

$$\Sigma_{a|b} := \Lambda_{aa}^{-1} \ (\neq \Sigma_{aa})$$

$$\mu_{a|b} \ ?$$

Now consider all of the terms in (2.70) that are linear in $\varkappa_a$

$$\varkappa_a^T \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (\varkappa_b - \mu_b) \}$$

where we have used $\Lambda_{ba}^T = \Lambda_{ab}$.

Since

$$\Sigma_{a|b}^{-1} \mu_{a|b} = \Lambda_{aa} \mu_a - \Lambda_{ab} (\varkappa_b - \mu_b),$$

$$\mu_{a|b} = \overset{\overset{\textcolor{red}{\Lambda_{aa}^{-1}}}{\textcolor{red}{\parallel}}}{\Sigma_{a|b}} \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (\varkappa_b - \mu_b) \}$$

$$= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\varkappa_b - \mu_b)$$

$$\therefore \quad \mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b), \qquad \Sigma_{a|b} = \Lambda_{aa}^{-1}$$

Let us find $\Lambda_{aa}$ and $\Lambda_{ab}$.

Recall $\Lambda = \Sigma^{-1}$,
$$\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

Use the following identification for the inverse of a partitioned matrix
$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CM BD^{-1} \end{pmatrix} \qquad (2.76)$$

where $M := (A - BD^{-1}C)^{-1}$

So we have

$$\Lambda_{aa} = \left(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}\right)^{-1}$$

$$\Lambda_{ba} = -\left(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}\right)^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}$$

and hence

$$\left(\begin{array}{l} \mu_{a|b} = \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b) \qquad\qquad \color{red}{x_b \ \text{fix}} \\[2em] \Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba} \end{array}\right.$$

Remark

– $\mu_{a|b}$ is a linear function of $x_b$

– $\Sigma_{a|b}$ is independent of $x_b$

# 2.3.2 Marginal Gaussian distributions

Consider the following marginal distribution

$$p(x_a) = \int p(x_a, x_b) \, dx_b$$

$$\mathcal{N}(x \mid \mu, \Sigma)$$

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \begin{matrix} M \\ D-M \end{matrix}$$

($x_b$ 로 주변화하고 남은 $x_a$는 어떤 분포인지)

Strategy : focus on quadratic form of exponent and identify

the mean vector and covariance matrix of $p(x_a)$

Recall (2.70)

$$\Sigma^{-1} = \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)$$

$$= -\frac{1}{2}(x_a - \mu_a)^T \Lambda_{aa}(x_a - \mu_a) - \frac{1}{2}(x_a - \mu_a)^T \Lambda_{ab}(x_b - \mu_b)$$

$$- \frac{1}{2}(x_b - \mu_b)^T \Lambda_{ba}(x_a - \mu_a) - \frac{1}{2}(x_b - \mu_b)^T \Lambda_{bb}(x_b - \mu_b)$$

In order to integrate out $x_b$, pick out those terms involving $x_b$

$$-\frac{1}{2} x_b^T \Lambda_{bb} x_b + x_b^T m_1 = -\frac{1}{2}(x_b - \Lambda_{bb}^{-1} m_1)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m_1)$$

(2.84)

$$+ \frac{1}{2} m_1^T \Lambda_{bb}^{-1} m_1 \qquad \text{(square expression)}$$

$$\underbrace{\phantom{+ \frac{1}{2} m_1^T \Lambda_{bb}^{-1} m_1}}_{\text{indep. of } x_b}$$

where $\quad m_1 := \Lambda_{bb} \mu_b - \Lambda_{ba}(x_a - \mu_a)$

For $\quad p(x_a) = \int p(x_a, x_b) \, dx_b$ ,

(2.86) $\quad \int \exp\left\{ -\frac{1}{2}(x_b - \Lambda_{bb}^{-1} m_1)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m_1) \right\} dx_b$

which is an inverse of the normalization coefficient.

As seen before, this coefficient is independent of mean.

Combining the last term $\left(\frac{1}{2} m^T \Lambda_{bb}^{-1} m\right)$ in (2.84) with remaining terms in (2.70) depending on $x_a$, we obtain

$$-\frac{1}{2} m^T \Lambda_{bb}^{-1} m - \frac{1}{2} x_a^T \Lambda_{aa} x_a + x_a^T (\Lambda_{aa} \mu_a + \Lambda_{ab} \mu_b) + \text{constant}$$

$$= -\frac{1}{2} \left[ \Lambda_{bb} \mu_b - \Lambda_{ba}(x_a - \mu_a) \right]^T \Lambda_{bb}^{-1} \left[ \Lambda_{bb} \mu_b - \Lambda_{ba}(x_a - \mu_a) \right]$$

$$+ x_a^T (\Lambda_{aa} \mu_a + \Lambda_{ab} \mu_b) + \text{constant}$$

$$= -\frac{1}{2} x_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) x_a + x_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mu_a$$

$$+ \text{constant}$$

Recall the exponent in D-dim Gaussian can be written

$$-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) = -\frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu + \text{constant}$$

Denote the covariance of $p(x_a)$ by $\Sigma_a$ and $\Sigma_a$ is

given by

$$\Sigma_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1}$$

$$\underbrace{\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}}_{(\Sigma_a)^{-1}}$$

Similarly, mean vector is given by

$$\Sigma_a \underbrace{(\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})}_{\Sigma_a^{-1}} \mu_a = \mu_a$$

To simplify
$$\Sigma_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1},$$

recall
$$\begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

and use (2.76) expression of the inverse of a partitioned

matrix
$$\Sigma_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1} = \Sigma_{aa}$$

Thus we have
$$\mathbb{E}[x_a] = \mu_a, \quad \mathrm{cov}[x_a] = \Sigma_{aa}$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} \end{pmatrix}$$

where
$$p(x_a) = \int p(x_a, x_b)\, dx_b$$

$$\mathcal{N}(x \mid \mu, \Sigma) \quad \text{with} \quad \Lambda := \Sigma^{-1} \qquad \text{\textcolor{red}{D-dim}} \qquad \textcolor{red}{x}$$

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$
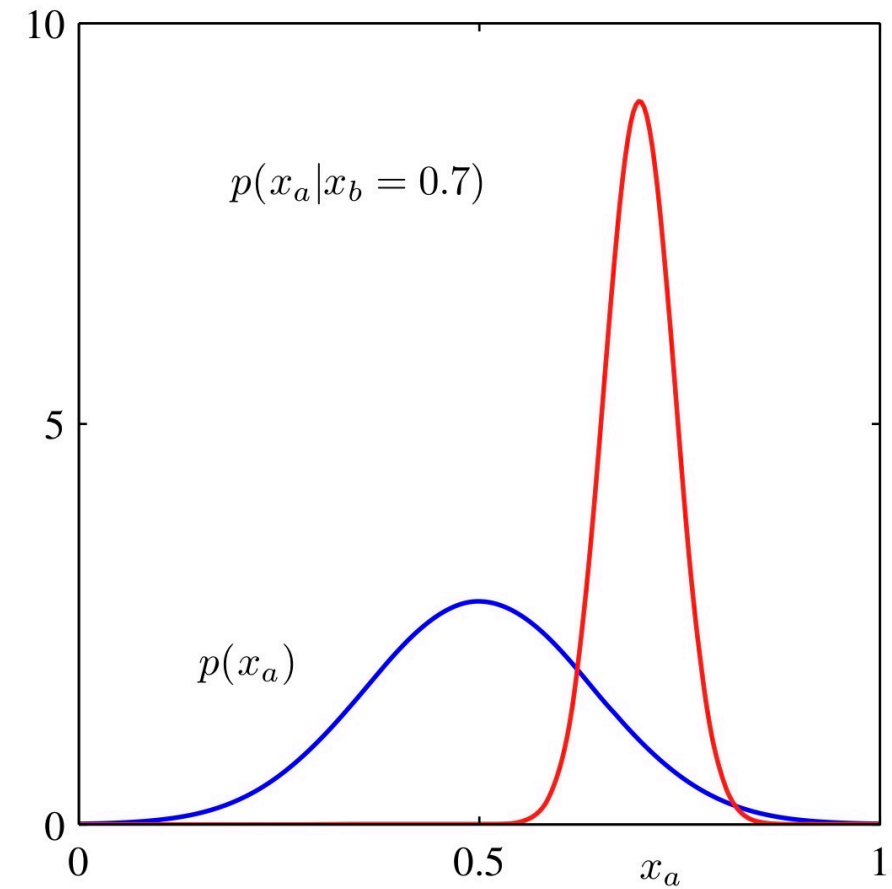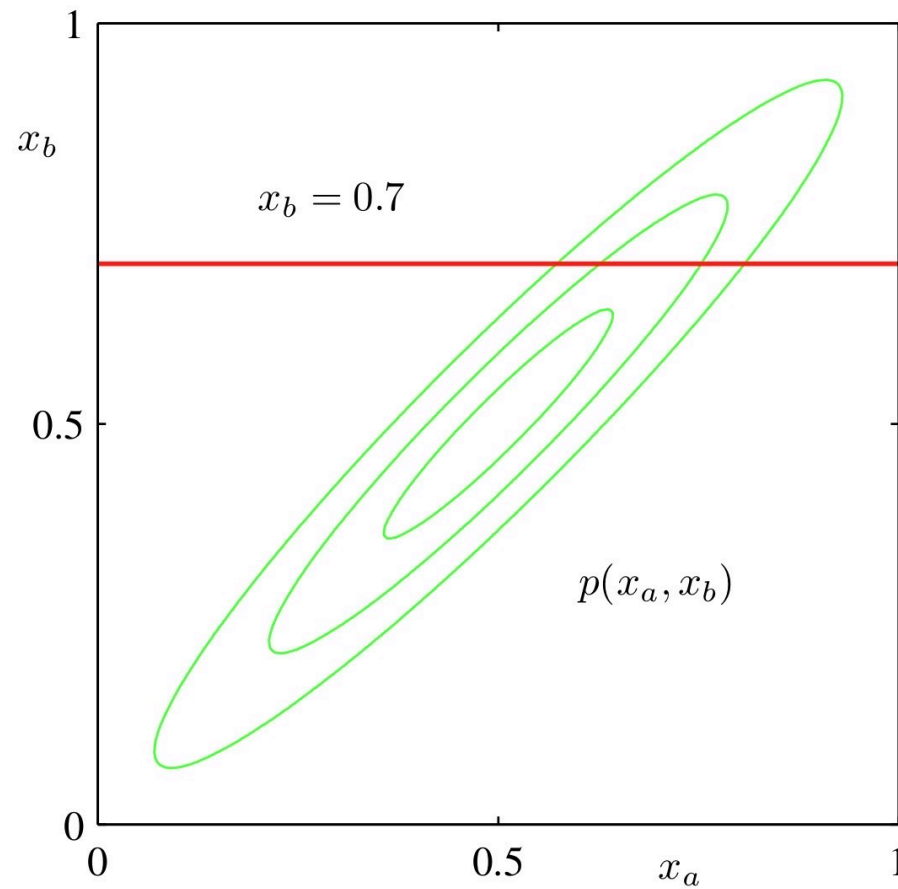
## Conditioned distribution

$$p(x_a \mid x_b) = \mathcal{N}(x_a \mid \mu_{a|b}, \Lambda_{aa}^{-1})$$

$$\mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b)$$

## Marginal distribution

$$p(x_a) = \mathcal{N}(x_a \mid \mu_a, \Sigma_{aa})$$

**Figure 2.9** The plot on the left shows the contours of a Gaussian distribution $p(x_a, x_b)$ over two variables, and the plot on the right shows the marginal distribution $p(x_a)$ (blue curve) and the conditional distribution $p(x_a|x_b)$ for $x_b = 0.7$ (red curve).

