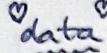


Statistical Inference & Multivariate AnalysisModule
No. MA3242024We have  data

We want to analyze, find distribution, etc.,

Graphical Methods

- Histogram
 - Symmetric
 - Peakness
 - Tail thickness
- Q-Q plot



Normal \rightarrow not good for modelling
non-negative values

{ Weibull, Gamma, Exp }

→ Empirical Sample estimator (To find empirical CDF from data)

i) Order sample $X_{(1)} \leq X_{(2)} \dots \leq X_{(n)}$

ii) $P[X_i = X_{(i)}] = 1/n$

$$\text{iii) } F(x) = P(X \leq x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

$$= \begin{cases} 0 & : x < X_{(1)} \\ \frac{i}{n} & : X_{(i)} \leq x < X_{(2)} \\ \dots & \dots \\ 1 & : x > X_{(n)} \end{cases}$$

Another method,

$$\hat{f}(x) \approx \frac{F(x+h) - F(x-h)}{2h}$$

$\hat{f}(x) \rightarrow$ Sreq density (konda like PDF)
~~at~~ at x

$$= \frac{1}{n} \sum_{i=1}^n \frac{I(x-h \leq X_i \leq x+h)}{2h}$$

$\{X_1, \dots, X_n\}$ each element is 1D
Assumptions

- univariate
- independent

Aim: What is best fitted distribution?

• Kernels \rightarrow some function $K(\cdot)$ following

\hookrightarrow used for approximating in an interval

$$1) K(\cdot) \geq 0$$

$$2) \int_{-\infty}^{\infty} K(x) dx = 1$$

Ex- II, $K(x_i - x)$

gaussian kernel

Kernel density estimation

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

Q-Q plot

Quantile Quantile plot

Characteristics of distributions

Skewness / Symmetric

Peakness of dist. / Kurtosis

Existence of outliers

Measure of heavy tail

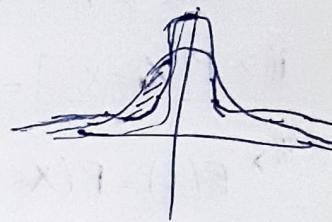
Dist

Normal

t-dist

Double exponential

$$\frac{1}{2b} e^{-\frac{|x-\mu|}{b}}$$



$$V(X) = E(X - \mu)^2$$

$$S(x) = E\left[\frac{X - \mu_x}{\sigma_x}\right]^3$$

$$K(x) = E(X - \mu)^4 \quad , \text{ If } K(x) > 3, \text{ high kurtosis}$$

Taking standard dist's,

$$\text{Normal} \rightarrow f(x) \propto e^{-\frac{x^2}{2}}$$

$$\text{DoubleExp} \rightarrow f(x) \propto e^{-|x|} \quad \rightarrow \text{Heavier tail than normal}$$

$$t\text{-dist} \rightarrow f(x) \propto \frac{1}{(1 + \frac{x^2}{\nu})^{\frac{\nu+1}{2}}} \approx \frac{1}{x^2} \quad \rightarrow \text{Heavier tail than normal}$$

$$\text{Pareto} \rightarrow f(x) \propto \frac{1}{(1+x)^{\alpha}} \quad \rightarrow \text{Heavier tail than double exp}$$

Heavy tail measurement
must be with reference
And,
Both mean & variance
of both must be same

Pareto



- Slowly Varying Function if $\frac{f(tx)}{f(x)} \rightarrow 1$ as $x \rightarrow \infty$

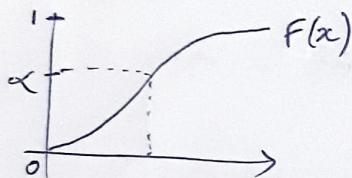
- Cauchy distribution

- Mixture Models

week 2

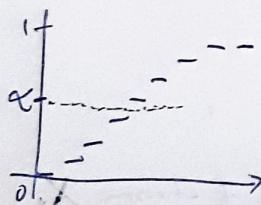
projects?
Quantile Regression
VaR

Quantile of a distribution



α^{th} Quantile point is $F^{-1}(\alpha) = \{x : F(x) = \alpha\}$
but only holds in case of cts functions

for discrete,



Then, α^{th} quantile pt. is $\hat{F}^{-1}(\alpha) = \inf \{x : \hat{F}(x) \geq \alpha\}$

for example: $\hat{F}(x) = \begin{cases} 0 & x < x_{(1)} \\ \frac{1}{n} & x_{(1)} \leq x < x_{(2)} \\ \frac{2}{n} & \\ \vdots & \\ 1 & x > x_{(n)} \end{cases}$

Then, for $0 < \alpha < \frac{1}{n}$, α^{th} quantile is $\inf \{x_{(1)}, x_{(2)}, \dots, x_{(n)}\} = x_{(1)}$

for any general α ,

$$F^{-1}(\alpha) = \begin{cases} x_{(\lfloor n\alpha \rfloor)} & , n\alpha \text{ is an integer} \\ x_{(\lfloor n\alpha \rfloor + 1)} & , n\alpha \text{ is not an integer} \end{cases}$$

