Chapter 2 Probability distributions Introduction to useful probability distributions < continuous Discuss key statistical concepts such as Bayesian inference Density estimation (find pc*)) independent and identically - lid assumption

- un supervised learning

- parametric vs non parametric Conjugate distributions (prior and posterior)

2.1 Discrete random variables concept, parameters Bernoulli, binomial, beta f distribution Bernoulli distribution R. V $X \in \{0, 1\}$ parameter M denote the prob. of X = 1pcx = 1 | M) = M - Prob.12420 parameter Bern $(x|\mu) = \mu^{x}(1-\mu)^{-x}$ (x < 10, 14)

Remark

-
$$p(x=0|M) + p(x=1|M) = (1-M) + M = [$$

- $\mathbb{E}[x] = 0 \cdot p(x=0|M) + 1 \cdot p(x=1|M) = M$
- $Var[x] = \mathbb{E}[x^{2}] - \mathbb{E}[x]^{2} = M - M^{2} = M(1-M)$
Bern(x(M)
Assume $fx_{1}, \dots x_{N}$ is drawn independently from some Bernoulli
Find a parameter $M (= p(x=1))$ in frequentist setting

See likelihood function $(\sigma f \mu) (D = 4x_{1,...} x_{n}t, x_{i} \in \{\circ, l\})$ $P(D|\mu) = \prod_{n=1}^{N} P(x_{n}|\mu) = \prod_{n=1}^{N} \mu^{x_{n}} (1-\mu)^{(1-x_{n})} (2.5)$

Estimate (find) a parameter M by maximizing (2.5)

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

ln PCDIM) only depends on
$$x_n (\Sigma x_n)$$
. Find a value for M
of $\frac{1}{dM} \ln PCD(M) = 0$

Chapter 2 Probability Distributions

Binomial distribution
$$Oliginary I$$

Let $x = 0$ or 1 and N be \approx trials.
R.v $m \in \{0, 1, ..., N\}$ to be \approx of $x = 1$.
From (2.5)
 m
binomial distribution ∞ $m^m(1-m)^{n-m}$
 $m = \approx$ of $x = 1$, M : the probability of $x = 1$

To normalize prob. dist. calculate all of possible
$$*$$

of obtaining m, $x = 1$. Denote by
Bin $(m | N, M) = \binom{N}{m} M^m (1-M)^{N-m}$
where $\binom{N}{m} := \frac{N!}{(N-m)! m!} = NCm$
We have
 $E[m] = \sum_{m=0}^{N} m Bin Cm | N, M)$
 $(III \ge x = 1 e! i! 4) = \sum_{m=0}^{N} m \binom{N}{m} M^m (1-M)^{N-m} = NM$

$$V_{\text{or}} [m] = \sum_{m=0}^{N} (m - E[m])^2 Bin (m | N, M) = N M (1-M)$$



Recall the
$$P(D|M)$$
 of Bernoulli distribution
likelihood $P(D|M) = \prod_{n=1}^{N} \mu^{x_n} (1-\mu)^{(1-x_n)} = x_n \in \{0, 1\}$
 $\mu^{x} (1-\mu)^{1-x}$ ($P(M) \ OH > 1 \ge 34 \ OP \ge 1>1 \ OCH \ge 01 \ M \ge 1 \ge 2 \ge 30 \ OH \ge 1>1$
To see Bayesian approach, we need to introduce $P(A)$.
beta distribution nomalization constant
 $P(A) = \frac{P(A+b)}{P(A+b)} \mu^{A-1} (1-\mu)^{b-1} \qquad (2.13)$
 $P(A) = \frac{P(A+b)}{P(A+b)} \mu^{A-1} (1-\mu)^{b-1} \qquad (2.13)$
where $P(x) := \int_{0}^{\infty} u^{x-1} e^{-u} du$ is the gamma function (1.414)

S Beta (Ma, b) dM = 1 Chapter 2 Probability Distributions

Remark

-
$$\int_{0}^{1} Beta(\mu | a, b) d\mu = 1$$

- $E[\mu] = \frac{a}{a+b}$, $Var[\mu] = \frac{ab}{(a+b)^{2}(a+b+1)}$ Excercise.
- a, b are called hyperparameters
The posterior dist. of μ has the form as
 $P(\mu | m, l, a, b) \propto \mu^{m+a-1} (1-\mu)^{l+b-1}$
where $l = N-M$. $M: \neq of x_{n}=1$

$$\Rightarrow posterior$$

$$\Rightarrow p(\mu | m, l, a, b) = \frac{\Gamma(m+a+l+b)}{\Gamma(m+a)} p(l+b) \qquad (1-\mu) \qquad (2.18)$$

a, b: parameters of prior) - > posterior m, l: result of observation

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

a, b can be interpreted as effective \times observations of x = 1 and x = 0.

Sequential approach (Bayesian view point)

Chapter 2 Probability Distributions



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.

Bern X,...XN observation

Let us predict the outcome of the next trial

$$\frac{p(x = 1 | D)}{p(x)} = \int_{0}^{1} p(x = 1, M | D) dM = \int_{0}^{1} \frac{p(x = 1, M, D)}{p(D)} dM$$

$$= \int_{0}^{1} \frac{p(x = 1, M, D)}{p(M, D)} \cdot \frac{p(M, D)}{p(D)} dM$$

$$= \int_{0}^{1} \frac{p(x = 1, M, D)}{p(x = 1, M, D)} \cdot \frac{p(M, D)}{p(D)} dM$$

$$= \int_{0}^{1} p(x = 1, M) \cdot \frac{p(M, D)}{p(M, D)} dM$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, D)} dM = \frac{p(A, B, C)}{p(B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, D)} dM = \frac{p(A, B, C)}{p(B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, D)} dM = \frac{p(A, B, C)}{p(A, B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, D)} dM = \frac{p(A, B, C)}{p(A, B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, D)} dM = \frac{p(A, B, C)}{p(A, B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, D)} dM = \frac{p(A, B, C)}{p(A, B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(M, C)} dM = \frac{p(A, B, C)}{p(A, B, C)}$$

$$= \int_{0}^{1} M \frac{p(M, D)}{p(A, C)} dM = \frac{p(A, B, C)}{p(A, C)}$$

Methaditer 2 Probability Distributions

$$PCx=|D| = \frac{m+a}{m+a+b+b} = \frac{m+a}{w+a+b}$$

$$(l=N-m)$$

a, b are hyperparameters of prior. \overline{N} m, l are from the result of experiment (M = 36 of x = 1)Remark

 $-m, l \rightarrow \infty$ (huge observations), then $ML \approx Bayesian$ $-a, b \rightarrow \infty$, then variance $\rightarrow 0$ (both prior and posterior)

2.2 1	Multinomia	Vor	lables				
Extend	Bernoulli	, binon	nial, beta	distri	butions		
Discrete	variable s	that	can take	on one	ट र्ज	K possible	case
1-of-K scheme C one hot encoding)							
	X orange apple apple grape	ota	nge apple 1 0 2 1 0 1 0 0	e grape O O I		$ \begin{array}{c} & & & & \\ & & (1,0,0) \\ & & (0,1,0) \\ & & (0,1,0) \\ & & (0,0,1) \end{array} $	
	ohange	2		0 		((, 0, 0)	F

Remark

- X can take K possible cases.

$$-\sum_{k=1}^{k} p(k) (M_{k}) = \sum_{k=1}^{k} M_{k} = 1$$

$$- \operatorname{EC} \times [M] = \sum_{k} PC \times (M) \times = (M_{1}, \dots, M_{k})^{T} = M$$



Multinomial distribution $M := (M_{1,...}, M_{k})^{T}$ Joint distribution of the quantities $M_{1,...}, M_{k}$ conditioned on M and on X of N total observations.

$$Mult Cm_1, m_2, \dots m_k | M_1, N) = \begin{pmatrix} N \\ m_1 m_2 \dots m_k \end{pmatrix} \frac{k}{\prod} M_k^{m_k} O \leq M_k \leq 1$$

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 洲 categorical 변名을 갖는 N 개 자료 에서 각 ド-class 가

Mr 低 가진 就是 (M) 子文の凝은 ccH)

2.2.1 Dirichlet distribution (multi-dim version of beta)
Consider the prior distributions for the parameters (Met
of multinomial distribution. (or categorical distribution)
Recall Mult (m_1,...
$$m_{E} | M, N$$
) or $\prod_{k=1}^{K} M_{E}^{m_{k}}$
 $\Rightarrow \qquad p(M | 0k) \circ \prod_{k=1}^{K} M_{E}^{m_{k}} \qquad Ma^{-1} \prod_{M_{1}} M^{-1}$
where $o \leq M_{E} \leq 1$, $\sum_{k=1}^{K} M_{E} = 1$. Here $o_{K} := (x_{1,...}, x_{E})^{T}$ is
the parameter
 $i M_{E}^{K}$ is confined to a simplex of dim $K-1$

Chapter 2 Probability Distributions

Dirichlet distribution

$$Dir(M(\alpha)) := \frac{\int C\alpha_0}{\int C\alpha_1 \cdots \int C\alpha_k} \frac{k}{\kappa} \frac{\alpha_k - 1}{M_k} \qquad o \leq M_k \leq 1$$

where
$$P(x)$$
 is the Gamma Sunction, $x_0 := \sum x_k$
So the posterior distribution for the parameters (Mri
P(MID, or) oc P(DIM) p(MIOr) oc $\prod_{k=1}^{k} M_k^{\alpha_k + M_k - 1}$
posterior likelihood
Posterior dist, again follow Dirichlet dist.