### 2.3.3 Bayes' theorem for Gaussian variables

Linear Gaussian model example

Gaussian marginal dist. $p(x)$, Gaussian conditional dist $p(y|x)$

$p(y|x)$ has a mean as a linear function of $x$ and a covariance which is independent of $x$.

i.e.
$$p(x) = \mathcal{N}(x \mid \mu, \Lambda^{-1})$$
$$p(y|x) = \mathcal{N}(y \mid Ax+b, L^{-1})$$

$x \in M-\text{dim}$

$y \in D-\text{dim}$

where $\mu$, $A$ and $b$ are parameters governing the means, and $\Lambda$ and $L$ are precision matrices.

We will find $p(y)$ and $p(x|y)$.

<span style="color:red">Known</span>
<span style="color:red">$p(x)$</span>
<span style="color:red">$p(y|x)$</span>

marginal ─── conditional

Let $z := \begin{pmatrix} x \\ y \end{pmatrix}$ and us consider the joint prob. dist

<span style="color:red">$p(z) = p(y|x)\, p(x)$</span>

$\ln p(z) = \ln p(x) + \ln p(y|x)$     (2.102)

$$= -\frac{1}{2}(x - \mu_1)^T \Lambda (x - \mu_1)$$

indep. of $x, y$.

$$-\frac{1}{2}(y - Ax - b)^T L (y - Ax - b) + \text{const}$$

This is a quadratic function of the component of $z$, $x, y$

hence $p(z)$ is a Gaussian

Consider the second term in (2.102)

$$-\frac{1}{2}x^T(\Lambda + A^TLA)x - \frac{1}{2}y^TLy + \frac{1}{2}y^TLAx + \frac{1}{2}x^TA^TLy$$

$$= -\frac{1}{2}\begin{pmatrix}x\\y\end{pmatrix}^T\begin{pmatrix}\Lambda + A^TLA & -A^TL\\-LA & L\end{pmatrix}\begin{pmatrix}x\\y\end{pmatrix} = -\frac{1}{2}z^TRz$$

$z$ has precision (inverse of covariance) matrix given by

$$z = \begin{pmatrix}x\\y\end{pmatrix} \qquad \begin{pmatrix}\Lambda + A^TLA & -A^TL\\-LA & L\end{pmatrix}$$

$$\Rightarrow \quad \text{cov}[z] = R^{-1} = \begin{pmatrix}\Lambda^{-1} & \Lambda^{-1}A^T\\A\Lambda^{-1} & L^{-1} + A\Lambda^{-1}A^T\end{pmatrix} \qquad (2.105)$$

Consider the linear term in (2.102)

$$x^T \Lambda \mu - x^T A^T L b + y^T L b = \binom{x}{y}^T \left( \frac{\Lambda \mu - A^T L b}{L b} \right)$$

$$\Rightarrow \quad \mathbb{E}[z] = R^{-1} \left( \frac{\Lambda \mu - A^T L b}{L b} \right)$$

$$z = \binom{x}{y} \qquad\qquad = \left( \frac{\mu}{A \mu + b} \right) \begin{matrix} x \\ y \end{matrix} \qquad (2.108)$$

$$A x + b$$

Using section 2.3.2 and $p(y) = \int p(z) \, dx$,

$$\mathbb{E}[y] = A \mu + b$$

$$\text{cov}[y] = L^{-1} + A \Lambda^{-1} A^T$$

Now we can find an expression for $p(x|y)$.

$$\mathbb{E}[x|y] = (\Lambda + A^T L A)^{-1} \{ A^T L (y - b) + \Lambda \mu \}$$

$$\text{cov}[x|y] = (\Lambda + A^T L A)^{-1}$$

## Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for $\mathbf{x}$ and a conditional Gaussian distribution for $\mathbf{y}$ given $\mathbf{x}$ in the form

$$
\begin{aligned}
p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) & (2.113) \\
p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) & (2.114)
\end{aligned}
$$

the marginal distribution of $\mathbf{y}$ and the conditional distribution of $\mathbf{x}$ given $\mathbf{y}$ are given by

$$
\begin{aligned}
p(\mathbf{y}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}) & (2.115) \\
p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) & (2.116)
\end{aligned}
$$

where

$$
\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}. \qquad (2.117)
$$

## 2.3.4 Maximum likelihood for the Gaussian

Data set $\mathbf{X} = (\mathbf{x}_1, \ldots, \mathbf{x}_N)^T$, $\{\mathbf{x}_n\}$ iid samples of $D$-dimensional Gaussian. The log likelihood function is given by

$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \Sigma) \qquad \color{red}{\mathbf{X} \quad N \times D \quad \text{matrix}}$$

$$= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Note that likelihood function depends only on the following two quantities

$$\sum_{n=1}^{N} \mathbf{x}_n , \qquad \color{magenta}{D \times D} \quad \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^T$$

These are known as sufficient statistics for Gaussian

$$\nabla_{\mu} \ln p(X \mid \mu, \Sigma) = \sum_{n=1}^{N} \Sigma^{-1} (x_n - \mu)$$

<span style="color:red">D - dim vector</span>

Set this gradient to zero vector, we obtain

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

solution of maximum likelihood estimator

MLE

(sample mean)

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

(sample covariance)

Remark

- $\Sigma_{ML}$ involves $\mu_{ML}$

- $\mu_{ML}$ is independent of $\Sigma_{ML}$

Evaluate the expectations of thes solutions under the true distribution. Then we obtain

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

$$\mathbb{E}[\Sigma_{ML}] = \frac{N-1}{N} \Sigma \qquad \text{biased}$$

## 2.3.5 Sequential estimation

### Sequential estimation for maximum likelihood

This method allows data points to be proceed one at time and then discarded and are important for on-line applications

Consider

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

which we will denote by $\mu_{ML}^{(N)}$ based on $N$ observations

$$\text{MI}_{ML}^{(N)} = \frac{1}{N} \sum_{n=1}^{N} x_n$$

$$= \frac{1}{N} x_N + \frac{1}{N} \sum_{n=1}^{N-1} x_n$$

$$= \frac{1}{N} x_n + \frac{N-1}{N} \text{MI}_{ML}^{(N-1)}$$

$$= \text{MI}_{ML}^{(N-1)} + \frac{1}{N} \underbrace{(x_N - \text{MI}_{ML}^{(N-1)})}_{\text{error signal}}$$

# General formulation of sequential learning (Robbins - Monro)

Two r.v. $z, \theta$ governed by a joint distribution $p(z, \theta)$

Define deterministic function $f(\theta)$ by

$$f(\theta) := \mathbb{E}_z[z|\theta] = \int z \, p(z|\theta) \, dz$$

e.g.

$$\mathbb{E}[t|x]$$

conditional expectation
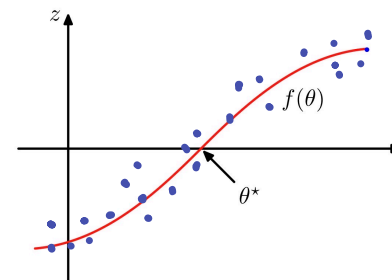
which is a function of $\theta$ (called regression function)

Find the root $\theta^*$ at which $f(\theta^*) = 0$
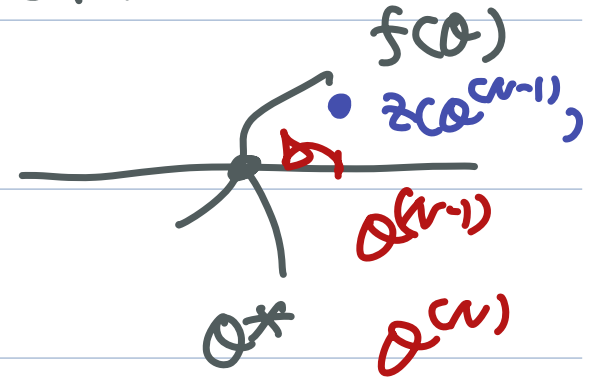
Suppose we can observe values of $z$ one at a time



**Figure 2.10** A schematic illustration of two correlated random variables $z$ and $\theta$, together with the regression function $f(\theta)$ given by the conditional expectation $\mathbb{E}[z|\theta]$. The Robbins-Monro algorithm provides a general sequential procedure for finding the root $\theta^*$ of such functions.