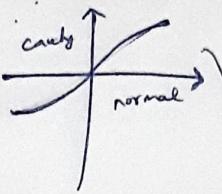
This is true when our 'origin' is $x_{(1)}$

if 'origin' is diff then,

$$F'(x) = \begin{cases} x_{(\lfloor n\alpha \rfloor + \text{origin} - 1)} & , n\alpha \text{ not} \\ x_{(\lfloor n\alpha \rfloor + \text{origin})} & , n\alpha \text{ not} \end{cases}$$

• QQ plot

→ quantiles of 2 diff distributions along diff axes

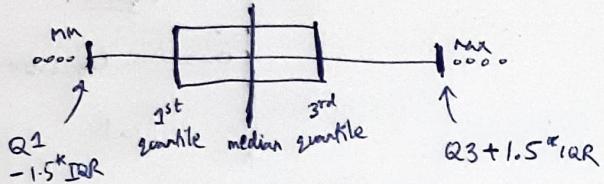


$\alpha = 0.99$	$\text{Normal}(0,1)$	$\text{Cauchy}(0,1)$
$\alpha = 0.98$		
$\alpha = 0.90$		
$\alpha = 0.85$		
\vdots		

Based on this,
can compare kurtosis/heavy tail ness
of 2 diff distributions

• Box Plot

$\xleftarrow{\text{IQR}}$ inter quartile range



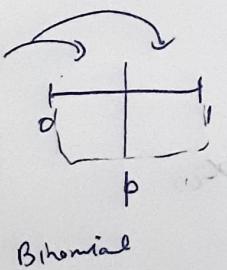
Goto from univariate → multivariate

- Extension of Binomial distribution in higher dimension → Multinomial

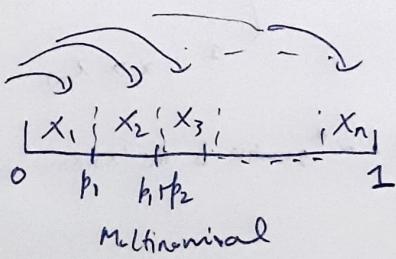
$$(X_1, X_2, \dots, X_K) \sim \text{Multinomial}(n, p_1, p_2, \dots, p_K)$$

$$\sum_{i=1}^K X_i = n \quad \sum_{i=1}^K p_i = 1$$

Generating elements from these dist's



Binomial



Multinomial

probability of choosing X_k

$$P\left(\sum_{i=1}^{k-1} p_i < v_i \leq \sum_{i=1}^k p_i\right)$$

$$= P[v_i \leq \sum_{i=1}^k p_i] - P[v_i < \sum_{i=1}^{k-1} p_i]$$

$$= \sum_{i=1}^k p_i - \sum_{i=1}^{k-1} p_i$$

$$= p_k$$

$$\bullet P[X_1 = x_1, \dots, X_k = x_k]$$

$$= \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$$

$$\bullet P[X_k = x_k] \sim \text{Bin}(n, p_k)$$

(proof using total prob thm)

gaze
gays
guage
guess
gauge

• Multinormal

$\{z_1, \dots, z_n\} \stackrel{iid}{\sim} N(0, I) \leftarrow n$ independent normals

$$f(\underline{z}) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n z_i^2}, \quad -\infty \leq z_i \leq \infty$$

multinormal

$$= \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \underline{z}' \underline{z}}$$

Note:- for normal dist, $\text{cov} = 0 \Leftrightarrow$ independence
But not in general

Let

$$\underline{z} \sim N(0, I)$$

$$\underline{x} \sim N(\underline{\mu}, \Sigma) \quad \text{then, } E(\underline{x}) = \underline{\mu}, \quad V(\underline{x}) = \Sigma$$

Since Σ is symmetric,

$$\text{we can write } \Sigma = \sum_{i=1}^n$$

:

Then, pdf of multivariate normal is, $\{z_1, z_2, \dots, z_d\} \stackrel{iid}{\sim} N(0, I)$

$$f(\underline{x}, \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu})}$$

$$\underline{x}, \underline{\mu} \rightarrow d \times 1$$

$$\Sigma \rightarrow d \times d$$

$$\underline{x} \sim N(\underline{\mu}, I)$$

$$\underline{x} = \underline{\mu} + \Sigma^{1/2} \underline{z}$$

$$\underline{z} \rightarrow p.d. \text{ matrix}$$

We see properties of Multivariate normal dist now,

i) Moment generating func

$$M_X(\underline{t}) = E_{\underline{X}}(e^{\underline{t}' \underline{X}})$$

=

$$\left| \begin{array}{l} \underline{X} = \underline{\mu} + \Sigma^{1/2} \underline{Z} \\ \underline{t}^* = \underline{t}' \Sigma^{1/2} \end{array} \right.$$