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

As binomial dist. with beta prior, we can interprete α_k of Dirichlet prior as an effective \ll of $\alpha_k = 1$.



2.3 The Gaussian distribution (a.k.a. normal dist.) Single real value $x \in IR$ $\mathcal{N}(x \mid \mathcal{M}, \sigma^2) := \frac{1}{(2\pi\sigma^2)^{k_2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mathcal{M})^2\right\}$

where M: mean, σ^2 : variance

$$D - \dim \text{ vector } & \in \mathbb{R}^{P}$$

$$\mathcal{N}(\mathcal{K}(\mathcal{M}, \Sigma) = \frac{1}{(2\pi)^{p_{2}}} \frac{1}{1\Sigma 1^{l_{2}}} \exp\left(1 - \frac{1}{2}(\mathcal{K} - \mathcal{M})^{T}\Sigma(\mathcal{K} - \mathcal{M})\right)$$

$$\text{IND DXD DX1}$$

where M: D-dim mean vector, $\Sigma: D \times D$ covariance matrix $|\Sigma|:$ determinant of Σ

- Remark: Why Graussian is important?
- Fits many natural phenomena
- Maximum entropy in continuous r.v
- Central limit theorem



Fix N ijy random samples of vector $X_1, X_2, ..., X_N$ are population with MI, D. fram from a $\overline{\mathbb{X}}_{N} := \frac{1}{N} \sum_{n=1}^{N} \mathbb{X}_{n}$ ん R.V. > $\mathcal{N}(\mathcal{M}, \frac{1}{\mathcal{N}}\Sigma)$ 3 N = 1N = 2N = 10Sample 2 2 2 mean 0.5 0.5 0.5 0

Figure 2.6 Histogram plot $Chapter 2^{n}$ Figure 2.6 Histogram plot $Chapter 2^{n}$ Histogram plot N. We observe that as N increases, the distribution tends towards a Gaussian.

WLOGT, assume I is symmetric (and real).

Consider the eigenvector equation of Z

$$\sum u_i = \pi_i u_i$$
 $j = 1 \dots D$

Eigenvalues
$$7_{1}$$
, 7_{p} are real and its eigenvectors can be
chosen to form an orthonormal set, so that
 $u_{l_{x}}^{T} u_{l_{y}} = I_{\lambda_{y}}$
where $I_{\lambda_{y}}$ is the $\lambda_{r, y}$ element of the identity matrix.
By eigenvector decomposition, Σ can be expressed as
 $\Sigma = \sum_{\substack{z=1\\z=1}^{p} 7_{\lambda_{z}} u_{\lambda_{z}} u_{\lambda_{z}}^{T}$ $DxI (xD)$
 DxD vector
and $\Sigma^{-1} = \sum_{\substack{z=1\\z=1}^{p} \frac{1}{7_{\lambda_{z}}} u_{\lambda_{z}}^{T} u_{\lambda_{z}}^{T}$ (inverse)

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 λ_i .





- If 7i > 0 $\forall i = 1..., D$, contour surface of Δ is ellipsoid. - center M, axes oriented along u_i and scalling factors are given by $7i^{\frac{1}{2}}$

WLOGT, assume all eigenvalues of
$$\Sigma$$
 are strictly positive
Otherwise the distribution cannot be normalized (see ch 12)
i.e. Σ is assume to be positive definite.
 $y := U(x - \mu)$
Now consider the Gaussian dist. in y_i coordinate system.
Jacobian matrix J with
 $J_{ij} = \frac{\partial \chi_i}{\partial y_j} = U_{ij}$
where U_{ji} are element of U^T
 $U^{y} + \mu I = x$

$$\int [J]^{2} = [U^{T}]^{2} = [U^{T}]U^{T} = [U^{T}]U$$

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

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

has
$$E[X] = M$$
, $Var[X] = \sigma^2$

Now we will interpret parameters
$$\mu$$
, Σ .
 $E[X] = \frac{1}{(2\pi)^{3}p} \frac{1}{1\Sigma 1^{\frac{1}{2}}} \int exp \left\{ -\frac{1}{2} (X - \mu)^{T} \Sigma^{-1} (X - \mu) \right\} \times dX$

$$= \frac{1}{(2\pi)^{2}} \frac{1}{|\Sigma|^{2}} \int \exp\left\{-\frac{1}{2} \frac{2}{2} \sum_{i=1}^{n} \frac{1}{2} \left(\frac{2}{2} \pm \mu \right) \frac{1}{2} \frac{1}{2} \frac{1}{2} \left(\frac{2}{2} \pm \mu \right) \frac{1}{2} \frac$$

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

Note that the exponent is even. So the term Z (~+,M) will vanish. in the factor $EC \times J = M$. \Rightarrow XGIR Now consider the second order moments of multivariate Graussian. In univariate case, the second order moment is given by ECX2]. In multivariate Graussian, there are D² second order moments given by ECX; X;]

 $\mathscr{K} = \begin{pmatrix} \mathfrak{X}_{i} \\ \vdots \\ \mathfrak{X}_{D} \end{pmatrix} \mathbb{E} \begin{bmatrix} \mathfrak{X}_{i} \mathfrak{X}_{j} \end{bmatrix}$

$$E[\Sigma \# \#^{T}] = \frac{1}{(2\pi)^{9_{2}}} \frac{1}{1\Sigma|^{1_{2}}} \int \exp\left\{-\frac{1}{2}(\#-M)^{T}\Sigma^{-1}(\#-M)\left[\# \#^{T}d\#\right]\right\}$$

$$= \frac{1}{(2\pi)^{9_{2}}} \frac{1}{1\Sigma|^{1_{2}}} \int \exp\left\{-\frac{1}{2}\frac{1}{2T}\Sigma^{-1}\frac{1}{2T}\right\} \frac{2}{(2T+M)(2T+M)^{T}} \frac{2}{2T}$$

$$= \frac{1}{(2\pi)^{9_{2}}} \frac{1}{1\Sigma|^{1_{2}}} \int \exp\left\{-\frac{1}{2}\frac{1}{2T}\Sigma^{-1}\frac{1}{2T}\right\} \frac{2}{(2T+M)(2T+M)^{T}} \frac{2}{2T}$$

$$\Rightarrow = U(2\#-M), \text{ rows of } U \text{ are eigenvectors of } \Sigma$$

$$\Rightarrow = U^{T}\Psi = (u_{1}...u_{p})\binom{Y_{1}}{1} = \sum_{i=1}^{p} Y_{i} u_{i}$$