Assume the conditional variance of $z$ is finite

$$\mathbb{E}_z[(z - f) | \theta] < \infty$$

$\theta^*$ solution

$f(\theta)$

$z(\theta^{(N-1)})$

$\theta^{(N-1)}$

$\theta^*$   $\theta^{(N)}$

and wlog $f(\theta) > 0$ for $\theta > \theta^*$ and $f(\theta) < 0$ for $\theta < \theta^*$

A sequence of successive estimates of the root $\theta^*$ given by

$$\theta^{(N)} := \theta^{(N-1)} - a_{N-1} \, z(\theta^{(N-1)}) \qquad (2.129)$$

where $z(\theta^{(N)})$ is an observed value of $z$ when $\theta = \theta^{(N)}$

$\{a_N\}$ represents a seq of positive numbers satisfying

$$\lim_{N \to \infty} a_N = 0$$

$$\sum_{N=1}^{\infty} a_N = \infty \qquad \text{e.g. } \frac{1}{N}$$

$$\sum_{N=1}^{\infty} a_N^2 < \infty$$

of $f(\theta)$

By [Robbins - Monro], (2.129) converges to the root with

probability one.

Remark

— Third condition ensures that the accumulated noise has

finite variance and hence does not spoil convergence.

## General Maximum likelihood problem

$$f(\theta) = \int z\, p(z|\theta)\, dz$$

$$\|$$

By definition of $\theta_{ML}$, $\theta_{ML}$ satisfies

$$E_x\left[\frac{\partial}{\partial\theta} \ln p(x|\theta)\right]$$

$$\frac{\partial}{\partial\theta}\left\{-\frac{1}{N}\sum_{n=1}^{N} \ln p(x_n|\theta)\right\}\Bigg|_{\theta_{ML}} = 0$$

Taking $N \to \infty$ and exchanging derivative and summation,

$$-\lim_{N\to\infty} \frac{1}{N}\sum_{n=1}^{N} \frac{\partial}{\partial\theta} \ln p(x_n|\theta) = E_x\left[-\frac{\partial}{\partial\theta} \ln p(x|\theta)\right]$$

$$E_x[f(x)] \approx \lim_{N\to\infty} \frac{1}{N}\sum_{n=1}^{N} f(x_n) \qquad \text{observations } x_n$$

I.e. find the root of a regression function

Apply Robbins - Monro procedure

$$\theta^{(N)} := \theta^{(N-1)} - a_{N-1} \frac{\partial}{\partial \theta^{(N-1)}} \left[ - \ln p(x_N | \theta^{(N-1)}) \right] \qquad (2.135)$$

Specific example: sequential estimation of the mean of Gaussian distribution

In this case $\theta^{(N)}$ is the $\mu_{ML}^{(N)}$ mean of the Gaussian and $z$ is given by

$$(2.136)$$

$$z = \frac{\partial}{\partial \mu_{ML}} \ln p(x | \mu_{ML}, \sigma^2) = - \frac{1}{\sigma^2} (x - \mu_{ML})$$

Substituting (2.136) into (2.135) with $a_N = \sigma^2/N-1$

then we obtain (2.126)

## 2.3.6 Bayesian inference for the Gaussian

MLE method gave point estimates for $\mu$, $\Sigma$ (section 2.3.4)

Now develop a Bayesian treatment

Single Gaussian random variable $x$. Suppose $\sigma^2$ is known.

Aim to inference $\mu$ given $N$ observations $X = \{x_1, .. x_N\}$

The likelihood function is given by

$$p(X|\mu) = \prod_{n=1}^{N} p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu)^2\right\}$$

Note that this function is the form of the exponential of

a quadratic form of $\mu$.

We will choose a prior $p(\mu)$ given by Gaussian because

the product of two exponentials of quadratic function of

$\mu$ will also be Gaussian

Take prior prob. $p(\mu)$ to be

$$p(\mu) = N(\mu \mid \mu_0, \sigma_0^2)$$

$\mu_0, \sigma_0^2$ hyperparameters

Posterior

$$p(\mu \mid x) \propto p(x \mid \mu) \cdot p(\mu) \overset{N(\mu \mid \mu_0, \sigma_0^2)}{{}''{}}$$

Exercise 2.38, we obtain

$$p(\mu \mid x) = \mathcal{N}(\mu \mid \mu_N, \sigma_N^2)$$

$$\frac{1}{N}\sum x_n$$

where
$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\,\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\,\mu_{ML}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

## Remark

- $\mu_N$ is a compromise between $\mu_0$ and $\mu_{ML}$

- Effect of change in value $N$

- Precision is additive, if $N \to \infty$, $\sigma_N^2 \to 0$

- When $N$ is finite, if $\sigma_o^2 \to \infty$, then the posterior mean reduces to $\mu_{ML}$ and variance $\sigma_N^2$ becomes $\frac{\sigma^2}{N}$

Sequential inference in Bayesian paradigm

$$P(\mu \mid \mathcal{X}) \propto \left[ P(\mu) \prod_{n=1}^{N-1} P(x_n \mid \mu) \right] P(x_N \mid \mu)$$

Posterior

$\propto$ posterior distribution after observing $N-1$ data

$\mathcal{X} = \{ x_1, \dots x_N \}$

Now we wish to infer the variance and assume mean is known.

Let the precision $\lambda := 1/\sigma^2$. The likelihood function of $\lambda$

$$p(x \mid \lambda) = \prod_{n=1}^{N} \mathcal{N}(x_n \mid \mu, \lambda^{-1}) \propto \lambda^{N/2} \exp\left\{ -\frac{\lambda}{2} \sum_{n=1}^{N} (x - \mu)^2 \right\}$$

$$\lambda > 0$$

i.e. the form of $\lambda^{N/2} \cdot \exp\left(-\frac{\lambda}{2}\right)$
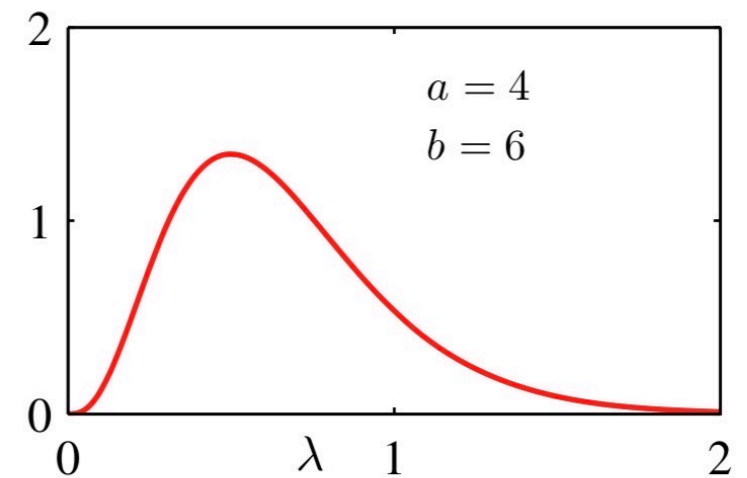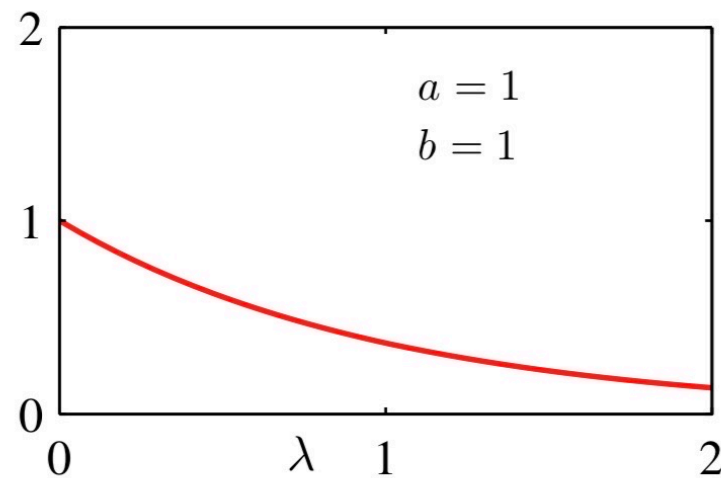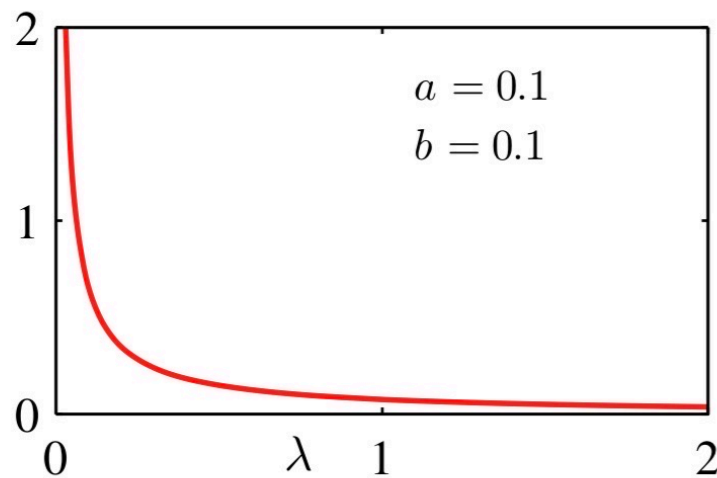
The corresponds to gamma distribution

$$a, b > 0$$

$$\text{Gam}(\lambda \mid a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda), \qquad \lambda > 0$$

Here $\Gamma(a)$ is a gamma function $\qquad \Gamma(a) = \int_0^\infty u^{a-1} e^{-u} du$

# Remark

- If $a > 0$, gamma distribution has finite integral.

- If $a \geq 1$, the distribution itself is finite.

- 
$$\mathbb{E}[\lambda] = \frac{a}{b} \qquad \text{Var}[\lambda] = \frac{a}{b^2}$$



**Figure 2.13** Plot of the gamma distribution $\text{Gam}(\lambda|a, b)$ defined by (2.146) for various values of the parameters $a$ and $b$.

Consider the prior dist. $\text{Gam}(\lambda \mid a_0, b_0)$. ($a_0, b_0$: hyperparameter)

The posterior dist. of $\lambda$ is as below

$$p(\lambda \mid x) \propto \underbrace{\lambda^{N/2} \exp\left\{ -\frac{\lambda}{2} \sum_{n=1}^{N} (x_n - \mu)^2 \right\}}_{\text{likelihood function of } \lambda} \cdot \underbrace{\text{Gam}(\lambda \mid a_0, b_0)}_{\text{prior of } \lambda}$$

$$\propto \lambda^{a_0 - 1} \lambda^{N/2} \exp\left\{ -b_0 \lambda - \frac{\lambda}{2} \sum_{n=1}^{N} (x_n - \mu)^2 \right\} \qquad (2.149)$$

$$\Rightarrow \quad p(\lambda \mid x) = \text{Gam}(\lambda \mid a_N, b_N) \quad \text{where}$$

$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \sum_{n=1}^{N} (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2$$

# Remark

- Effect of observing $N$ data points

  - increases the value of $a$ by $\frac{N}{2}$

  - " $b$ by $\frac{N}{2} \sigma_{ML}^2$

- We can interprete the parameter $a_0$ in terms of $2a_0$

'effective' prior observations.