If $\underline{X} \sim N(\underline{\mu}, \Sigma)$

$M_X(\underline{t}) = e^{\underline{t}' \underline{\mu} + \frac{1}{2} \underline{t}' \Sigma \underline{t}}$

$\underline{t}^* = \underline{t}' \Sigma^{1/2}$

If $\underline{X} \sim N(\underline{\mu}, \Sigma)$

$$M_X(\underline{t}) = e^{\underline{t}' \underline{\mu} + \frac{1}{2} \underline{t}' \Sigma \underline{t}}$$

$\underline{t}^* = \underline{t}' \Sigma^{1/2}$

If $\underline{X} \sim N(\underline{\mu}, \Sigma)$

MGF of \underline{X} is $E(e^{\underline{t}' \underline{X}}) = e^{\underline{t}' \underline{\mu} + \frac{1}{2} \underline{t}' \Sigma \underline{t}}$

• $\underline{Y} = A \underline{X} \sim N(A \underline{\mu}, A \Sigma A')$

Then, $E[t' Y] =$

Ex: If $\underline{X}^{d \times 1} = \begin{pmatrix} X_1^{d_1 \times 1} \\ X_2^{(d-d_1) \times 1} \end{pmatrix}$, then find distⁿ of \underline{X}_1 .

Ans: Let $\underline{X}_1 = \underbrace{(I^{d_1 \times d_1}, 0)}_{\text{Take as } A} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$

Take $\underline{\mu} = \begin{pmatrix} \underline{\mu}_1^{d_1 \times 1} \\ \underline{\mu}_2^{(d-d_1) \times 1} \end{pmatrix}$.

$\Sigma = \begin{pmatrix} \Sigma_{11}^{d_1 \times d_1} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$

Then, $\underline{X}_1^{d_1 \times 1} \sim N(\underline{\mu}_1^{d_1 \times 1}, \Sigma_{11})$

ii) \underline{x}_1 and \underline{x}_2 are ind $\Leftrightarrow \Sigma_{12} = \Sigma_{21} = 0$

$$\begin{aligned} M_X(t) &= E(e^{t'X}) = E(e^{t_1'x_1 + t_2'x_2}) \\ &= e^{(t_1' + t_2')(\mu)} + (t_1', t_2') \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} \\ &= M_{X_1}(t_1) \cdot M_{X_2}(t_2) \end{aligned}$$

- $(X_1, X_2) \sim \text{BVN} (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$

then, $X_1 | X_2 = x_2 \sim N(\mu_1 + \frac{\sigma_1}{\sigma_2}(x_2 - \mu_2), \sigma_1^2(1 - g^2))$

Given $(X_1, X_2) \sim N(\mu, \Sigma)$

$$\Rightarrow (X_1 | X_2 = x_2) \sim N(\mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

Proof:

$$W = X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2$$

$$Y = X_2$$

Then, $\begin{pmatrix} W \\ Y \end{pmatrix} = A \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} I^{d_1 \times d_1} & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I^{d-d_1 \times d-d_1} \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$

Var! $\begin{pmatrix} W \\ Y \end{pmatrix} = A\Sigma A^T = \begin{pmatrix} I - \Sigma_{12}\Sigma_{22}^{-1} \\ 0 \end{pmatrix} \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} = \begin{pmatrix} \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$

$$E\begin{pmatrix} W \\ Y \end{pmatrix} = \begin{pmatrix} \mu_1 - \Sigma_{12}\Sigma_{22}^{-1}\mu_2 \\ \mu_2 \end{pmatrix}$$

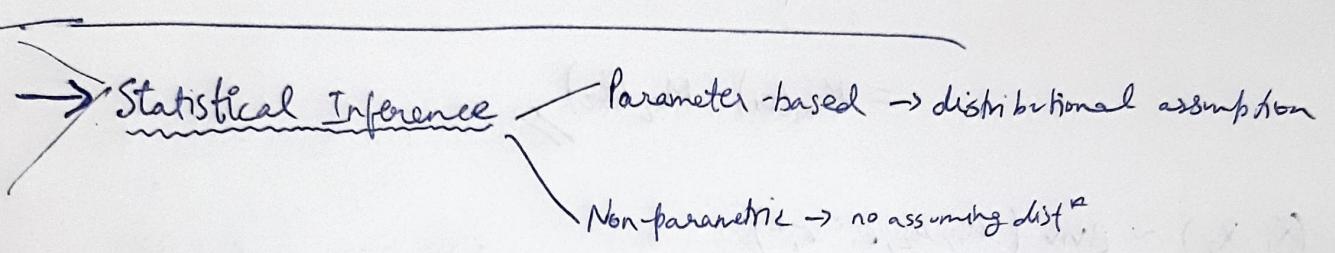
Since, W and Y are ind, dist of $W|Y$ is same as W .

so, $(W|Y) \sim N(\mu_1 - \Sigma_{12}\Sigma_{22}^{-1}\mu_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$