Chapter 2 Probability Distributions
$$\frac{1}{(2\pi)^{\vartheta_{2}}} \frac{1}{|\Sigma|^{\vartheta_{2}}} \int \exp\left\{-\frac{1}{2} \overset{2}{\approx}^{T} \Sigma^{T} \overset{2}{\approx}\right\} \overset{2}{\approx} \overset{2}{\approx}^{T} d \overset{2}{\approx}$$

$$= \frac{1}{(2\pi)^{\vartheta_{2}}} \frac{1}{|\Sigma|^{\vartheta_{2}}} \sum_{\tilde{k}=1}^{\tilde{p}} \sum_{j=1}^{\tilde{p}} u_{j} u_{j}^{T} \int \exp\left\{-\frac{p}{k} \frac{y_{k}^{2}}{2\pi_{k}}\right\} y_{\lambda} y_{j} d y$$

$$= \sum_{\tilde{\lambda}=1}^{\tilde{p}} u_{1_{\tilde{\lambda}}} u_{\lambda}^{T} \pi_{\tilde{\lambda}} = \sum$$
we have used $|\Sigma| = \prod_{\tilde{\lambda}=1}^{\tilde{p}} \pi_{\tilde{\lambda}}$ and $\frac{1}{(2\pi)^{\vartheta_{2}}} \frac{1}{\pi^{\vartheta_{2}}} \exp\left\{-\frac{y^{2}}{2\pi_{k}}\right\} \sim N(0, \pi)$
e.g. $p=2$

$$\sum_{\tilde{\lambda}=j}^{\tilde{\lambda}} \iint \exp\left\{-\frac{y_{1}^{2}}{2\pi_{k}}\right\} \exp\left\{-\frac{y_{2}^{2}}{2\pi_{k}}\right\} y_{\lambda} y_{j} d y_{j} d y_{j}$$

will vonish Chapter Probability Distributions

Thus we obtain
$$\mathbb{E}[\mathbb{X} \times \mathbb{F}^T] = \mathbb{M} \times \mathbb{M}^T + \Sigma$$
 (DXD matrix)
and covariance of \mathbb{X} can be obtained by
 $\mathbb{C} \times \mathbb{E}[\mathbb{X}] = \mathbb{E}[\mathbb{X} - \mathbb{E}[\mathbb{X}])(\mathbb{X} - \mathbb{E}[\mathbb{X}])^T]$
 $= \Sigma$

RemarkDD
$$\Sigma D$$
Sym. real- $\overset{\times}{\times}$ of parameters: $\frac{D(D+3)}{2}$ quadratic- $\Sigma = \text{diag}(\sigma_{\lambda}^{2})$ or $\Sigma = \sigma^{2}I$ deep dive2DD+1into Gaussian

- Unimodal (single maximum)

Figure 2.8 Contours of constant x_2 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.



2.3.1 Conditional Gaussian distributions

$$X : D-dimensional vector with NCX(M, \Sigma)$$
 which is
partitioned into X_a , X_b with $X_a \in IR^M$, $X_b \in IR^{D-M}$
 $X = \begin{pmatrix} X_a \\ X_b \end{pmatrix}$

 Let $\Lambda := \Sigma^{-1}$. Inverse of covariance matrix, precision matrix $M \times M$ $\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} P - M \times P - M$ Note that Λ_{aa} , Λ_{bb} are symmetric, $\Lambda_{ba} = \Lambda_{ab}^{T}$. $\Lambda_{aa} \neq \Sigma_{aa}^{-1}$

Find conditional distribution PC XalXb). Fix Xb. Consider the quadratic form in the exponent.

$$-\frac{1}{2}(\cancel{x}-\cancel{\mu})^{T}\Sigma^{T}(\cancel{x}-\cancel{\mu})) \xrightarrow{[x P]} (\cancel{x} P \xrightarrow{P\times 1} (\cancel{x} F \xrightarrow{fixed}))$$

$$= -\frac{1}{2}(\cancel{x}a-\cancel{\mu}a)^{T}\wedge_{aa}(\cancel{x}a-\cancel{\mu}a) - \frac{1}{2}(\cancel{x}a-\cancel{\mu}a)^{T}\wedge_{ab}(\cancel{x}b-\cancel{\mu}b))$$

$$-\frac{1}{2}(\cancel{x}b-\cancel{\mu}b)^{T}\wedge_{ba}(\cancel{x}a-\cancel{\mu}a) - \frac{1}{2}(\cancel{x}b-\cancel{\mu}b)^{T}\wedge_{bb}(\cancel{x}b-\cancel{\mu}b))$$

$$(2.70)$$

First, P(XalXb) will be M-dim Graussian, because density Sunction is a quadratic form of exponent. Now we are going to find its mean vector and covariance (M D-4) by "completeting the square" $(M | P-H) \begin{pmatrix} O & O \\ O & O \end{pmatrix} \begin{pmatrix} M \\ \overline{P-H} \end{pmatrix}$

Chapter 2 Probability Distributions

DXD

DXI

42

E.g. in case pcz) of
$$\exp(ax^2 + bx + c) \rightarrow x \ge gaussian$$

$$\Rightarrow \qquad pcx) of exp \left\{ a\left(x^2 + \frac{b}{a}x + \frac{b^2}{4a^2}\right) - \frac{b^2}{4a} + c \right\}$$

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

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

(2) Since
$$N(x|\mu,\sigma^2) = O(exp(-\frac{1}{2\sigma^2}(x^2-2\mu x+\mu^2)))$$

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

 $\Rightarrow \quad a = -\frac{1}{26^2}, \quad b = \frac{M}{6^2}, \quad \mu = -\frac{b}{2a}, \quad \sigma^2 = -\frac{1}{2a}, \quad Chapter 2 Probability Distributions$

43

Likewise, the exponent in
$$D - \dim$$
 Gaussian can be written
 $-\frac{1}{2}(* - \mu I)^T \Sigma^{-1}(* - \mu I) = (-\frac{1}{2})*^T \Sigma^{-1} * + *^T \Sigma^{-1} \mu I + \text{constant}$

In view of (2.70), (second order in *a)

So we obtain covariance of
$$P(\text{#al#b})$$
 is given by

$$\Sigma_{alb} := \Lambda_{aa}^{-l} (\neq \Sigma_{aa})$$

$$\bigwedge_{alb}^{m} ?$$

Now consider all of the terms in (2.70) that are linear in X_a

$$\mathcal{X}_{a}^{T} \left\{ \Lambda_{aa} \mathcal{M}_{a}^{I} - \Lambda_{ab} \left(\mathcal{X}_{b} - \mathcal{M}_{b} \right) \right\}$$

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

Since

$$\Sigma_{alb}^{-1} \mathcal{M}_{alb} = \Lambda_{aa} \mathcal{M}_{a} - \Lambda_{ab} (\mathcal{K}_{b} - \mathcal{M}_{b}),$$

$$\mathcal{M}_{alb}^{-1} = \Sigma_{alb} \left\{ \Lambda_{aa} \mathcal{M}_{a} - \Lambda_{ab} (\mathcal{K}_{b} - \mathcal{M}_{b}) \right\}$$

$$= \mathcal{M}_{a} - \Lambda_{a}^{-1} \Lambda_{ab} (\mathcal{K}_{b} - \mathcal{M}_{b})$$

Chapter 2 Probability Distributions

.

$$\therefore \quad M_{alb} = M_a - \Lambda_{aa}^{-1} \Lambda_{ab} (X_b - M_b) \qquad \Sigma_{alb} = \Lambda_{aa}^{-1}$$

Let us find have and hab.

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

Use the following identification for the inverse of a partitioned matrix $\begin{pmatrix} A B \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1}FCMB^{-1} \end{pmatrix}$ (2.06)

where M:= (A - Chapter 2 Probability Distributions

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

Remark

2.3.2 Marginal Gaussian distributions
Consider the following marginal distribution
$$\mathcal{N}(\mathcal{X}|\mathcal{M}, \Sigma)$$

 $P(\mathcal{X}_{a}) = \int P(\mathcal{X}_{a}, \mathcal{X}_{b}) d\mathcal{X}_{b} \qquad \mathcal{X} = \begin{pmatrix} \mathcal{X}_{a} \\ \mathcal{X}_{b} \end{pmatrix}^{M}$
 $(\mathfrak{X}_{b} \ge \mathfrak{A} \not \bowtie \mathfrak{b} \mathfrak{M} \mathcal{A} \quad \mathfrak{G} \not \bowtie \not \bowtie \mathfrak{G} \not \varkappa \mathfrak{G} \not \bowtie \mathfrak{G} \not \varkappa \mathfrak{G} \not \bowtie \mathfrak{G} \not \varkappa \mathfrak{G} \not \bowtie \mathfrak{G} \not \varkappa \mathfrak{G} \varkappa \mathfrak{G} \not \varkappa \mathfrak{G} \varkappa \mathfrak{G}$

In order to integrate out \$\$, pick out those terms involving *b $-\frac{1}{2} \times_{b}^{T} \wedge_{bb} \times_{b} + \times_{b}^{T} M = -\frac{1}{2} (\times_{b} - \Lambda_{bb}^{-1} m)^{T} \wedge_{bb} (\times_{b} - \Lambda_{bb}^{-1} m)$ (2.84) $+\frac{1}{2}$ m^T \wedge_{bb} m (square expression) intep. of **b where $m_{1} := \Lambda_{bb} M_{b} - \Lambda_{ba} (X_{a} - M_{a})$