- $\mathbb{E}[\lambda | \ast] = \dfrac{a_N}{b_N} = \dfrac{2a_0 + N}{2b_0 + N \sigma_{ML}^2}$

$\text{var}[\lambda | \ast] = \dfrac{a_N}{b_N^2}$

$\mathbb{E}[\lambda] \longrightarrow \dfrac{1}{\sigma_{ML}^2}$

$\lambda = \dfrac{1}{\sigma^2}$

Now suppose that the both $\mu$ and $\lambda$ are unknown

Consider the dependence of the likelihood function on $\mu$ and $\lambda$

$$p(x \mid \mu, \lambda) = \prod_{n=1}^{N} \left( \frac{\lambda}{2\pi} \right)^{1/2} \exp\left\{ -\frac{\lambda}{2}(x_n - \mu)^2 \right\}$$

$$\lambda, \mu \qquad \propto \left[ \lambda^{1/2} \exp\left( -\frac{\lambda \mu^2}{2} \right) \right]^{N} \exp\left[ \lambda \mu \sum_{n=1}^{N} x_n - \frac{\lambda}{2} \sum_{n=1}^{N} x_n^2 \right]$$

Thus the prior distribution should take the form

$$p(\mu, \lambda) \propto \left[ \lambda^{1/2} \exp\left( -\frac{\lambda \mu^2}{2} \right) \right]^{\beta} \exp\{ c\lambda\mu - d\lambda \}$$

(2.153)

$$= \exp\left\{ -\frac{\beta\lambda}{2}\left( \mu - \frac{c}{\beta} \right)^2 \right\} \underbrace{\lambda^{\beta/2} \exp\left\{ -\left( d - \frac{c^2}{2\beta} \right)\lambda \right\}}_{\lambda}$$

$$\underbrace{\phantom{= \exp\left\{ -\frac{\beta\lambda}{2}\left( \mu - \frac{c}{\beta} \right)^2 \right\}}}_{\mu, \lambda}$$

where $c$, $d$ and $\beta$ are constants. Use $p(\mu, \lambda) = p(\mu | \lambda) \, p(\lambda)$.

$p(\mu | \lambda)$: a Gaussian whose precision is a linear function of $\lambda$

$p(\lambda)$: a gamma distribution. So we take a prior

$$p(\mu, \lambda) = N(\mu | \mu_0, (\beta \lambda)^{-1}) \; \text{Gam}(\lambda | a, b) \qquad (2.154)$$

where $\mu_0 := c/\beta$, $a := (1+\beta)/2$, $b := d - \dfrac{c^2}{2\beta}$

(2.154) is called normal gamma or Gaussian gamma.

Note that it is not the simply the product of an independent Gaussian prior and gamma prior.

Multivariate Gaussian $N(x \mid \mu, \Lambda^{-1})$ for D-dim $x$

First, when precision matrix $\Lambda$ is known, the conjugate prior

distribution is again a Gaussian.

Second, for known mean and unknown precision matrix $\Lambda$,    D×D

the conjugate prior distribution is the Wishart distribution

given by

$$\mathcal{W}(\Lambda \mid W, \nu) = B |\Lambda|^{(\nu - D - 1)/2} \exp\left(-\frac{1}{2} \mathrm{Tr}(W^{-1} \Lambda)\right)$$

(D×D for $W$)     (trace of matrix over Tr)

where $\nu$ is called the number of degrees of freedom

$W$ is a $D \times D$ scale matrix. The nomalization constant $B$ is given by

$$B(W, \nu) = |W|^{-\nu/2} \left( 2^{\nu D/2} \pi^{D(D-1)/4} \prod_{i=1}^{D} \Gamma\left( \frac{\nu + 1 - i}{2} \right) \right)^{-1}.$$

If both the mean and precision are unknown, the conjugate prior is given by

$$p(\mu, \Lambda \mid \mu_0, \beta, W, \nu) = N(\mu \mid \mu_0, (\beta \Lambda)^{-1}) \, W(\Lambda \mid W, \nu)$$

which is known as the normal - Wishart or Gaussian - Wishart

## 2.3.7 Student's t-distribution

Conjugate prior for the precision of a Gaussian is given by a gamma distribution.

Consider univariate Gaussian $N(x \mid \mu, \tau^{-1})$ with Gamma prior $Gam(\tau \mid a, b)$. Integrate out the precision

$$p(x \mid \mu, a, b) = \int_0^\infty N(x \mid \mu, \tau^{-1}) \, Gam(\tau \mid a, b) \, d\tau \qquad \tau > 0$$

$$= \int_0^\infty \frac{b^a e^{(-b\tau)} \tau^{a-1}}{\Gamma(a)} \left( \frac{\tau}{2\pi} \right)^{1/2} \exp\left\{ -\frac{\tau}{2}(x-\mu)^2 \right\} d\tau$$

$$= \frac{b^a}{\Gamma(a)} \left( \frac{1}{2\pi} \right)^{1/2} \left[ b + \frac{(x-\mu)^2}{2} \right]^{-a-1/2} \Gamma(a+1/2)$$

where we have made the change of variable $z = \tau[b + (x-\mu)^2/2]$

Define new parameters $\nu = 2a$ and $\lambda = a/b$.

$$St(x \mid \mu, \lambda, \nu) = \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi\nu}\right)^{1/2} \left[1 + \frac{\lambda(x-\mu)^2}{\nu}\right]^{-\nu/2 - 1/2}$$

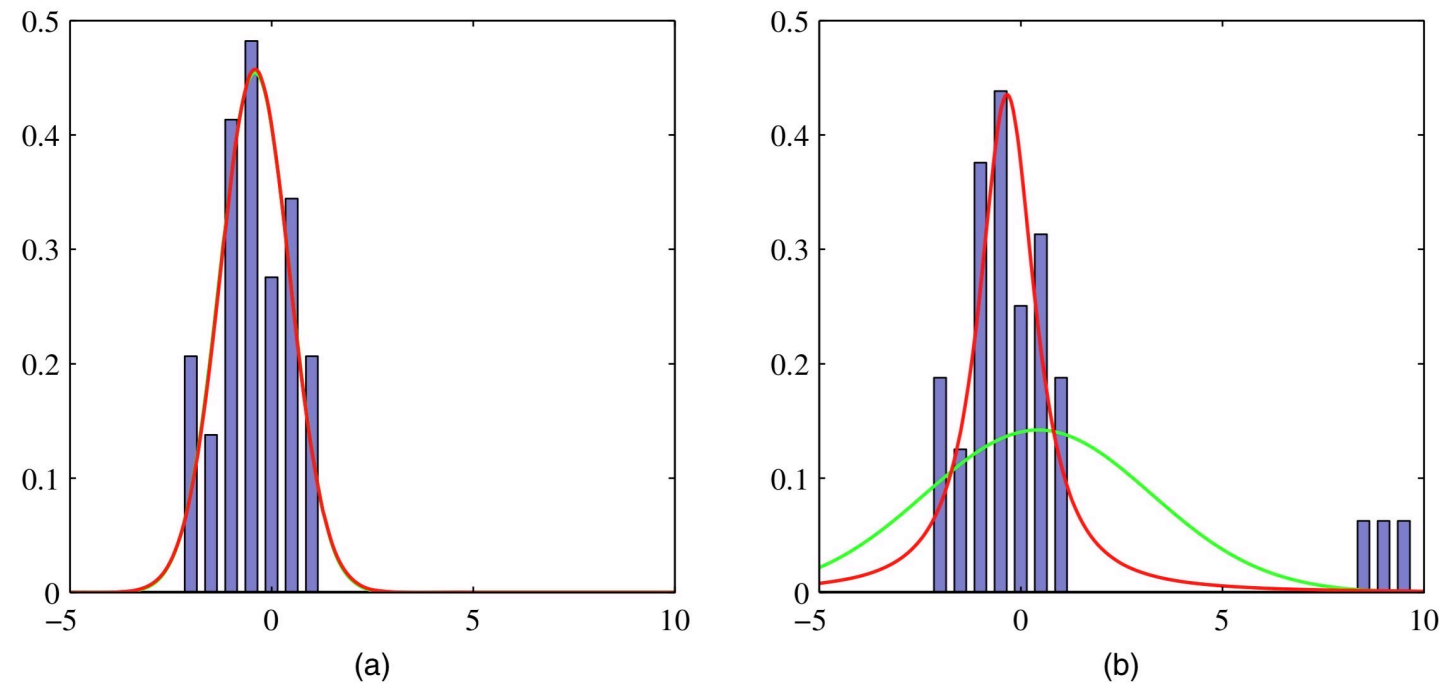known as Student's $t$-distribution. $\lambda$ is called precision and $\nu$ is called the degree of freedom.

When $\nu = 1$, $t$-distribution reduces to the Cauchy dist.

While in the limit $\nu \to \infty$, $t$-distribution becomes Gaussian $N(x \mid \mu, \lambda^{-1})$

# Remark

– t - dist. can be interpreted infinite mixture of Gaussian

– Longer tail, robustness property



**Figure 2.16** Illustration of the robustness of Student's t-distribution compared to a Gaussian. (a) Histogram distribution of 30 data points drawn from a Gaussian distribution, together with the maximum likelihood fit obtained from a t-distribution (red curve) and a Gaussian (green curve, largely hidden by the red curve). Because the t-distribution contains the Gaussian as a special case it gives almost the same solution as the Gaussian. (b) The same data set but with three additional outlying data points showing how the Gaussian (green curve) is strongly distorted by the outliers, whereas the t-distribution (red curve) is relatively unaffected.