Putting. $X_2 = x_2$,

$$(X_1 \mid X_2 = x_2) \sim N(\mu_1 - \Sigma_{12} \Sigma_{22}^{-1} \mu_2 + \Sigma_{12} \dots, \dots)$$

$$\Rightarrow \sim N(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21})$$



Parameter-based approaches

Frequentist based
(θ is fixed)

- Maximum likelihood estimate
- method of moments
- minimax algo
- Uniformly minimum variance
- Unbiased estimate

Bayesian approach
(θ is random)

Hogg & Craig
Lee notes

• Loss/Risk for estimation

$$L(\hat{\theta}, \theta)$$

↑
loss

for example:

$$L(\hat{\theta}, \theta) = (\hat{\theta}(x_1, \dots, x_n) - \theta)^2$$
$$|\hat{\theta} - \theta| = \begin{cases} 1, & |\hat{\theta} - \theta| > \varepsilon \\ 0, & |\hat{\theta} - \theta| \leq \varepsilon \end{cases}$$

$$\text{Risk} = E(L(\hat{\theta}, \theta))$$

$$\{X_1, \dots, X_n\} \sim F(x, \theta) \quad \begin{matrix} \downarrow \\ \text{obs} \end{matrix} \quad \begin{matrix} \theta \rightarrow \text{fixed} & (\text{Frequentist setup}) \\ \text{we need to guess } \theta. \end{matrix}$$

$$\begin{aligned} \text{Risk} &= \mathbb{E}_{\underline{x}} [L(\hat{\theta}(\underline{x}), \theta)] \\ &= \iint_{\underline{x}} L(\hat{\theta}(\underline{x}), \theta) f(\underline{x}|\theta) d\underline{x} \end{aligned}$$

~~$\hat{\theta}(\underline{x})$~~
func of θ

But when $\theta \rightarrow \text{random}$, how to take average of loss (since we have 2 r.v's, X and θ)

• Bayesian Approach

$$P(B_i|A) \propto P(A|B_i) P(B_i)$$

\downarrow posterior $\underbrace{\downarrow}_{\text{likelihood function}}$ \downarrow prior $\xleftarrow{\text{relate this with}}$ $\prod f(x_i|\theta)$
 [Chances of observing the available info]

Bayes Thm

$$P(B_i|A) = \frac{P(A|B_i) P(B_i)}{\sum P(A|B_i) P(B_i)}$$

• Maximum Likelihood Estimate (MLE)

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} [\ln L(\theta)]$$

\downarrow
 $p(\theta | \underline{x})$
 \uparrow
 $p(\underline{x} | \theta)$

Let $\pi(\theta)$: prior dist.

$$p(\theta | \underline{x}) = \frac{f(\underline{x}|\theta) \pi(\theta)}{\int f(\underline{x}|\theta) \pi(\theta) d\theta}$$

$$\propto f(\underline{x}|\theta) \pi(\theta)$$

\uparrow likelihood \uparrow prior

$$\text{Risk} = \underset{(\theta|\underline{x})}{\mathbb{E}} [L(\hat{\theta}, \theta)] \quad (\because \theta \text{ is random})$$

Now, we can minimise risk and find a $\hat{\theta}$.

If loss is MSE, (and under frequentist setup)

$$\begin{aligned} E(L(\hat{\theta}, \theta)) &= E_x (\hat{\theta} - \theta)^2 \\ &= E_x (\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta)^2 \\ &= E (\hat{\theta} - E(\hat{\theta}))^2 + (E(\hat{\theta}) - \theta)^2 = \text{Variance} + \text{Bias}^2 \end{aligned}$$

↓

Bias = 0 \Rightarrow set of unbiased estimator

↓
Choose $\hat{\theta}$ \rightarrow variance becomes minimum

↓

Uniformly minimum variance unbiased estimator (UMVUE)

• Maximum likelihood

If $\{X_1, \dots, X_n\} \sim i.i.d. N(\mu, 1)$

$$\begin{aligned} \text{Then } L(\theta) &= f(x_1, \dots, x_n | \mu) \\ &= \prod f(x_i | \mu) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2} \end{aligned}$$

Ex:

$$\{X_1, \dots, X_n\} \sim f(x; \theta)$$

$$\text{After time } \rightarrow Y_i = \begin{cases} c, & X_i > c \\ X_i, & X_i \leq c \end{cases} \quad \text{Find maximum likelihood of } Y.$$

$$L(\theta) = \prod_{i=1}^n f(x_i, \theta)^{\delta_i} P[X_i > c]^{1-\delta_i}$$

$$\therefore \delta_i = \begin{cases} 0, & X_i > c \\ 1, & X_i \leq c \end{cases}$$

!! this is truncation on time,

consider ex. of
truncation on # of obs.

$$\text{Then } Y_i = \begin{cases} X_i, & i < n \\ X_n, & i \geq n \end{cases}$$

$$\frac{\partial \ln L(\theta)}{\partial \mu} = \frac{\partial}{\partial \mu} \left[-\frac{1}{2} \sum (x_i - \mu)^2 \right] = 0$$

$$\Rightarrow \hat{\mu} = \bar{x}$$

Ex: Weibull - β, θ

$$\hat{(\beta, \theta)} = \underset{\beta, \theta}{\operatorname{argmax}} \ln L(\beta, \theta)$$