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

(2.86) $\int \exp \left\{ -\frac{1}{2} \left(\chi_{b} - \Lambda_{bb}^{-1} m_{l} \right)^{T} \Lambda_{bb} \left(\chi_{b} - \Lambda_{bb}^{-1} m_{l} \right) \right\} d\chi_{b}$ Chapter 2 Probability Distributions

which is an inverse of the normalization coefficient. As seen before, this coefficient is independent of mean. Combining the last term ($\frac{1}{2}$ m^T \wedge_{bb} m) in (2.84) with remaining. terms in (2.70) depending on *a, we obtain - I MIT Abb MI - I Xa Aaa Xa + Xa (Aaa Ma + Aab Mb) + constant $= -\frac{1}{2} \left[\Lambda_{bb} M_{b} - \Lambda_{ba} (X_{a} - M_{a}) \right]^{T} \Lambda_{bb} \left[\Lambda_{bb} M_{b} - \Lambda_{ba} (X_{a} - M_{a}) \right]$ + $\mathcal{K}_{a}^{T} (\Lambda_{aa} \mathcal{M}_{a} + \Lambda_{ab} \mathcal{M}_{b}) + constant$ $= -\frac{1}{2} \#_{a}^{T} \left(\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} \right) \#_{a} + \#_{a}^{T} \left(\Lambda_{aa}^{-1} \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} \right) \#_{a}$ + constant **Chapter 2 Probability Distributions** 50

Recall the exponent in
$$D - \dim$$
 Graussian can be written
 $-\frac{1}{2}(x - \mu i)^{T} \Sigma^{-1}(x - \mu i) = -\frac{1}{2}x^{T} \Sigma^{-1}x + x^{T} \Sigma^{-1} \mu i + \text{constant}$
Denote the covariance of $p(x_{a})$ by Σ_{a} and Σ_{a} is
given by $\Sigma_{a} = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb} \Lambda_{ba})^{-1}$
 $\Sigma_{a} = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb} \Lambda_{ba})^{-1}$
Similarly, mean vector is given by
 $\Sigma_{a} (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb} \Lambda_{ba}) M_{a} = M_{a}$
 Σ_{a}^{-1}

Chapter 2 Probability Distributions

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

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

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

matrix $\sum_{a} = \left(\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb} \Lambda_{ba} \right)^{-1} = \sum_{aa}$ $\mathbb{E}[\mathbb{X}_{a}] = M_{a}, \quad \operatorname{cov}[\mathbb{X}_{a}] = \mathbb{Z}_{aa}$ Thus we have

 $P(\mathcal{H}_{\alpha}) = \int P(\mathcal{H}_{\alpha}) P(\mathcal{H}_{\alpha})$ where

$$\mathcal{N}(\mathcal{K}|\mathcal{M}, \Sigma) \quad \text{with} \quad \Lambda := \Sigma^{-1} \qquad \mathcal{D}-\dim \quad \mathcal{K}$$

$$\mathcal{K} = \begin{pmatrix} \mathcal{K}_{a} \\ \mathcal{K}_{b} \end{pmatrix} \qquad \mathcal{M} = \begin{pmatrix} \mathcal{M}_{a} \\ \mathcal{M}_{b} \end{pmatrix} \qquad \Sigma = \begin{pmatrix} \Sigma_{aa} \quad \Sigma_{ab} \\ \Sigma_{ba} \quad \Sigma_{bb} \end{pmatrix} \qquad \Lambda = \begin{pmatrix} \Lambda_{aa} \quad \Lambda_{ab} \\ \Lambda_{ba} \quad \Lambda_{bb} \end{pmatrix}$$

Conditioned distribution

$$P(\mathcal{X}_{a}|\mathcal{X}_{b}) = \mathcal{N}(\mathcal{X}_{a}|\mathcal{M}_{a|b}, \Lambda_{aa})$$
$$\mathcal{M}_{a|b} = \mathcal{M}_{a} - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathcal{X}_{b} - \mathcal{M}_{b})$$

Marginal distribution

$$p(x_{\alpha}) = N(x_{\alpha}|M_{\alpha}, \Sigma_{\alpha\alpha})$$



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 Graussian model example Gaussian marginal dist. pcx) Gaussian conditional dist pcy1x) pcylx) has a mean as a linear function of x and a covariance which is independent of X. *G M-dim i.e. $P(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathcal{M}, \Lambda^{-1})$ YE D-dim $p(\gamma|\chi) = \mathcal{N}(\gamma|A \times +b, L')$ where MI, A and b are parameters governing the means, and Λ and L are precision matrices.

We will find p(y) and p(x(1y)).
manginal conditional
$$P(x(1y))$$
.
Let $z_{1} = \begin{pmatrix} x \\ y \end{pmatrix}$ and us consider the joint prob. dist
 $p(z_{1}) = p(x_{1}, y) = p(y(x_{1})p(y)$
 $p(z_{1}) = p(x_{1}, y) = p(x_{1}, y)$
 $p(z_{1}) = p(x_{1}, y)$
 $p($

Consider the second term in (2.102)

$$-\frac{1}{2} \times^{T} (\Lambda + A^{T} L A) \times -\frac{1}{2} \times^{T} L \times + \frac{1}{2} \times^{T} L \times + \frac{1}{2} \times^{T} A^{T} L \times$$

$$= -\frac{1}{2} \begin{pmatrix} * \\ * \end{pmatrix}^{T} \begin{pmatrix} \Lambda + A^{T} L A \\ -L A \end{pmatrix} \begin{pmatrix} -A^{T} L \end{pmatrix} \begin{pmatrix} * \\ * \end{pmatrix} = -\frac{1}{2} \overleftarrow{z}^{T} R \overleftarrow{z}$$

 \Rightarrow has precision (inverse of covariance) matrix given by $\frac{1}{2} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 + A^T L A - A^T L \\ -L A L \end{pmatrix}$

$$\Rightarrow \text{ COV } \square \exists = R^{-1} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^{T} \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1} A^{T} \end{pmatrix} \quad (2.105)$$

Consider the linear term in (2.102)

 $\widehat{}$

$$\times^{\mathsf{T}} \wedge \mathfrak{M} - \times^{\mathsf{T}} A^{\mathsf{T}} \sqcup \mathbb{B} + \mathbb{Y}^{\mathsf{T}} \sqcup \mathbb{B} = \begin{pmatrix} \times \\ & \end{pmatrix}^{\mathsf{T}} \begin{pmatrix} \wedge \mathfrak{M} - A^{\mathsf{T}} \sqcup \mathbb{B} \\ & \sqcup \mathbb{B} \end{pmatrix}$$

$$\mathbb{E}[\mathcal{F}] = R^{-1} \begin{pmatrix} \lambda \mu - A^{T} L B \end{pmatrix}$$

$$\mathcal{L} = \begin{pmatrix} \mathscr{K} \\ \mathscr{Y} \end{pmatrix} = \begin{pmatrix} \mathscr{M} \\ A \mathscr{M} + \mathscr{W} \end{pmatrix} \qquad (2.108)$$

$$\mathcal{L} = \begin{pmatrix} \mathscr{K} \\ A \mathscr{M} + \mathscr{W} \end{pmatrix} \qquad (2.108)$$

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

$$E[Y] = AM + Ib$$

$$cov [Y] = L' + AA' A^{T}$$

Now we can find an expression for $PC \times (Y)$. $EC \times (Y) = (\Lambda + A^{T}LA)^{-1} \{A^{T}L(Y-B) + \Lambda \mu \}$ $COV [\times (Y) = (\Lambda + A^{T}LA)^{-1}$

Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
 (2.113)

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}+\mathbf{b},\mathbf{L}^{-1})$$
 (2.114)

the marginal distribution of y and the conditional distribution of x given y are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
 (2.115)

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

where

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

2.3.4 Maximum likelihood for the Gaussian
Data set
$$\chi = (\chi_{1,...} \chi_{N})^{T}$$
, $\{\chi_{n}\}$ iii samples of D-dimensional
Gaussian. The log likelihood function is given by
 $\ln p(\chi_{1/M}, \Sigma) \qquad \qquad \chi \qquad N \times D \qquad motrix$
 $= -\frac{NP}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^{N} (\chi_{n} - M)^{T} \Sigma^{-1} (\chi_{n} - M)$
Note that likelihood function depends only on the following two
quantities $\sum_{n=1}^{N} \chi_{n} \qquad \sum_{n=1}^{N} \chi_{n} \qquad \chi \qquad X_{n} \qquad X_{n}$