# Multivariate Student's $t$-distribution

$$St(x \mid \mu, \Lambda, \nu) = \int_0^\infty N(x \mid \mu, (\eta \Lambda)^{-1}) \, Gam(\eta \mid \nu/2, \nu/2) \, d\eta$$

$$= \frac{\Gamma(D/2 + \nu/2)}{\Gamma(\nu/2)} \frac{|\Lambda|^{1/2}}{(\pi \nu)^{D/2}} \left[ 1 + \frac{\Delta^2}{\nu} \right]^{-D/2 - \nu/2}$$

where

$$\Delta := (x - \mu)^T \Lambda (x - \mu)$$

## Remark

- $\mathbb{E}[x] = \mu$      if    $\nu > 1$

- $cov[x] = \dfrac{\nu}{(\nu - 2)} \Lambda^{-1}$    if    $\nu > 2$

- $mode[x] = \mu$

## 2.3.8 Periodic variables

Consider an angular (polar) coordinate $0 \leq \theta < 2\pi$ and the problem of evaluating the mean of observations $D = \{\theta_1, \dots \theta_N\}$.

Simple average $(\theta_1 + \dots \theta_N)/N$ is strongly coordinate dependent.

Set angular observations as points on unit circle.

Let $x_i$ be two-dim vector with $x_i = (\cos \theta_i, \sin \theta_i)$

Average the vectors $\{\mathbf{x}_n\}$ instead to give
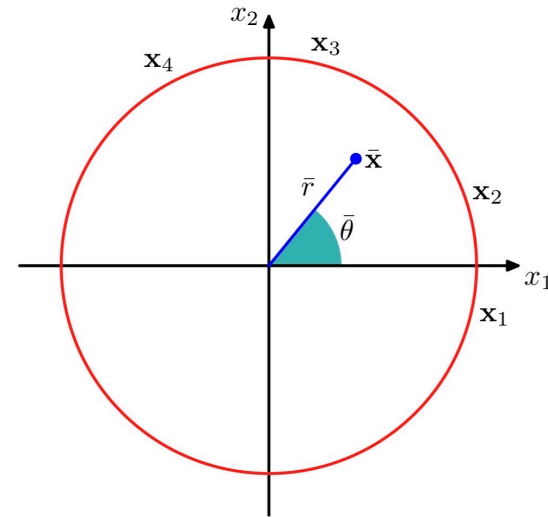
$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n$$



**Figure 2.17** Illustration of the representation of values $\theta_n$ of a periodic variable as two-dimensional vectors $\mathbf{x}_n$ living on the unit circle. Also shown is the average $\bar{\mathbf{x}}$ of those vectors.

$$= \bar{r} (\cos\bar{\theta}, \sin\bar{\theta})$$

i.e. $\quad \bar{r}\cos\bar{\theta} = \frac{1}{N} \sum_{n=1}^{N} \cos\theta_n, \quad \bar{r}\sin\bar{\theta} = \frac{1}{N} \sum_{n=1}^{N} \sin\theta_n$

Thus we can solve for $\bar{\theta}$ to give

$$\bar{\theta} = \tan^{-1}\left\{ \frac{\sum_n \sin\theta_n}{\sum_n \cos\theta_n} \right\}$$
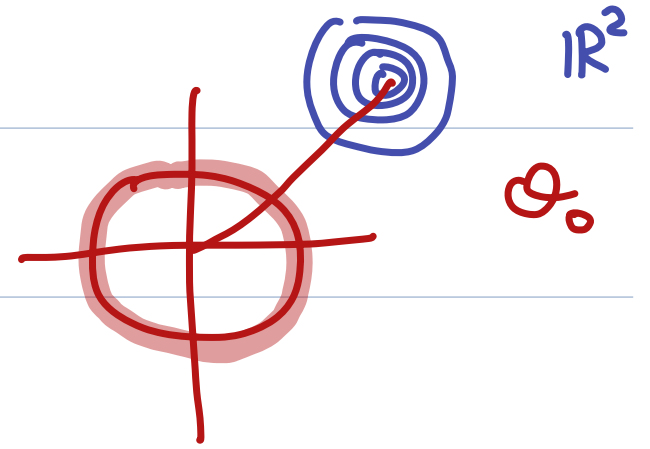
Consider $p(\theta)$ that have period $2\pi$ and must satisfies

$$p(\theta) \geq 0$$

$$\int_0^{2\pi} p(\theta) \, d\theta = 1$$

$$p(\theta + 2\pi) = p(\theta)$$

We can easily obtain a Gaussian-like distribution.

Consider a Gaussian over $x = (x_1, x_2)$ having mean $\mu = (\mu_1, \mu_2)$

and covariance matrix $\Sigma = \sigma^2 I$ so that

$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2}\right\} \quad (2.173)$$

Map $\mathbb{X} = (x_1, x_2)$ and $\mu$ into polar coordinates

$$x_1 = r\cos\theta, \qquad x_2 = r\sin\theta$$

$$\Bigg( \qquad \mu_1 = r_0\cos\theta_0, \qquad \mu_2 = r_0\sin\theta_0 \qquad \text{fixed } r_0, \theta_0$$

Substitute these transformation into (2.173) with $r=1$ condition

The exponent in (2.173)

$$-\frac{1}{2\sigma^2}\left\{ (r\cos\theta - r_0\cos\theta_0)^2 + (r\sin\theta - r_0\sin\theta_0)^2 \right\}$$

$$(r=1)$$

$$= -\frac{1}{2\sigma^2}\left\{ 1 + r_0^2 - 2r_0\cos\theta\cos\theta_0 - 2r_0\sin\theta\sin\theta_0 \right\}$$

$$= \frac{r_0}{\sigma^2}\cos(\theta - \theta_0) + \text{const}$$

Define $m = r_0 / \sigma^2$. Then we obtain the expression for the distribution of $p(\theta)$ along unit circle

$$p(\theta \mid \theta_0, m) = \frac{1}{2\pi} \frac{1}{I_0(m)} \exp\{ m \cos(\theta - \theta_0) \} \qquad 0 \le \theta \le 2\pi$$

which called von Mises distribution. Here $\theta_0$ represents the mean and $m = r_0 / \sigma^2$ is called concentration parameter.

$I_0(m)$ : zeroth-order Bessel function of the first kind

$$I_0(m) := \frac{1}{2\pi} \int_0^{2\pi} \exp\{ m \cos\theta \} \, d\theta$$

Now consider the maximum likelihood for $\theta_0$ and $m$

Observations $D = \{\theta_1, \ldots \theta_N\}$ is given

$$\ln P(D \mid \theta_0, m) = \prod_{n=1}^{N} P(\theta_n \mid \theta_0, m) \tag{2.181}$$

$$= -N \ln(2\pi) - N \ln I_0(m) + m \sum_{n=1}^{N} \cos(\theta_n - \theta_0)$$

Set the derivative w.r.t $\theta_0$ equal to zero gives

$$\sum_{n=1}^{N} \sin(\theta_n - \theta_0) = 0$$

Thus, we obtain $\quad \theta_0^{ML} = \tan^{-1} \left\{ \dfrac{\sum \sin \theta_n}{\sum \cos \theta_n} \right\}$

Similarly maximizing (2.181) w.r.t m. Set the derivative of (2.181) w.r.t m, then we have

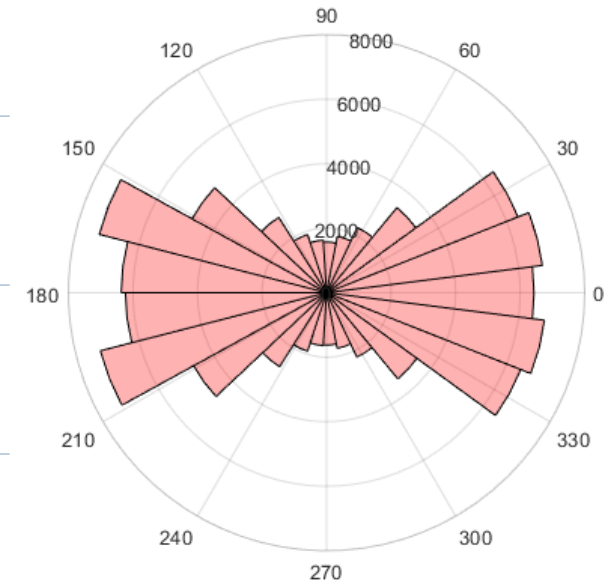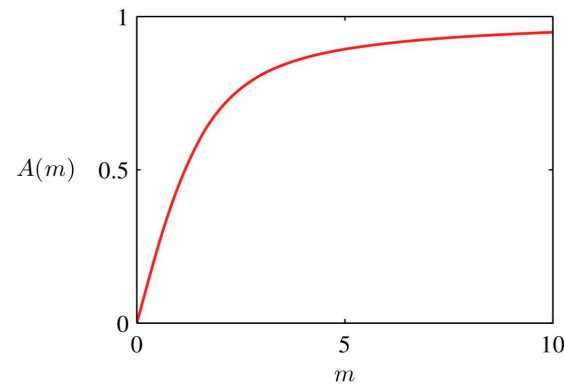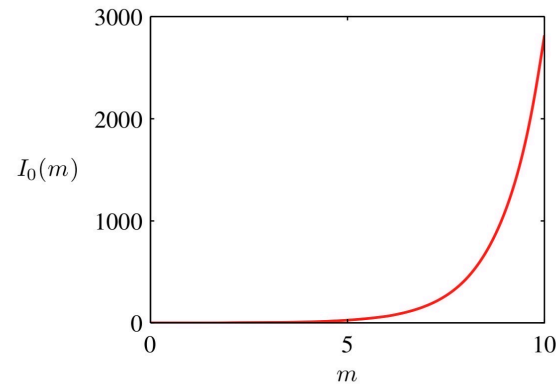$$A(m) = \frac{1}{N} \sum_{n=1}^{N} \cos(\theta_n - \theta_0^{ML}) \qquad (2.185)$$

where we used $I_0'(m) = I_1(m)$ and have defined

$$A(m) := \frac{I_1(m)}{I_0(m)}$$

We can rewrite (2.185) in the form

$$A(m_{ML}) = \left(\frac{1}{N} \sum_{n=1}^{N} \cos\theta_n\right) \cos\theta_0^{ML} + \left(\frac{1}{N} \sum_{n=1}^{N} \sin\theta_n\right) \sin\theta_0^{ML}$$

Here $A(m)$ can be inverted unmerically.