$\ln L(f(x_i, \beta), \beta) = \text{profile log likelihood for } \beta$

Ex:

$$f(x; \theta) = \begin{cases} 1/\theta, & 0 \leq x_i \leq \theta \\ 0, & \text{o.w.} \end{cases}$$

order as $x_{(1)} < \dots < x_{(n)}$

Then,

$$L(\theta) = \prod_{i=1}^n f(x_i, \theta) = \begin{cases} 1/\theta^n, & 0 \leq x_{(1)} \leq \dots \leq x_{(n)} \leq \theta \\ 0, & \text{o.w.} \end{cases}$$

Now, $L(\theta)$ is ↓ sing func.

So, max likelihood is at min θ .

To see clearly, let $I(a, b) = \begin{cases} 1, & a \leq b \\ 0, & \text{o.w.} \end{cases}$ b.e.z $0 \leq x_{(1)} \Leftrightarrow x_i \geq 0 + i$
 $x_{(n)} \leq \theta \Leftrightarrow x_i \leq \theta + n$

Then, $L(\theta) = \prod_{i=1}^n f(x_i, \theta) = \frac{1}{\theta^n} I(0, x_{(1)}) I(x_{(n)}, \theta)$

Ex:

$$f(x, \theta) = \begin{cases} 1, & \theta \leq x \leq \theta+1 \\ 0, & \text{o.w.} \end{cases}$$

$$\prod_{i=1}^n f(x_i, \theta) = \begin{cases} 1, & \theta \leq x_{(1)} < \dots < x_{(n)} \leq \theta+1 \\ 0, & \text{o.w.} \end{cases}$$

$$= I(\theta; x_{(1)}) I(x_{(n)}; \theta+1)$$

So, MLE is at, $x_{(n)} - 1 \leq \theta \leq x_{(n)}$ \rightarrow so MLE may not be unique.

Ex: $X \sim \text{bin}(n, p)$

Ther, $L($



• EM algorithm

① Introduce 'missing observations' within the system

$$② E_{\underline{z} | \underline{x}, \theta^t} f_\theta(\underline{x}_i, \underline{z}_i) = \underset{\text{pseudo}}{\cancel{f_\theta(\underline{x}_i, \underline{z}_i)}} \text{ likelihood func}$$

\uparrow
unknown / hidden

$$= f(\underline{x}, \theta, \theta^t)$$

Arabin is
always correct
- Anant

depending on
evolution

Expectation step

$$Q(\theta, \theta^t) = E_{\underline{z} | \underline{x}, \theta^t} (f(\underline{x}, \underline{z} | \theta))$$

$\hat{\theta} \rightarrow \theta_{\text{true}}$

Maximization step

$$\hat{\theta}_{\text{EM}} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^t)$$

$f(x, \theta)$

So, $L(\theta) = \ln f(\underline{x}, \theta)$ = likelihood based on incomplete data

$$\ln L(\theta) - \ln L(\theta^t) = \ln \frac{f(\underline{x} | \theta)}{f(\underline{x} | \theta^t)}$$

reg show this ≥ 0 at $\theta = \hat{\theta}$

$$\begin{aligned} & \ln \frac{\int f(\underline{x}, \underline{z} | \theta) d\underline{z}}{f(\underline{x} | \theta^t)} \\ &= \ln \int \frac{f(\underline{x}, \underline{z} | \theta)}{f(\underline{x} | \theta^t)} d\underline{z} \\ &= \ln \int \left[\frac{f(\underline{x}, \underline{z} | \theta)}{f(\underline{x}, \underline{z} | \theta^t)} \frac{f(\underline{x}, \underline{z} | \theta^t)}{f(\underline{x}, \underline{z} | \theta^t)} \right] d\underline{z} \\ &\geq \int \ln \left[\frac{f(\underline{x}, \underline{z} | \theta)}{f(\underline{x}, \underline{z} | \theta^t)} \frac{f(\underline{x}, \underline{z} | \theta^t)}{f(\underline{x}, \underline{z} | \theta^t)} \right] d\underline{z} \end{aligned}$$

(By Jensen's inequality)
Since log is concave func

$$\Rightarrow \ln L(\theta) - \ln L(\theta^*) \geq \Delta L(\theta, \theta^*)$$

Say $g \propto f(x, z|\theta)$

Then,

$$\int \left[\ln \frac{f(x, z|\theta)}{f(x, z|\theta^*)} \right] f(z|x, \theta^*) dz$$

$$= \int \ln f(x, z|\theta) g(z|x, \theta^*) dz - \int \ln f(x, z|\theta^*) g(z|x, \theta^*) dz$$

\Rightarrow

↑
This term θ is not random,
so, when taking derivative, it
goes to 0.