These are known as sufficient statistics for Goussian

$$\nabla_{\mu 1} \ln p(X | \mu, \Sigma) = \sum_{n=1}^{N} \Sigma^{-1} (X_n - \mu) \qquad P - dim vector$$
Set this gradient to zero vector, we obtain

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^{N} X_n \qquad \text{solution of maximum} \\
\qquad \text{(sample mean)} \qquad \text{MLE}$$

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

$$(\text{somple covariance})$$

Remark

- \sum_{ML} involves M_{ML} - M_{ML} is independent of \sum_{ML} F(F) = F(F) = F(F)Evaluate the expectations of thes solutions under the true distribution. Then we obtain

 $\mathbb{E}\mathbb{C}\mathcal{M}_{ML}] = \mathcal{M} \qquad \text{unbiased estimate}$ $\mathbb{E}\mathbb{C}\mathbb{Z}_{ML}] = \frac{N-1}{N}\mathbb{Z}$ $\int_{\mathcal{M}} \mathbb{E}\mathbb{C}\mathbb{C}_{ML} = \frac{N-1}{N}\mathbb{Z}$

2.3.5 Sequential estimation Sequential estimation for maximum likelihoot This method allows data points to be proceed one at time and then discarded and one important for on-line applications $\mathcal{M}_{ML} = \frac{1}{N} \sum_{n=1}^{N} \mathcal{K}_{n}$ Consider

which we will tenote by MML baset on N observations



Assume the conditional variance of
$$2$$
 is finite
 $E[(2-5)[0] < \infty$
 ∂^{*} solution
and wlog $f(0) > 0$ for $0 > 0^{*}$ and $f(0) < 0$ for $0 < 0^{*}$
A sequence of successive estimates of the root 0^{*} given by
 $\partial^{(W)} := \partial^{(W-1)} - \alpha_{N-1} \ge (\partial^{(W-1)})$
(2.129)
where $\ge (\partial^{(W)})$ is an observed value of \ge when $0 = \partial^{(W)}$

$$\begin{cases} Q_{N} i \ \text{represents a seq of positive numbers satisfying}} \\ Q_{im} Q_{N} = 0 \\ \sum_{N=0}^{\infty} Q_{N} = \infty \\ \sum_{N=1}^{\infty} Q_{N}^{2} < \infty \\ \sum_{N=1}^{\infty} Q_{N}^{2} < \infty \\ \text{of } f(0) \\ \text{By [Robbins - Monro], (2.129) converges to the root with probability one.} \\ \text{Remark} \end{cases}$$

- Third condition ensures that the accumulated noise has finite variance and Chapter 2 Probability Distributions spoil convergence. 68

General Maximum litelihood problem

$$f(\theta) = \int \frac{1}{2} p(2\theta) d\theta$$
By definition of O_{ML} , O_{ML} satisfies

$$\frac{\partial}{\partial \theta} \left\{ -\frac{1}{N} \sum_{n=1}^{N} e_n p(x_n | \theta) \right\} = 0$$

$$F_x \left[\frac{\partial}{\partial \theta} e_n p(x | \theta) \right]$$
Taking $N \to \infty$ and exchanging derivative and summation,

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

$$\mathbb{E}_x \left[\frac{\partial}{\partial \theta} e_n p(x | \theta) \right]$$

$$\mathbb{E}_x \left[\frac{\partial}{\partial \theta} e_n p(x | \theta) \right]$$
I.e. find the root of a regression function

Apply Robbins - Monro procedure $Q^{(N)} := Q^{(N-1)} - Q_{N-1} \frac{\partial}{\partial Q^{(N-1)}} \left[-l_n P^{(X_N)} Q^{(N-1)} \right]$ (2.135) Specific example: sequential estimation of the mean of Gaussian distribution In this case of is the MML mean of the Graussian and 2 is given by (2.136) $z = \frac{\partial}{\partial M_{ML}} \ln p(x|M_{ML},\sigma^2) = -\frac{1}{\sigma^2} (x - M_{ML})$

Substituting (2.136) into (2.135) with $\alpha_{N} = \frac{\sigma^{2}}{N-1}$

then we obtain (2.126)

2.3.6 Bayesian inference for the Gaussian
MLE method gave point estimates for
$$M$$
, Σ (section 2.3.4)
Now develop a Bayesian treatment
Single Gaussian random variable (Σ) Suppose σ^2 is known
Aim to inference M given N observations $X = \{X_1, ..., X_N\}$
The likelihood function is given by
 $p(X|M) = \prod_{n=1}^{N} p(X_n|M) = \frac{1}{(2\pi \sigma^2)^N \Sigma} \exp\{1 - \frac{1}{2\sigma^2} \sum_{n=1}^{N} (X_n - M)^2\}$
Note that this function is the form of the exponential of
a quadratic form of Chapter 2 Probability Distributions 72
We will choose a prior
$$p(\mu)$$
 given by Graussian because
the product of two exponentials of quadratic function of
 μ will also be Graussian
Take prior prob. $p(\mu)$ to be
 $p(\mu) = N(\mu | M_0, \sigma_0^2)$
 M_0, σ_0^2 hyperparameters

.

Posterior

Exercise 2.38. we obtain posterior parameters

$$P(\mathcal{M}|\mathcal{K}) = \mathcal{N}(\mathcal{M}|\mathcal{M}_{N}, \sigma_{N}^{2}) \qquad \qquad \frac{1}{N} Z x_{n}$$
where $\mathcal{M}_{N} = \frac{\sigma^{2}}{\mathcal{N}\sigma_{0}^{2} + \sigma^{2}} \mathcal{M}_{0} + \frac{\mathcal{N}\sigma_{0}^{2}}{\mathcal{N}\sigma_{0}^{2} + \sigma^{2}} \mathcal{M}_{ML}$

$$\frac{1}{\sigma_{N}^{2}} = \frac{1}{\sigma_{0}^{2}} + \frac{\mathcal{N}}{\sigma^{2}}$$
Remark
$$- \mathcal{M}_{N} \text{ is a compromise between } \mathcal{M}_{0} \text{ and } \mathcal{M}_{ML}$$

$$- \text{Effect of change in value } \mathcal{N}$$

- Precision is additive, if
$$N \rightarrow \infty$$
, $\sigma_N^2 \rightarrow 0$
- When N is finite, if $\sigma_n^2 \rightarrow \infty$, then the poterior mean reduces to M_{ML} and variance σ_N^2 becomes $\frac{\sigma_n^2}{N}$



Now we wish to infer the variance and assume mean is known.
Let the precision
$$\pi := \sqrt{5^2}$$
. The likelihood function of π
 $p(x|\pi) = \prod_{n=1}^{\infty} N (x_n|\mu, \pi^n) \circ C = \pi^{N_2} \exp\left(-\frac{\pi}{2} \sum_{n=1}^{\infty} (x - \mu)^2\right)^n$
i.e. the form of $\pi^{N_2} \cdot \exp\left(-\frac{\pi}{2}\right)$
The corresponds to gamma distribution
 $Gam(\pi|\alpha,b) = \frac{1}{\Gamma(\alpha)} b^{\alpha} \pi^{\alpha-1} \exp(-b\pi), \quad \pi > 0$
Here $\Gamma(\alpha)$ is a gamma function $\Gamma(\alpha) = \int_{0}^{\infty} u^{\alpha-1} e^{-u} du$
 $Chapter 2 Probability Distributions$

Remark



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

Consider the prior dist. Gram
$$(7 \mid a_0, b_0)$$
. $(a_0, b_0: hyperportanteter)$
The posterior dist. of 7 is as below
 $p(7 \mid x) \propto \frac{\pi^{N_2}}{1 \text{ likelihood function of 7}} \stackrel{?}{=} \frac{c_1}{2} \sum_{n=1}^{\infty} (x_n - \mu)^2 \frac{1}{r} \cdot \frac{c_1}{r} (x_n - \mu)^2 \frac{1}{r} \cdot \frac{c_1}{r} (x_n - \mu)^2 \frac{1}{r} (x_n - \mu)^2 \frac{1}{r$