**Figure 2.20** Plot of the Bessel function $I_0(m)$ defined by (2.180), together with the function $A(m)$ defined by (2.186).
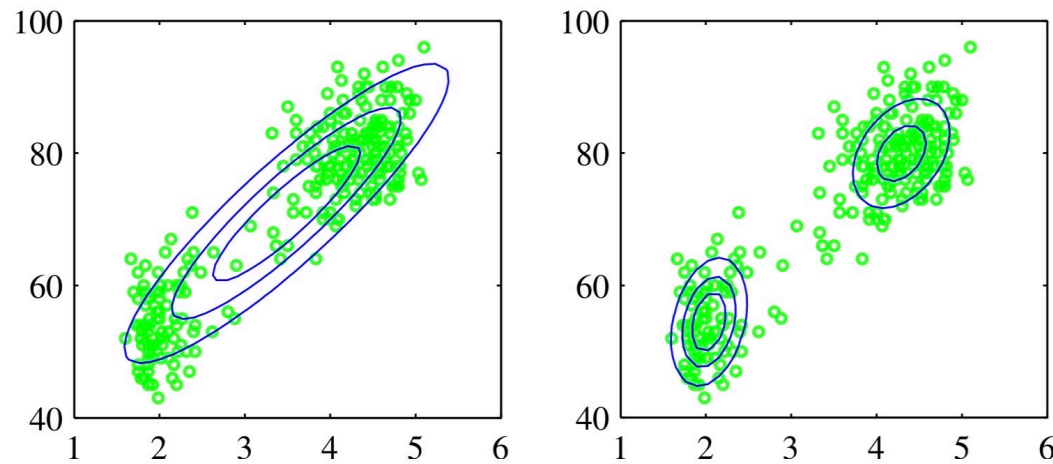
Remark: other techniques to construct periodic variable

- Histogram in polar coordinates

- Mixtures of von Mises distributions

# 2.3.9 Mixtures of Gaussians

## Limitations of a Gaussian (unimodal)

**Figure 2.21** Plots of the 'old faithful' data in which the blue curves show contours of constant probability density. On the left is a single Gaussian distribution which has been fitted to the data using maximum likelihood. Note that this distribution fails to capture the two clumps in the data and indeed places much of its probability mass in the central region between the clumps where the data are relatively sparse. On the right the distribution is given by a linear combination of two Gaussians which has been fitted to the data by maximum likelihood using techniques discussed Chapter 9, and which gives a better representation of the data.



Mixture distribution : linear combinations of basic distributions.

Mixture of Gaussians : superposition of K Gaussians

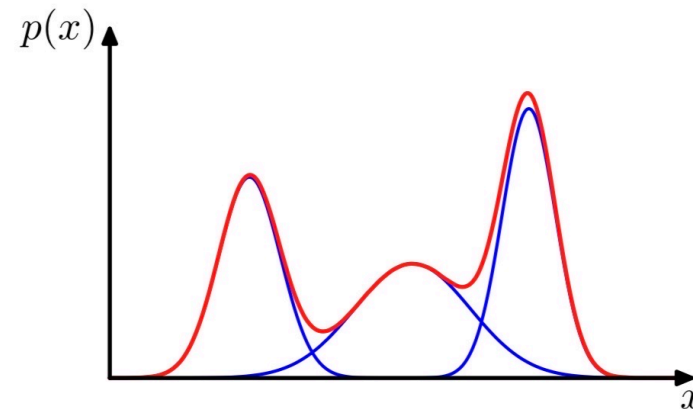$$p(x) = \sum_{k=1}^{K} \pi_k \, \mathcal{N}(x \mid \mu_k, \Sigma_k)$$

Each $\mathcal{N}(x \mid \mu_k, \Sigma_k)$ is called a component of mixture

The parameter $\pi_k$ are called mixing coefficients and satisfies

$$\sum_{k=1}^{K} \pi_k = 1 \qquad , \qquad 0 \leq \pi_k \leq 1$$

**Figure 2.22** Example of a Gaussian mixture distribution in one dimension showing three Gaussians (each scaled by a coefficient) in blue and their sum in red.

$p(x)$ can be rewrite in the form

$$p(x) = \sum_{k=1}^{K} p(k)\, p(x \mid k)$$

$\pi_k = p(k)$ prior probability of picking the $k^{th}$ component

$\mathcal{N}(x \mid \mu_k, \Sigma_k) = p(x \mid k)$ probability of $x$ conditioned on $k$.

Consider the posterior $P(K|X)$ a.k.a. responsibilities

$$\gamma_K(x) := P(K|X) = \frac{P(K,X)}{P(X)} = \frac{P(K)\,P(X|K)}{\sum_\ell P(\ell)\,P(X|\ell)}$$

$$= \frac{\pi_K\,\mathcal{N}(X|\mu_K,\Sigma_K)}{\sum_\ell \pi_\ell\,\mathcal{N}(X|\mu_\ell,\Sigma_\ell)}$$

Gaussian mixture is governed by $\pi := \{\pi_1 \dots \pi_K\}$, $\mu := \{\mu_1, \dots \mu_K\}$

and $\Sigma := \{\Sigma_1, \dots \Sigma_K\}$

One way to set these parameters is to use maximum likelihood.

$$\ln p(X \mid \pi, \mu, \Sigma) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \, N(x_n \mid \mu_k, \Sigma_k) \right\}$$

Maximum likelihood solution for the parameters no longer has a closed-form analytical solution.

Expectation maximization (chapter 9)

## 2.4 The exponential family

Broad class of distributions called the exponentials family

Random vector $x$, parameters $\eta$ (called natural parameters)

$$p(x \mid \eta) := h(x) \, g(\eta) \, \exp\{ \eta^T u(x) \} \tag{2.194}$$

Here $u(x)$ is some function of $x$ and $g(\eta)$ can be interpreted

as the normalization coefficient. i.e.

$$g(\eta) \int h(x) \, \exp\{ \eta^T u(x) \} \, dx = 1 \tag{2.195}$$

Recall    Bernoulli   distribution

$$p(x \mid \mu) = \mathrm{Bern}(x \mid \mu) = \mu^{x}(1-\mu)^{1-x} \qquad x = 0, 1$$

$$= \exp\{ x \ln \mu + (1-x) \ln(1-\mu)\}$$

$$= (1-\mu) \exp\left\{ \ln\left(\frac{\mu}{1-\mu}\right) x \right\}$$

Comparison   with  (2.194)

$$\eta = \ln\left(\frac{\mu}{1-\mu}\right)$$

Solve for $\mu$ to give

$$\mu = \sigma(\eta) := \frac{1}{1 + \exp(-\eta)} \qquad \text{logistic sigmoid}$$

Thus, Bernoulli distribution can be rewrited in the form

$$p(x|\eta) = \sigma(-\eta) \exp(\eta x)$$

we have used $\sigma(-\eta) = 1 - \sigma(\eta)$. Comparison with (2.194)

$$u(x) = x$$
$$h(x) = 1$$
$$g(\eta) = \sigma(-\eta)$$

Next consider the multinomial distribution

$$p(x \mid \mu) = \prod_{k=1}^{M} \mu_k^{x_k} = \exp\left\{ \sum_{k=1}^{M} x_k \ln \mu_k \right\}$$

where $x = (x_1, \dots x_M)^T$ (one-hot vector)

The standard representation (2.149) so that

$$p(x \mid \eta) = \exp(\eta^T x)$$

where $\eta = (\eta_1, \dots \eta_M)^T$ with $\eta_k = \ln \mu_k$. i.e.

$$u(x) = x, \qquad h(x) = 1, \qquad g(\eta) = 1$$

Since $\sum_{k=1}^{M} \mu_k = 1$, parameters $\eta_k$ are not independent.

i.e. $\mu_M = 1 - \sum_{k=1}^{M-1} \mu_k$. left $M-1$ parameters

$$0 \leq \mu_k \leq 1, \qquad \sum_{k=1}^{M-1} \mu_k \leq 1$$

So the multinomial distribution becomes $\overset{1 - \sum_{k=1}^{M-1} x_k}{\underset{/\!/}{}}$

$$\exp\left\{ \sum_{k=1}^{M} x_k \ln \mu_k \right\} = \exp\left\{ \sum_{k=1}^{M-1} x_k \ln \mu_k + x_M \underset{\underset{\ln(1 - \sum_{k=1}^{M-1} \mu_k)}{/\!/}}{\ln \mu_M} \right\}$$

$$= \exp\left\{ \sum_{k=1}^{M-1} x_k \ln\left( \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \right) + \ln\left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right\}$$

Now identify

$$\ln \left( \frac{\mu_k}{1 - \sum_j \mu_j} \right) = \eta_k$$

and

$$\mu_k = \frac{\exp(\eta_j)}{1 + \sum_j \exp(\eta_j)}$$

soft max

normalized exponential

In this representation, multinomial distribution

$$p(x \mid \eta) = \left( 1 + \sum_{k=1}^{M-1} \exp(\eta_k) \right)^{-1} \exp(\eta^T x), \qquad x \in \mathbb{R}^{M-1}$$

$$\eta = (\eta_1 \dots \eta_{M-1})^T.$$

$$u(\mathbf{x}) = \mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_{M-1} \end{pmatrix} \qquad h(\mathbf{x}) = 1, \qquad g(\eta)) = \left( 1 + \sum_{k=1}^{M-1} \exp(\eta_k) \right)^{-1}$$

Finally, consider the univariate Gaussian distribution.

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

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

$$\eta = \begin{pmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{pmatrix}, \qquad u(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}, \qquad h(x) = (2\pi)^{-\frac{1}{2}} \qquad g(\eta)) = (-2\eta_2)^{\frac{1}{2}} \exp\left( \frac{\eta_1^2}{4\eta_2} \right)$$

## 2.4.1 Maximum likelihood and sufficient statistics

Consider the exponential family of distributions over $x$

$$p(x \mid \eta) = h(x)\, g(\eta)\, \exp\{\eta^T u(x)\} \qquad (2.194)$$

Taking the gradient of both side of

① $\quad g(\eta) \int h(x)\, \exp\{\eta^T u(x)\}\, dx = 1 \qquad (2.195)$

w.r.t $\eta$, we have

$$\nabla g(\eta) \int h(x)\, \exp\{\eta^T u(x)\}\, dx$$

$$+ \, g(\eta) \int h(x)\, \exp\{\eta^T u(x)\}\, u(x)\, dx = 0$$

Using (2.195) then

$$-\frac{1}{g(\eta)} \nabla g(\eta) = g(\eta) \int h(x) \exp\{\eta^T u(x)\} u(x) dx = \mathbb{E}[u(x)]$$

We therefore obtain the result

$$-\nabla \ln g(\eta) = \mathbb{E}[u(x)]$$

Now consider iid samples denoted by $X = \{x_1 \dots x_N\}$

for which likelihood

$$p(X|\eta) = \left(\prod_{n=1}^{N} h(x_n)\right) g(\eta)^N \exp\left\{\eta^T \sum_{n=1}^{N} u(x_n)\right\}$$

Setting the gradient of $\ln P(X|\eta)$ w.r.t $\eta$ to zero.