- Gaussian Mixture model as clustering framework
- K-Means
- ↓ Spectral based cluster will provide assigned
can explore beyond also } for project

→ EM algo in "estimating the parameters of two mixture of normal distribution
↳ to figure out $\beta, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2$.

$$f(x; \beta, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2) = \beta f_1(x; \mu_1, \sigma_1^2) + (1-\beta) f_2(x; \mu_2, \sigma_2^2)$$

$X \rightarrow$ may be from f_1 or f_2

Let $Z_i = \begin{cases} 1, & X_i \sim f_1 \\ 0, & X_i \sim f_2 \end{cases}$ Then, $(X_i, Z_i), i=1, 2, \dots, n$
complete observations

everything
is concrete
- Me

Then,

$$f(x_i, z_i) = f_1^{z_i}(x_i) f_2^{(1-z_i)}(x_i)$$

\leftarrow but this is not the likelihood

the chance of observing
 x_i

So,

$$\begin{aligned} P[x_i \leq X_i \leq x_i + \Delta x_i, z = z_i] &= P[x_i \leq X \leq x_i + \Delta x_i | z_i = z] \\ &= p_i f(x_i) \end{aligned}$$

So, actual likelihood is,

b

$$f(x_i, z_i) = (p f_1(x_i))^{z_i} ((1-p) f_2(x_i))^{(1-z_i)}$$

$$\text{Complete likelihood is } (\text{for } i=1 \text{ to } n) = \prod_{i=1}^n f(x_i, z_i) = \prod_{i=1}^n \left[(p f_1(x_i))^{z_i} ((1-p) f_2(x_i))^{(1-z_i)} \right]$$

Then taking log,

$$\begin{aligned} \ln L &= \sum_{i=1}^n z_i [\ln p + \ln f_1(x_i; \mu_1, \sigma_1^2)] \\ &\quad + (1-z_i) [\ln(1-p) + \ln f_2(x_i; \mu_2, \sigma_2^2)] \end{aligned}$$

$$\text{Pseudolikelihood} = E_{(z_i | X, \theta^t)} (\ln L(\theta))$$

$$\begin{aligned} &= \sum_{i=1}^n (E_{z_i | X, \theta^t} (z_i))^I + \sum_{i=1}^n E_{z_i | X, \theta^t} (z_i (\ln f_1(x_i; \mu_1, \sigma_1^2)))^{II} \\ &\quad + \sum_{i=1}^n \ln(1-p) E_{z_i | X, \theta^t} (1-z_i) + \sum_{i=1}^n E_{z_i | X, \theta^t} ((1-z_i) \ln f_2(x_i; \mu_2, \sigma_2^2))^{IV} \end{aligned}$$

Then,

$$\begin{aligned} (E(z_i))^I &= \sum_{i=1}^n z_i P[z_i | X, \theta^t] \leftarrow (\text{from total probability theorem}) \\ &= E[\sum_i z_i P[z_i | X_i, \theta^t]] \text{ (depends on only } X_i, \text{ independent of others)} \end{aligned}$$

$$= \sum_i P[z_i = 1 | X_i = x_i, \theta^t]$$

$$= \sum_i \frac{p f_1(x_i; \mu_1, \sigma_1^2)}{p f_1(x_i; \mu_1, \sigma_1^2) + (1-p) f_2(x_i; \mu_2, \sigma_2^2)} \rightarrow m_i^t$$

(Here $\theta^t = (p^t, \mu_1^t, \sigma_1^{2t}, \mu_2^t, \sigma_2^{2t})$)

all are known
So, this expectation is known

$$\text{So, pseudolikelihood} = \sum_{i=1}^n m_i^t \ln p + \sum_{i=1}^n m_i^t [\ln f_1(\cdot)] \\ + \sum_{i=1}^n (1-m_i^t) \ln(1-p) + \sum_{i=1}^n (1-m_i^t) \ln f_2(\cdot)$$

Now, take derivative w.r.t $p, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2$, and maximise to get,

$$\hat{p}^{t+1} = \frac{1}{n} \sum_{i=1}^n m_i^t$$

$$\hat{\mu}_1^{t+1} = \frac{1}{\sum m_i^t} \sum_{i=1}^n m_i^t x_i$$

$$\hat{\sigma}_1^{t+1} = \frac{1}{\sum m_i^t} \sum_{i=1}^n m_i^t (x_i - \hat{\mu}_1^{t+1})^2$$

$$\hat{\mu}_2^{t+1} = \frac{1}{\sum (1-m_i^t)} \sum_{i=1}^n (1-m_i^t) x_i$$

$$\hat{\sigma}_2^{t+1} = \frac{1}{\sum (1-m_i^t)} \sum_{i=1}^n (1-m_i^t) (x_i - \hat{\mu}_2^{t+1})^2$$

} This can be used for estimating & clustering
one iteration formula
Has been converted
to single dimensional problem

This EM algo will
converge rapidly even
for far-away starting points

Arabir
↑
catch

sleeker
mature
scratches

- Calculation of MLE in case of Multivariate Normal dist^A

$$X \sim N(\mu, \Sigma), \quad D = \{x_1^{dx \times 1}, \dots, x_n^{dx \times 1}\} \text{ are iid}$$