Remark

- Effect of observing N data points - increases the value of a by $\frac{7}{2}$ b by $\frac{N}{2} O_{ML}^{2}$ 11 - We can interprete the parameter as in terms र्ण 200 'effective' prior observations. $- \operatorname{E[7]} = \frac{a_N}{b_N} = \frac{2a_0 + N}{2b_0 + N \sigma_{ML}^2}$ 3= --- $\operatorname{Vor} [7] [\%] = \frac{c_{W}}{b_{1}^{2}}$ $\mathcal{E}[\mathcal{N}] \longrightarrow$

Now suppose that the both
$$\mu$$
 and $\overline{\lambda}$ are unknown
Consider the dependence of the likelihood function on μ and $\overline{\lambda}$
 $p(x \mid \mu, \overline{\lambda}) = \prod_{n=1}^{N} \left(\frac{\lambda}{2x}\right)^{k_{2}} \exp\left\{-\frac{\overline{\lambda}}{2}(x_{n}-\mu)^{2}\right\}^{n}$
 $\overline{\lambda}, \mu$ $C\left[\overline{\lambda}^{k_{2}} \exp\left(-\frac{\overline{\lambda}\mu^{2}}{2}\right)\right]^{N} \exp\left[\overline{\lambda}\mu \sum_{n=1}^{N} x_{n} - \frac{\overline{\lambda}}{2} \sum_{n=1}^{N} x_{n}^{2}\right]$
Thus the prior distribution should take the form
 $p(\mu, \overline{\lambda}) \circ C\left[\overline{\lambda}^{k_{2}} \exp\left(-\frac{\overline{\lambda}\mu^{2}}{2}\right)\right]^{\beta} \exp\left\{C\overline{\lambda}\mu - d\overline{\lambda}\right\}$
 $= \exp\left\{-\frac{\beta\overline{\lambda}}{2}(\mu - \zeta_{\beta})^{2}\right\} \overline{\lambda}^{\beta/2} \exp\left\{-\left(d - \frac{C^{2}}{2\beta}\right)\overline{\lambda}\right\}$
(2.153)
Manager 2 Probability Distributions

where
$$c, d$$
 and β are constants. Use $p(\mu, \pi) = p(\mu|\pi) p(\pi)$.
 $p(\mu|\pi)$: a Graussian whose precision is a linear sunction of π
 $p(\pi)$: a gamma distribution. So we take a prior
 $p(\mu, \pi) = \mathcal{N}(\mu|\mu_0, (\beta\pi)^{-1})$ Gram $(\pi|\alpha, b)$ (2.154)
where $\mathcal{M}_0 := \frac{c}{\beta}$, $\alpha := \frac{(1+\beta)}{2}$, $b := d - \frac{c^2}{\beta p}$
(2.154) is called normal gamma or Graussian gamma.
Note that it is not the simply the product of an independent
Graussian prior and gamma prior.

Multivariate Graussian
$$\mathcal{N}(\mathcal{X}|\mathcal{M}, \Lambda^{-1})$$
 for D -dim \mathcal{X}
First, when precision matrix Λ is known, the conjugate prior
distribution is again a Graussian.
Second, for known mean and unknown precision matrix Λ ,
the conjugate prior distribution is the Wishart distribution
given by trace of matrix
 $\mathcal{W}(\Lambda|\mathcal{W},\mathcal{P}) = B|\Lambda|^{(\mathcal{V}-\mathcal{P}-1)/2} \exp\left(-\frac{1}{2}\operatorname{Tr}(\mathcal{W}^{-1}\Lambda)\right)$
where \mathcal{V} is called the number of degrees of freedom

If both the mean and precision are unknown, the conjugate prior is given by $P(\mathcal{M}, \Lambda \mid \mathcal{M}_0, \beta, w, \nu) = \mathcal{N}(\mathcal{M} \mid \mathcal{M}_0, (\beta \Lambda)^{-1}) \mathcal{W}(\Lambda \mid w, \nu)$

which is known as the probability Distributions. It or Groussian - Wishartes

2.3.7 Student's t-distribution Conjugate prior for the precision of a Gaussian is given by a gamma distribution. Consider univariate Gaussian NCXIM, 2") with Gamma prior Gram (C/a, b). Integrate out the precision $P(X|M,a,b) = \int_{a}^{b} N(X|M,t^{-1}) Gram(t|a,b) dt$ $= \int_{0}^{\infty} \frac{b^{\alpha} c^{-b^{2}} z^{\alpha-1}}{(2\pi)^{2}} \left(\frac{z}{2\pi}\right)^{2} \exp\left\{-\frac{z}{2}(2\pi-M)^{2}\right)^{2} dz$ $= \frac{b}{\Gamma(\alpha)} \left(\frac{1}{2\pi}\right)^{2} \left[b + \frac{(x-M)^{2}}{Chapter 2 Probability Distributions}\right]^{-\alpha-1/2} \Gamma(\alpha+1/2)$

84

where we have made the change of variable
$$z = C [b + (x - M)^{2}]$$

Define new parameters $V = 2a$ and $\overline{A} = 9b$.
St($x(M, \overline{A}, V) = \frac{P(W_{2} + \frac{1}{2})}{P(W_{2})} \left(\frac{\overline{A}}{\overline{X}V}\right)^{\frac{1}{2}} \left[1 + \frac{\overline{A}(x - M)^{2}}{V}\right]^{-\frac{1}{2}-\frac{1}{2}}$

known as Student's t - distribution. 7 is called precision and V is called the degree of freedom. When V = I, t - distribution reduces to the Couchy dist. While in the limit $V \rightarrow \infty$, t - distribution becomes Gaussian $N(X|M, \pi^{-1})$ Remark



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 + - distribution

$$St(X|M(, \Lambda \nu)) = \int_{0}^{\infty} N(X|M((1\Lambda)^{-1}) \operatorname{Gram}(1|V_{2}, V_{2}) d\eta$$
$$= \frac{P(V_{2} + V_{2})}{P(V_{2})} \frac{1\Lambda V_{2}}{(\pi\nu) V_{2}} \left[1 + \frac{\Delta^{2}}{\nu}\right]^{-V_{2} - V_{2}}$$

where

$$\Delta := (X - M)^{T} \wedge (X - M)$$

Remark

-
$$E[X] = M$$
 if $V > 1$

-
$$COV[X] = \frac{V}{(V-2)} \land if V > 2$$

mode [*] = MI

2.3.8 Periodic variables Consider an angular (polar) coordinate 050<27 and the problem of evaluating the mean of observations $D = \{Q_1, \dots, Q_N\}$ ∇ Simple average $(O_1 + \cdots + O_N)/N$ is strongly coordinate dependent. 2 Set angular observations as points on unit circle. X_1 be two-dim vector with $X_2 = (\cos Q_2, \sin Q_2)$ Let embedding

Average the vectors
$$\{X_n\}$$
 instead to give

$$\overline{X} = \frac{1}{N} \sum_{n=1}^{N} X_n$$

$$= \overline{\Gamma} (\cos \overline{\Theta}, \sin \overline{\Theta})$$
i.e. $\overline{\Gamma} \cos \overline{\Theta} = \frac{1}{N} \sum_{n=1}^{N} \cos \alpha_n$, $\overline{\Gamma} \sin \overline{\Theta} = \frac{1}{N} \sum_{n=1}^{N} \sin \alpha_n$
Thus we can solve for $\overline{\Theta}$ to give
$$\overline{\Theta} = \tan^{-1} \left\{ \frac{\sum_n \sin \alpha_n}{\sum_n \cos \alpha_n} \right\}$$

Consider
$$p(\theta)$$
 that have period 2π and must satisfies

$$p(\theta) \ge 0$$

$$\int_{0}^{2\pi} p(\theta) d\theta = 1$$

$$P(\theta + 2\pi) = p(\theta)$$
We can easily obtain a Gaussian - like distribution.
Consider a Gaussian over $x = Cx_1, x_2$ having mean $\mu = (\mathcal{M}, \mathcal{M}_2)$
and covariance matrix $\Sigma = \sigma^2 I$ so that

$$P(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp \left\{ -\frac{(x_1 - \mathcal{M}_1)^2 + (x_2 - \mathcal{M}_2)^2}{2\sigma^2} \right\}$$

$$(2.173)$$

Map $X = (X_1, X_2)$ and M into polar coordinates $\begin{array}{ccc} x_1 = r \cos \varphi, & x_2 = r \sin \varphi \\ M_1 = r_0 \cos \varphi, & M_2 = r_0 \sin \varphi_0 & fixed r_0, \varphi_0 \end{array} \end{array}$ Substitute these transformation into (2.173) with r=1 condition The exponent in (2.173) $-\frac{1}{15^2}\left(\left(r\cos\theta - r_0\cos\theta_0\right)^2 + \left(r\sin\theta - r_0\sin\theta_0\right)^2\right)$ (1 = 1) $= -\frac{1}{2\sigma^{2}} \left\{ 1 + r_{0}^{2} - 2r_{0} \cos \theta \cos \theta - 2r_{0} \sin \theta \sin \theta \right\}$ $=\frac{r_0}{\sigma^2}\cos(\theta-\theta_0)+\cos(\theta-\theta_0)$

Define $M = \frac{10}{52}$. Then we obtain the expression for the distribution of p(O) along unit circle $P(O \mid O_0, m) = \frac{1}{2\pi} \frac{1}{I_0 \text{ cm}} \exp\{m \cos(O - O_0)\}$ which called von Mises distribution. Here Oo represents the mean and $m = \frac{r_0}{\sigma^2}$ is called concentration parameter.