We get the following condition to be satisfied by $\eta_{ML}$

$$- \nabla \ln g(\eta_{ML}) = \frac{1}{N} \sum_{n=1}^{N} u(x_n)$$

Note that MLE depends on the data only through $\sum_n u(x_n)$

called sufficient statistic of $(2.194)$

Do not need to store the entire data.

E.g.  Bernoulli    $u(x) = x$ .         sum of $\{x_n\}$

Gaussian    $u(x) = (x, x^2)^T$.  sum of $\{x_n\}$ and $\{x_n^2\}$.

## 2.4.2 Conjugate prior

For a given prob. density $p(x|\eta)$, seek a prior $p(\eta)$
that is conjugate to the likelihood function.
( the posterior has the same functional form as the prior)

For exponential family (2.194), $\exists$ conjugate prior of $\eta$

$$p(\eta | \chi, \nu) = f(\chi, \nu) \, g(\eta)^{\nu} \exp\{\nu \eta^T \chi\}$$

where $f(\chi, \nu)$ is a normalization coefficient and $g(\eta)$ is
the same function in (2.194)

# The posterior

$$p(\eta \mid \mathcal{X}, x, \nu) \propto p(\mathcal{X} \mid \eta) \cdot p(\eta \mid x, \nu)$$

$$\propto g(\eta)^{\nu+N} \exp\left\{ \eta^T \left( \sum_{n=1}^{N} u(x_n) + \nu x \right) \right\}$$

## 2.4.3 Noinformative priors

Intend to have as litte influence on the posterior as possible.

Let density or likelihood is given by $p(x|\lambda)$.

Consider noninformative prior $p(\lambda)$

First, $p(\lambda) = $ constant

- If the domain of $\lambda$ is unbounded, prior cannot be normalized. Such prior is called improper.

- Transformation behavior of density under a nonlinear change of variables

# Example 1.

Density of $x$ takes the form

$$p(x \mid \mu) = f(x - \mu)$$

$\mu$ is known as location parameter. E.g. $N(x \mid \mu, \sigma^2)$

## Translation invariance

If $x \to \hat{x} := x + c$, then

$$\hat{p}(\hat{x} \mid \hat{\mu}) = f(\hat{x} - \hat{\mu})$$

where we have defined $\hat{\mu} := \mu + c$

Thus $p(x|\mu) = p(\hat{x}|\hat{\mu})$ so density is independent of origin.

Prior distribution should satisfy this translation invariance property.

$$\Rightarrow \int_A^B p(\mu)\, d\mu = \int_{A-c}^{B-c} p(\mu)\, d\mu = \int_A^B p(\mu-c)\, d\mu \qquad \forall A, B$$

So we have $p(\mu - c) = p(\mu)$

Example of location parameter: mean of a Gaussian

The conjugate prior for $\mu$ is again Gaussian $p(\mu | \mu_0, \sigma_0^2)$

and we obtain noninformative prior by taking $\sigma_0^2 \to \infty$.

# Example 2.

Density of $x$ takes the form

$$p(x \mid \sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right) \qquad \sigma > 0$$

$\sigma$ is known as scale parameter.    E.g. $N(x \mid \mu, \sigma^2)$

## Scale invariance

If $\quad x \rightarrow \hat{x} := cx$, then

$$\hat{p}(\hat{x} \mid \hat{\sigma}) = \frac{1}{\hat{\sigma}} f\left(\frac{\hat{x}}{\hat{\sigma}}\right)$$
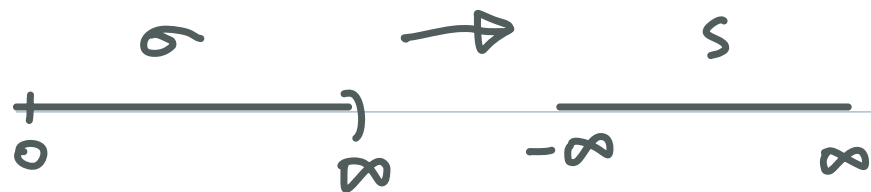
where we have defined $\hat{\sigma} = c\sigma$

So this transformation corresponds to a change of scale.

Prior distribution should satisfy this scale invariance property.

$$\Rightarrow \quad \int_A^B p(\sigma)\,d\sigma = \int_{A/c}^{B/c} p(\sigma)\,d\sigma = \int_A^B p(\tfrac{1}{c}\sigma)\tfrac{1}{c}\,d\sigma \quad \forall \; A, B$$

So we have $p(\sigma) = p(\tfrac{1}{c}\sigma)\tfrac{1}{c}$ and hence $p(\sigma) \propto \tfrac{1}{\sigma}$

Note that this is an improper prior because of $0 < \sigma < \infty$

$$\sigma \quad \longrightarrow \quad s$$

$$\underset{0}{\vdash} \quad \quad \underset{\infty}{)} \quad \quad \overline{\underset{-\infty}{\quad\quad} \underset{\infty}{\quad\quad}}$$

$$S(\sigma) = \ln\sigma \quad\quad ds = \tfrac{1}{\sigma}\,d\sigma$$

$$P_s(\ln\sigma) = \underbrace{P(\sigma)}_{\propto \frac{1}{\sigma}} \cdot \left|\frac{d\sigma}{ds}\right| = \text{constant}$$

Convenient to think of prior for scale parameter in terms of the density of the log of the parameter.

Example of scale parameter: standard deviation $\sigma$ of a Gaussian

$$N(x \mid \mu, \sigma^2) \propto \sigma^{-1} \exp\{-(\tilde{x}/\sigma)^2\}$$

where $\tilde{x} := x - \mu$

More convenient to work in terms of the precision

$\lambda = 1/\sigma^2$ rather than $\sigma$ itself

$d\lambda = \dfrac{1}{\sigma^3} d\sigma$

$\sigma \longrightarrow \lambda = \dfrac{1}{\sigma^2}$

$$\underset{0 \qquad\qquad\qquad \infty}{\vdash\!\!-\!\!-\!\!-\!\!-\!\!-\!\!-\!\!)} \qquad \underset{0 \qquad\qquad\qquad \infty}{\vdash\!\!-\!\!-\!\!-\!\!-\!\!-\!\!-\!\!}$$

$$\overset{\propto \frac{1}{\sigma}}{\underset{P_\lambda(\lambda) = P(\sigma) \cdot \left| \frac{d\sigma}{d\lambda} \right| \propto 1/\lambda}{}}$$

We have seen the conjugate prior for $\lambda$ was $\text{Gam}(\lambda \mid a_0, b_0)$

The noninformative is obtained as the special case $a_0 = b_0 = 0$.

## 2.5 Nonparametric Methods

Approaches to density modeling

Parametric vs Nonparametric (few assumptions)

Histogram method for density estimation

Single continuous random variable $x$. Partition $x$ into distinct bins of width $\Delta_i$ and then count the number $n_i$ of observations of $x$ falling in bin $i$
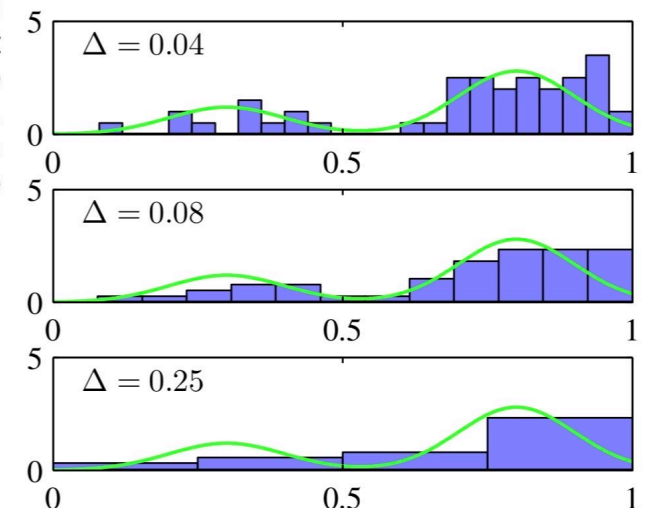
We obtain probability values for each bin given by

$$P_i = \frac{n_i}{N \Delta_i}$$

where $N$ is # of total observations. So a model for the density $p(x)$ is piecewise constant over the width $\Delta_i$ of each bin, and often the bins are chosen to have the same width $\Delta_i = \Delta$

Figure 2.24 An illustration of the histogram approach to density estimation, in which a data set of $50$ data points is generated from the distribution shown by the green curve. Histogram density estimates, based on (2.241), with a common bin width $\Delta$ are shown for various values of $\Delta$.

# Remark

- Effect of a choice of width $\Delta$ (smoothing parameter)

- After compting histogram, the data set can be discarded.

- Useful tool for a quick visualization of 1-d or 2-d data

- Limitation of high dimensional data $M^D$, $M$ bins in D-dim

## 2.5.1 Kernel density estimator

Estimate unknown probability density $p(x)$ in $D$-dim space.

Consider some small region $R$ containing $x$.

Then the probability mass associated with this $R$

$$P = \int_R p(x)\, dx \qquad (\text{true prob.})$$

Suppose we observed $N$ data set drawn from $p(x)$.