\uparrow
d^{2+d} unknown
parameters

(at least these many samples needed)

Vector calculus
tangent for 30 mins
 $f: \mathbb{R}^n \rightarrow \mathbb{R}, f(x) = x^T A x$

$$\frac{df(x)}{dx^{n \times 1}} = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]$$

$$\frac{\partial f}{\partial x_i} = (x^T A + x^T A^T) e_i = x^T A + x^T A^T \\ = 2x^T A \quad (A \text{ is sym})$$

Matrix calculus

$$f: \mathbb{R}^{mn} \rightarrow \mathbb{R}$$

Derivative will be,

Result 1 $\frac{d \text{tr}(AB)}{dA} = B$

Result 2 If A is symmetric, $\frac{d \ln|A|}{dA} = A^{-1}$

• MLE for multivariate normal distⁿ

→ $f(x_i) = \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$

$$L(\mu, \Sigma; x) = \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$$

$$\ln L(\mu, \Sigma) = \ln \left[\frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)} \right]$$

wrt μ ,

$$\Rightarrow \frac{\partial \ln L(\mu, \Sigma)}{\partial \mu} = 0 \quad \begin{matrix} \text{wrt } \Sigma, & \text{take derivative wrt } \Sigma^{-1} \text{ instead} \\ \Rightarrow \frac{\partial \ln L(\mu, \Sigma)}{\partial \Sigma^{-1}} = 0 \end{matrix}$$

$$\Rightarrow \frac{1}{2} \times 2 \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} = 0 \quad \begin{matrix} \Rightarrow \frac{N}{2} \Sigma - \frac{d}{d \Sigma^{-1}} \frac{1}{2} \sum_{i=1}^N \text{tr}(\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T) = 0 \end{matrix}$$

$$\Rightarrow \hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \quad \begin{matrix} \Rightarrow \frac{N}{2} \Sigma - \frac{1}{2} \times \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T = 0 \end{matrix}$$

$$\Rightarrow \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

==

→ Gaussian mixture Model (GMM)

$$x \sim \sum_{k=1}^K \pi_k f(x; \mu_k, \Sigma_k)$$

• For univariate GMM,

$$f(x; \pi, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2) = \pi f_1(x; \mu_1, \sigma_1^2) + (1-\pi) f_2(x; \mu_2, \sigma_2^2)$$

For $x = \{x_1, \dots, x_n\}$

$\tilde{x} = \{x_1, \dots, x_n\}$ Then $\hat{\pi}_i, \hat{\mu}_i, \hat{\sigma}_i^2, \dots$

For higher dimensional GMM's,

$$\pi_k = \frac{1}{N} \sum_{i=1}^N m_{ik}^t, \quad z_{ik} = \begin{cases} 1, & z_i \sim \delta_k \\ 0, & \text{o.w.} \end{cases}$$

$$\hat{\mu}_k = \frac{1}{\sum_{i=1}^N m_{ik}^t} \sum_{i=1}^N m_{ik}^t z_i$$

$$\hat{\Sigma}_k = \frac{1}{\sum_{i=1}^N m_{ik}^t} \sum_{i=1}^N m_{ik}^t (z_i - \hat{\mu}_k)(z_i - \hat{\mu}_k)^T$$

→ GMM can be used for clustering also.

→ K-Means Algorithm

→ assumption: # of clusters = K

$$x_i \begin{bmatrix} 1 & 2 & \dots & K \end{bmatrix} \rightarrow \text{clusters}$$

$$x_1 \begin{bmatrix} 0 & 1 & \dots & 0 \end{bmatrix}, \quad x_2 \begin{bmatrix} 1 & 0 & \dots & 0 \end{bmatrix}, \quad \vdots$$

$$x_n \begin{bmatrix} 0 & 0 & \dots & 0 & 1 \end{bmatrix}$$

Let $\gamma_{ik} = \begin{cases} 1, & \text{if } x_i \text{ is in } k^{\text{th}} \text{ cluster} \\ 0, & \text{o.w.} \end{cases}$

Then objective func is,

$$J = \frac{1}{N} \sum_i \sum_k \gamma_{ik} \|x_i - \mu_k\|^2$$

Algo

2. Fix γ , get best μ to minimize J

1. Fix μ , compute γ (by minimizing J)

{

To get μ, γ

$$\frac{\partial J}{\partial \mu_k} = 2 \sum_{i=1}^N \hat{\gamma}_{ik} (x_i - \hat{\mu}_k)^T = 0 \Rightarrow \hat{\mu}_k = \frac{1}{\sum_{i=1}^N \hat{\gamma}_{ik}} \sum_{i=1}^N \hat{\gamma}_{ik} x_i$$

Exam
How is K-means
Connected
to EM algo/
Same MLE?