I. (m): zeroth - order Bessel function of the first kind

$$I_{o}(m) := \frac{1}{2\pi} \int_{0}^{2\pi} \exp\{m \cos \theta + d\theta$$

Now consider the maximum likelihood for 00 and m Observations D= { Q1,... Qx } is given $l_n P(P \mid Q_0, m) = \frac{N}{\Pi} P(Q_n \mid Q_0, m)$ (2.181) $= -N l_n(2\pi) - N l_n T_o(m) + m \sum_{n=0}^{\infty} cos(Q_n - Q_0)$ Set the derivative w.r.t Oo equal to zero gives $\sum_{n=0}^{N} \sin(\theta_n - \theta_0) = 0$ 121

Thus we obtain



Similarly maximizing (2.181) w.r.t m. Set the derivative of (1.181) w.r.t m, then we have $A(m) = \frac{1}{N} \sum_{n=1}^{N} \cos(\theta_n - \theta_n^{ML}) \qquad (2.185)$

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

$$A(m_{ML}) = \left(\frac{1}{N}\sum_{n=1}^{N} \cos \Theta_{n}\right) \cos \Theta_{0}^{ML} + \left(\frac{1}{N}\sum_{n=1}^{N} \sin \Theta_{n}\right) \sin \Theta_{0}^{ML}$$

Chapter 2 Probability Distributions

94

Here Acm, 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 Graussians

Limitations of a Graussian (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 Graussians: superposition of K Gaussians $p(x) = \sum_{k=1}^{K} \pi_k \mathcal{N}(X | \mathcal{M}_k, \Sigma_k)$

Each NCX (Mr, Zr) is called a component of mixture The parameter Tk are called mixing coefficients and satisfies $\sum_{k=1}^{k} \pi_k = 1 \qquad o \in \pi_k \leq 1$ **Figure 2.22** Example of a Gaussian mixture distribution p(x) 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 $p(x) = \sum_{k=1}^{k} p(k) p(x(k))$

 $\pi_k = pc_k$) prior probability of picking the kth component $\mathcal{N}(\mathcal{K}(\mathcal{M}_k, \mathcal{Z}_k)) = P(\mathcal{K}(k))$ probability of \mathcal{K} conditioned on k. Consider the posterior P(F|X) a.k.a. responsibilities $\delta_{K}(X) := P(F|X) = \frac{P(F,X)}{P(X)} = \frac{P(F)P(X|F)}{\Sigma_{Q}P(Q)P(X|Q)}$ $= \frac{\pi_{K}N(X|M_{K},\Sigma_{K})}{\Sigma_{Q}\pi_{L}N(X|M_{Q},\Sigma_{Q})}$ Groussian mixture is governed by $\pi := \{\pi_1, \dots, \pi_k\}, M := \{M_1, \dots, M_k\}$ and $\Sigma := \{\Sigma_{1,\dots}, \Sigma_k\}$ One way to set these parameters is to use maximum likelihood.

$$l_n P(X|\Pi, M, \Sigma) = \sum_{n=1}^{N} l_n \left\{ \sum_{k=1}^{K} \pi_k N(X_n|M_k, \Sigma_k) \right\}$$

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 XX, parameters
$$\eta$$
 (called natural parameters)
 $p(X | \eta) := h(X) g(\eta) exp{\eta^T} u(X)$ (2.194)
Here $u(X)$ is some function of XX and $g(\eta)$ can be interpreted
as the normalization coefficient. i.e.

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

Recall Bernoulli distribution

$$P(X|\mu) = Bern(X|\mu) = \mu^{X}(1-\mu)^{1-X}$$

$$x = 0, 1$$

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

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

Comparison with (2.194)

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

Solve for μ to give $\mu = \sigma(\eta) := \frac{1}{1 + \exp(-\eta)}$ logistic sigmoid

Thus, Bernoulli distribution can be rewrited in the form

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

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

u(x) = x h(x) = 1 $g(x) = \sigma(-1)$ Chapter 2 Probability Distributions Next consider the multinomial distribution

$$p(x(|\mu|) = \prod_{k=1}^{M} \mu_{k}^{x_{k}} = exp \left\{ \sum_{k=1}^{M} x_{k} l_{n} \mu_{k} \right\}$$

where
$$X = (X_1, ..., X_M)^T$$
 (one-hot vector)

The standard representation (2.149) so that

$$P(x|\eta) = exp(\eta^{T}x)$$

where
$$\eta = (\eta_1, \dots, \eta_M)^T$$
 with $\eta_k = ln \mathcal{M}_k$. i.e.

$$u(c_{x}) = x$$
, $h(c_{x}) = ($, $g(c_{1}) = ($

 $\sum_{k=1}^{M_k} M_k = 1$, parameters N_k are not independent. Since i.e. $M_{\rm M} = 1 - \sum_{k=1}^{{\rm M}-1} M_k$ left M-1 parameters $0 \leq M_{\rm E} \leq 1$ $M_{\rm E} \leq 1$ $M_{\rm E} \leq 1$ $-\sum \chi_{E} l_{n} (1-\Sigma M_{E})$ Z XE laME So the multinomial distribution becomes 1- 2 xc $\exp \left\{ \sum_{k=1}^{M} \chi_{k} \ln M_{k} \right\} = \exp \left\{ \sum_{k=1}^{M-1} \chi_{k} \ln M_{k} + \chi_{M} \ln M_{M} \right\}$ $ln(1-\Sigma ME)$

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

Now identify

$$ln\left(\frac{M_{k}}{1-\Sigma_{j}M_{j}}\right) = \eta_{k}$$

and
$$M_{k} = \frac{\exp(n_{j})}{1 + \Sigma_{j} \exp(n_{j})}$$
 softmax
normalized exponential

In this representation, multinomial distribution

$$P(\mathcal{X}|\eta) = \left(1 + \sum_{k=1}^{M-1} \exp(\eta_k)\right)^{-1} \exp(\eta^T \mathcal{X})$$

 $N = (N_1, ..., N_{M-1}, O)^T$ (M-dimensional)

$$u(x) = X$$
, $h(x) = 1$, $g(n) = (1 + \sum_{k=1}^{M-1} exp(n_k))^{-1}$

Finally, consider the univariate Graussian distribution.

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

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

$$\eta = \begin{pmatrix} M_{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix} \quad \text{urc x} = \begin{pmatrix} x \\ x^2 \end{pmatrix} \quad h(x) = (2\pi)^{-1/2} \quad g(\eta) = (-2\eta_1)^{1/2} \exp\left(\frac{\eta_1^2}{4\eta_1}\right)$$

2.4.1 Maximum likelihood and sufficient statistics
Consider the exponential family of distributions over
$$x$$

 $P(x(n) = h(x) g(n) \exp \{nTu(x)\}$ (2.194)

Taking the gradient of both side of
$$g(n) \int h(x) exp(n^T u(x)) dx = 1$$

w.r.t n, we have

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

+ $g(n) \int h(x) exp \{ n^T u(x) \} u(x) dx = 0$
Chapter 2 Probability Distributions
Using (2.195) then

$$-\frac{1}{g(n)}\nabla g(n) = g(n) \int h(x) \exp \left(nTu(x)\right) u(x) dx = \mathbb{E}[u(x)]$$

We therefore obtain the result

 $-\nabla h g(n) = E[u(x)]$

Now consider it samples tenoted by $X = \{x_1, \dots, x_n\}$ for which likelihood

$$P(X(n)) = \left(\frac{N}{\Pi} h(x_n) \right) g(n)^N \exp \left\{ n^T \sum_{n=1}^{N} u(x_n) \right\}$$