Since each point has a probability $P$ of falling within $R$

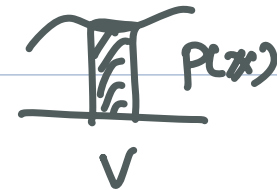$$\text{Bin}(k \mid N, P) = \binom{N}{k} P^k (1-P)^{N-k} \qquad k = 0, 1, \dots N$$

$$\Rightarrow \quad \mathbb{E}\left[\frac{K}{N}\right] = P, \quad \text{Var}\left[\frac{K}{N}\right] = P(1-P)/N$$

For large $N$,

$$K \simeq NP$$

If the region $R$ is sufficiently small that $p(x)$ is roughly constant over $R$, then

$$P \simeq p(x) V$$



where $V$ is the volume of $R$.

Combining these expressions we obtain the density estimate

$$p(x) = \frac{k}{NV} \qquad\qquad (2.246)$$

R 영역

Remark: two contradictory assumptions on R and K

We can exploit (2.246) in two different ways

- kernel density estimator (fix V)

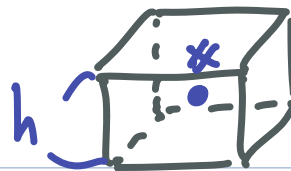- K nearest neighbour method (fix K)

Kernel method in detail (fix R and V)

Take the region R to be a small hypercube centered on $x$

To count the number K of points falling within R, define

$$k(u) := \begin{cases} 1 & |u_i| \leq \frac{1}{2} \quad i=1, \dots D \\ 0 & \text{otherwise} \end{cases}$$
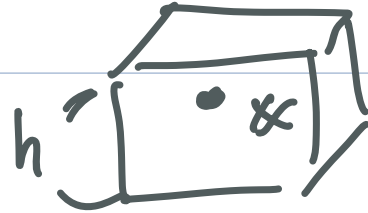
$k(u)$ is an example of kernel function, i.e. the quantity

$k((x - x_n)/h) = 1$   if $x_n$ lies in a cube of side $h$ centered

on $x$   otherwise $0$.

The total number of points lying in the cube

$$K := \sum_{n=1}^{N} k\left(\frac{x - x_n}{h}\right)$$



Substituting this expression into (2.246),

$$p(x) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{h^D} k\left(\frac{x - x_n}{h}\right)$$

where we have used $V = h^D$. (example of kernel density estimator)

We can obtain a smoother density model (Gaussian kernel)

$$p(x) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{ -\frac{\| x - x_n \|^2}{2h} \right\}$$

where $h$ represents the standard deviation of Gaussian component

and plays the role of a smoothing parameter.

Generally, we can choose any other kernel function $K(u)$
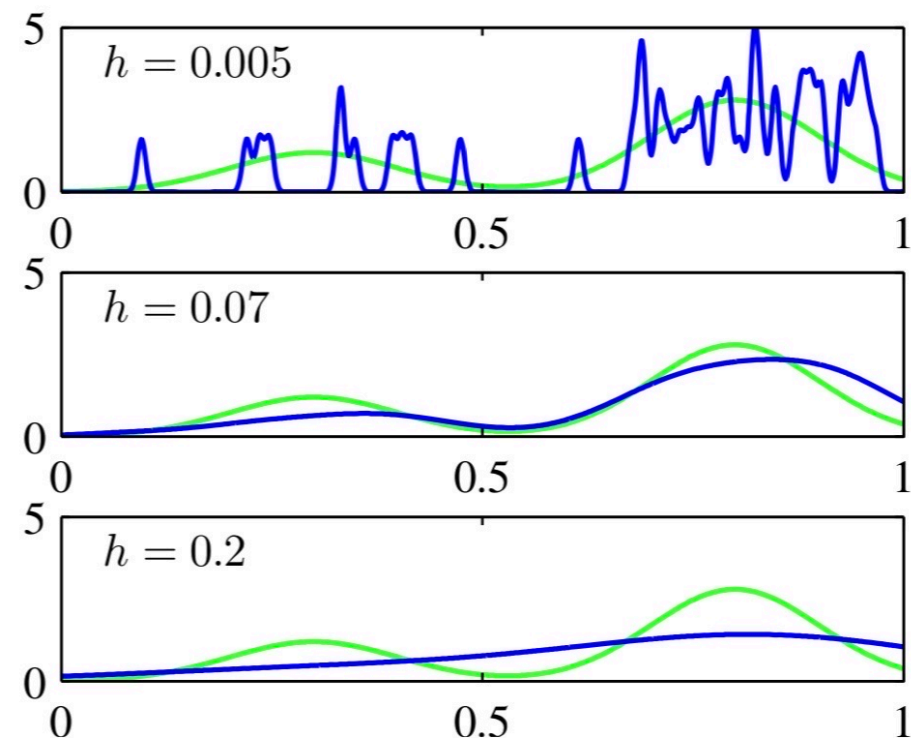
subject to

$$K(u) \geq 0$$

$$\int K(u) \, du = 1$$

# Remark

- No computation involved in the training phase

- Computational cost grows linearly with the data size.

**Figure 2.25** Illustration of the kernel density model (2.250) applied to the same data set used to demonstrate the histogram approach in Figure 2.24. We see that $h$ acts as a smoothing parameter and that if it is set too small (top panel), the result is a very noisy density model, whereas if it is set too large (bottom panel), then the bimodal nature of the underlying distribution from which the data is generated (shown by the green curve) is washed out. The best density model is obtained for some intermediate value of $h$ (middle panel).

## 2.5.2 Nearest - neighbour methods

### K nearest neighbours

For local density estimation, fix value of $K$ and use the data to find an appropriate value for $V$.

Consider a sphere centered on $x$ and allow the radius to grow until it contains $K$ data points. i.e. the radius is not determined (fixed)

The value of $K$ governs the degree of smoothing.

Use (2.246) $p(x) = \dfrac{K}{NV}$ with KNN method for density estimation

KNN method can be extended to classification.

Apply KNN to each class separately and then make use of Baye's theorem.

$N$: total # of data set, $N_k$: # of points in $C_k$

i.e. $\sum_{k=1}^{K} N_k = N$

New point $x$ (fixed). Draw a sphere centered on $x$ containing precisely $K$ points irrespective of their class. This sphere has the volume $V$ and contains $K_k$ points from the class $C_k$.
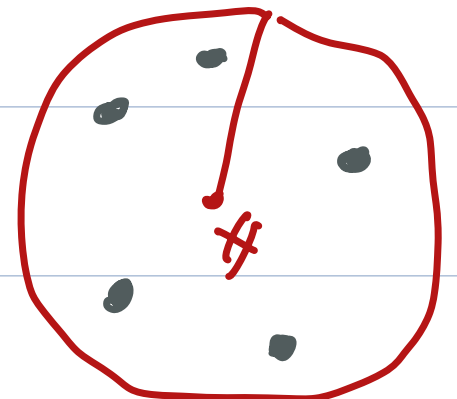
$$P(x \mid C_k) = \frac{K_k}{N_k V}$$

Similarly, the unconditioned density is given by

$$P(x) = \frac{K}{N V}$$
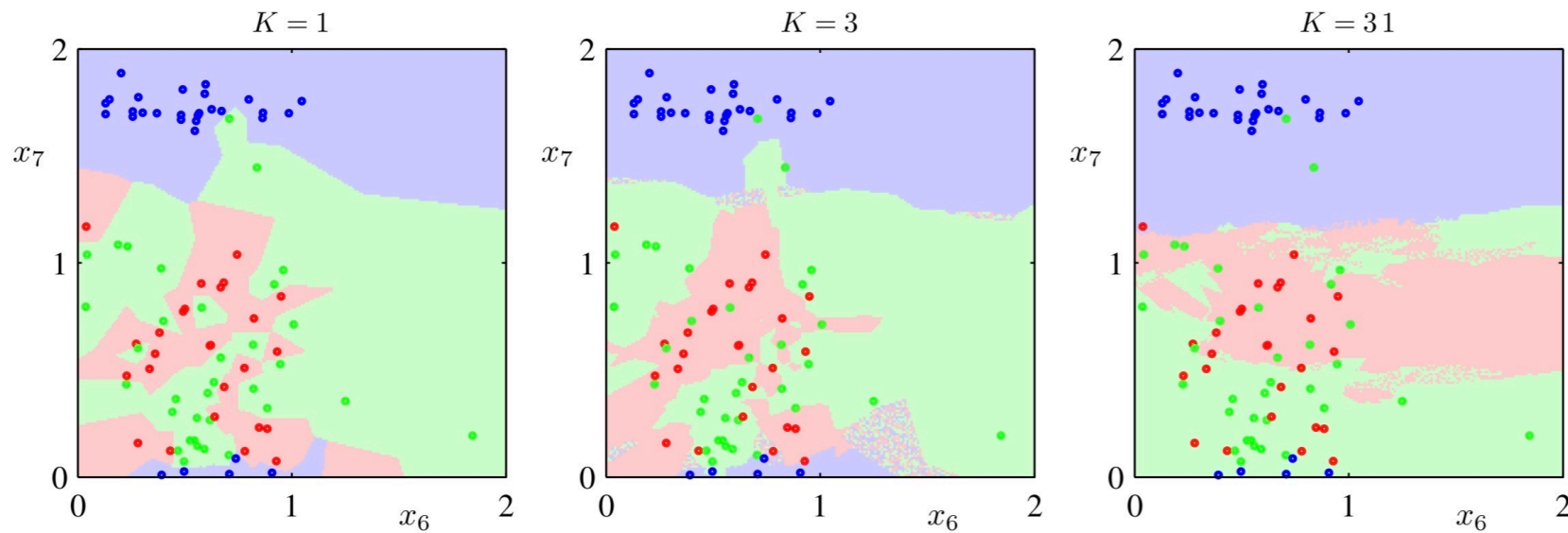
and class prior $P(C_k) = N_k / N$

Combining these equations and using Baye's theorem

$$P(C_k \mid x) = \frac{P(x \mid C_k) P(C_k)}{P(x)} = \frac{K_k}{K}$$

To minimize the probability of misclassification, assign $x$ to the class having the largest posterior probability $K_k/K$

The particular case of $K=1$ is called nearest-neighbour



**Figure 2.28** Plot of 200 data points from the oil data set showing values of $x_6$ plotted against $x_7$, where the red, green, and blue points correspond to the 'laminar', 'annular', and 'homogeneous' classes, respectively. Also shown are the classifications of the input space given by the $K$-nearest-neighbour algorithm for various values of $K$.

What happen if $K = N$.

# Remark

- K controls the degree of smoothing

- kNN and kernel density methods require the entire data set to stored