↑

(expression is
similar to GMM,
where $m_{ij} = \gamma_{ij}$)

GMM → soft thresholding
K-means → 0/1

So, in GMM, can
take cluster as
maximum M_{ik} value

book - Bishop

K-Means

- ① Minimizing MSE
- ② Hard thresholding
- ③ Good for spherical clusters

GMM

← for project also

minimizing log likelihood
 Soft
 Hard thresholding

Good for non-spherical clusters ← important point for project

(Extend more on this)
 ↴ elliptic shapes, spectral clustering, etc.
 * (will do later)

• MLE

Thm { $\hat{\theta}$ is an MLE of θ
 $g(\cdot)$ is an onto function
 Then, $g(\hat{\theta})$ is an MLE for $g(\theta)$ }

~~Statistical Inference~~ (for Bernoulli) (see later)

Proof: (Imp)

Since $g(\cdot)$ is onto,

$\theta \in H \rightarrow$ domain

$$\Lambda_y = \{ \theta : g(\theta) = y \} , H = \underline{\bigcup \Lambda_y}$$

Now,
 $\hat{y} = \underset{y}{\operatorname{argmax}} \left[\underset{\theta \in \Lambda_y}{\operatorname{Max}} L(\theta) \right] = \underset{\theta \in H}{\operatorname{argmax}} L(\theta)$

∴ $\hat{\theta}$ is an MLE of $\theta \Rightarrow \hat{\theta} = \underset{\theta \in H}{\operatorname{argmax}} L(\theta)$

Now,

$$\hat{y} = \hat{g}(\hat{\theta}) \xleftarrow{\text{same as}} g(\hat{\theta}) = y^* \Rightarrow \hat{\theta} \in \Lambda_{y^*}$$

↑

where $\Lambda_{y^*} = \{ \theta : g(\theta) = y^* \}$

claim: $L(\hat{\theta}) = \underset{\theta \in \Lambda_{y^*}}{\operatorname{Max}} L(\theta)$

, then $\Rightarrow \hat{\theta}$ which is arbitrary point in Λ_{y^*} gives the max of L

↳ this shows $y^* = \hat{y}$, so then $g(\hat{\theta}) = y^*$

Proof of claim ①,

$$L(\hat{\theta}) \geq \max_{\theta \in \Lambda_{y^*}} L(\theta) \quad \leftarrow \text{from definition of } L(\theta) \text{ and } \hat{\theta}$$

~~MLE~~

$$L(\hat{\theta}) \leq \max_{\theta \in \Lambda_{y^*}} L(\theta) \quad \leftarrow \text{since } \hat{\theta} \text{ is an arbitrary pt in } \Lambda_{y^*}$$

• Information Function

$$I(\theta) = E_x \left[-\frac{\partial^2}{\partial \theta^2} \ln f(x; \theta) \right] = E \left(\frac{\partial}{\partial \theta} \ln f(x, \theta) \right)^2 = V \left(\frac{\partial}{\partial \theta} \ln f(x, \theta) \right)$$

= expected value of hessian of $\ln f(x; \theta)$

Finding,

$$E \left(\frac{\partial}{\partial \theta} \ln f(x; \theta) \right)$$

$\frac{\partial}{\partial \theta} \ln f(x, \theta)$
called 'score' function
of f

$$= \int_x \left[\frac{\partial}{\partial \theta} \ln f(x; \theta) \right] f(x; \theta) dx$$

$$= \int_x \frac{1}{f(x; \theta)} \left(\frac{\partial}{\partial \theta} f(x; \theta) \right) f(x; \theta) dx$$

$$= \frac{\partial}{\partial \theta} \left(\int_x f(x; \theta) dx \right) \quad \begin{array}{l} (\text{under some regularity (?) conditions, we can interchange}) \\ \downarrow \\ \text{S and } \frac{\partial}{\partial x} \end{array}$$

$$= 0$$

① $f(x; \theta)$ is bounded func in x and domain of x
doesn't depend on θ

② $f(x; \theta)$ has as support, cts diff^A and integral uniformly
converges for all θ

Then, $\int_x \left[\frac{\partial^2}{\partial \theta^2} \ln f(x; \theta) \right] f(x; \theta) dx =$

• If $X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta)$

Then, $I_n(\theta) = n I(\theta)$

Proof:

$$\begin{aligned}
 I(\theta) &= E \left[\frac{-\partial^2}{\partial \theta^2} \sum_{i=1}^n \ln f(x_i; \theta) \right] \\
 &= \sum_{i=1}^n E \left(\frac{-\partial^2}{\partial \theta^2} \ln f(x_i; \theta) \right) \\
 &\stackrel{\wedge}{=} \sum_{i=1}^n I(\theta) = n I(\theta) //
 \end{aligned}$$

$V(\hat{\theta})_{MLE} \approx \frac{1}{n I(\theta)}$ for large n

$$\hat{T}(x) = \hat{\theta}_{MLE}$$

$$E(T(x)) = g(\theta) + o$$

\nearrow

$T(x)$ is an unbiased estimator

Then, $V(T(x)) \geq \frac{g'(\theta)^2}{n I(\theta)}$ (inequality true only in case of MLE)