Chapter 2 Probability Distributions

Setting the gradient of $\ln P(X|n)$ w.r.t n to zero. We get the following condition to be satisfied by n_{ML} $-\nabla \ln g(n_{ML}) = \frac{1}{N} \sum_{n=1}^{N} u(cx_n)$

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

Chapter 2 Probability Distributions

2.4:2 Conjugate prior
For a given prob. density
$$p(X|n|)$$
, seek a prior $p(n)$
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 η
 $p(n|1 \times, \nu) = f(x, \nu) g(n)^{\nu} \exp \{\nu \eta^{T} \times \}$
where $f(x, \nu)$ is a normalization coefficient and $g(n)$ is
the same function in Chapter 2. Propagations 111

The posterior

PCNIX, X, V) OC PCXIN) · PCNIX, V)

oc g cni)
$$exp \left\{ \eta_{1}^{T} \left(\sum_{n=1}^{N} u(c \times n) + V \times \right) \right\}$$

2.4.3 No informative priors to have as little influence on the posterior as Intend possible. Let density or likelihood is given by P(X(7)) Consider noninformative prior (דא (First $p(\pi) = constant$ - If the domain of 7 is unbounded, prior cannot be normalized. Such prior is called improper - Transformation behavior of tensity under a nonlinear change **Chapter 2 Probability Distributions** of variables

Example 1.

Density of x takes the form $p(x|\mu) = f(x - \mu)$

Translation invariance

If
$$x \rightarrow \hat{x} := x + c$$
, then
 $\hat{p}(\hat{x}|\hat{\mu}) = f(\hat{x} - \hat{\mu})$

where we have tesine Hapter 21 Probability Distributions

Thus
$$p(x|\mu) = p(x|\mu)$$
 so density is independent of origin.
Prior distribution should satisfy this translation invariance property.

$$\int_{A}^{B} p(\mu) d\mu = \int_{A-c}^{B-c} p(\mu) d\mu = \int_{A}^{B} p(\mu-c) d\mu \quad \forall A, B$$
So we have $p(\mu-c) = p(\mu)$

Example of location parameter: mean of a Gaussian The conjugate prior for M is again Gaussian $p(M \mid M_0, \sigma_0^2)$ and we obtain noninformative prior by taking $\sigma_0^2 \rightarrow \infty$. Example 2.

Density of x takes the form $p(x(\sigma)) = -\frac{1}{\sigma}f(\frac{x}{\sigma})$ $\sigma > 0$

σ is known as scale parameter. E.g. N(x)μ,σ²)

Scale invariance

If $x \rightarrow \hat{x} := cx$, then $\hat{p}(\hat{x} \mid \hat{\sigma}) = \frac{1}{\hat{\sigma}} f(\frac{\hat{x}}{\hat{\sigma}})$

where we have tesine thapter Probability Distributions

So this transformation corresponds to a change of scale.
Prior distribution should satisfy this scale invariance property.

$$\int_{A}^{B} p(\sigma) d\sigma = \int_{A_{c}}^{B_{c}} p(\sigma) d\sigma = \int_{A}^{B} p(\frac{1}{c}\sigma) \frac{1}{c} d\sigma \quad \forall A, B$$
So we have $p(\sigma) = p(\frac{1}{c}\sigma) \frac{1}{c}$ and hence $p(\sigma) \alpha C \frac{1}{c}$
Note that this is an improper prior because of $0 < \sigma < \infty$

$$\int_{0}^{\infty} \frac{1}{\rho \alpha} = \frac{1}{\rho \alpha} \frac{1}{\rho} \int_{0}^{\sigma} \frac{1}{\rho$$

118

where
$$\tilde{X} := X - M$$

More convenient to work in terms of the precision $7 = \frac{1}{6^2}$ rather than σ itself $d_7 = \frac{1}{6^3} d_6$ $\sigma \rightarrow 7 = \frac{1}{6^2}$ $rather than <math>\sigma$ itself $\sigma \rightarrow 7 = \frac{1}{6^2}$ $rather than <math>\sigma$ itself $\sigma \rightarrow 7 = \frac{1}{6^2}$ $rather than <math>\sigma$ itself $\sigma \rightarrow 7 = \frac{1}{6^2}$ $rather than <math>\sigma$ itself $\sigma \rightarrow 7 = \frac{1}{6^2}$ $rather than <math>\sigma$ itself rather than rates the precision that the precision the precisi We have seen the conjugate prior for $71 \text{ uas } \text{Gram}(7120, b_0)$ The noninformative is obtained as the special case $20 = b_0 = 0$. 2.5 Non parametric Methods Approaches to density modeling Parametric vs Non parametric (Sew assumptions) method for density estimation Histogram Single continuous random variable X. Partition X into bins of width Δ_{λ} and then count the number distinct n: of observations of x falling in bin *ì* **Chapter 2 Probability Distributions**

120

We obtain probability values for each bin given by $P_{\vec{x}} = \frac{n_{\vec{x}}}{N \Delta_{\vec{x}}}$



Chapter 2 Probability Distributions

Remark

- Effect of a choice of width \triangle (smoothing parameter)
- After compting histogram, the data set can be discarded.
- Useful tool for a quick visualization of 1-4 or 2-4 data
- Limitation of high dimensional tata M, M bins in D-dim

2.5.1 Kernel density estimator
Estimate unknown probability density PCR) in D-dim space.
Cansider some small region R containing. X.
Then the probability mass associated with this R

$$P = \int_{R} PCR dx$$
 (true prob.)
Suppose we observed N data set drawn from PCR.
Since each point has a probability P of falling within R
Bin(KIN, P) = $\binom{N}{k} P^{k} (I-P)^{N-k}$ $k=0,1,...N$

k=0,1,... N

$$\Rightarrow E[\frac{5}{N}] = P, \quad Var[\frac{5}{N}] = P(1-P)/N$$

For large N,

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

$$P \simeq P(x) \vee$$

where V is the volchapter 2 Probability Distributions

Combining these expressions we obtain the density estimate $p(x) = \frac{k}{NV}$ (2.246) $R = \frac{q}{q}$

Remark: two contradictory assumptions on R and K

We can exploit (2.246) in two different ways

- K nearest neighbour method (fix K)

- kernel density estimator (fix V)

Kernel method in detail
Take the region R to be a small hypercube centered on X
To count the number k of points falling, within R, define

$$F(UI) := \begin{cases} 1 & |U|_{1}| \le \frac{1}{2} & \overline{x}=1, ..., D\\ 1 & 0 & 0 \text{ therwise} \end{cases}$$

KCUD is an example of kernel function, i.e. the quantiby
 $F(C(X-Xn)/h) = 1$ if Xn lies in a cube of side h centered
on X otherwise 0.

Chapter 2 Probability Distributions

The total number of points lying in the cube $k := \sum_{n=1}^{N} k \left(\frac{x - x_n}{h}\right) \qquad h = 1$

Substituting this expression into (2.246),

$$P(x) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{h^{p}} \kappa\left(\frac{x-x_{n}}{h}\right)$$

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

We can obtain a smoother density model (Gaussian kernel) $P(x) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{(2x h^2)^{\theta/2}} \exp \left\{ -\frac{11 x - x_n 11^2}{2h} \right\}$

where h represents the standard deviation of Graussian component and plays the role of a smoothing parameter. Generally, we can choose any other Kernel function K(UI) subject to

$$k(u) \ge 0$$

 $(k(u) \ge 0$

Chapter 2 Probability Distributions

Remark



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 5 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

local density estimation, fix value of k and For use the tata to find an appropriate value for V. a sphere centered on * and allow the radius Consider until it contains K tata points. i.e. the radius to grow is not determined (fixed) The value of K governs the tegree of smoothing. $P(x) = \frac{K}{NV}$ with KNN method for density Vse (2.246) **Chapter 2 Probability Distributions** 130 estimation

can be extended to classification. KNN method to each class separately and then make use KNN Apply Baye's theorem. र्भ × of data set, Nr: N: fotal × of points in Ck i.e. $\sum_{k=1}^{K} N_k = N$ Draw a sphere centered point on X (fixed) X New K points irrespective of their containing precisely class has the volume V and contains Kr points from sphere This class Cr. the



Similarly, the unconditioned density is given by $P(x) = \frac{k}{NV}$





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 tegree of smoothing
- KNN and kernel density methods require the entire

tata set to stored