

Uncertainty

- 1) Aleatoric Uncertainty: Uncertainty is due to randomness in the outcome.
- 2) Epistemic Uncertainty: Uncertainty arises due to lack of knowledge.

$P(y|\theta) \rightarrow$ Likelihood

$P(\theta|y) \rightarrow$ Posterior

Book: Bayesian Data Analysis (CRC Press)
- Gelman et al

Doing Bayesian Data Analysis
- John Kruschke

Assessment

End Sem - 45%

Mid Sem - 30%

Quiz + Assignment - 20% + 5% (Turing).

Bayesian Inference

- Process of fitting a probability model to a set of data and summarizing the results by a prob. distribution over the parameters of the model and on the unobserved quantities. like prediction outcome.

Process of Bayesian Data Analysis (BDA)

- 1) Setting up a full model
 - ↳ Joint probability distribution of the observed and the unobserved quantities in the ~~data~~ problem: θ^* .
- 2) Conditioning on the observed data $P(\theta|y)$,
 $P(\tilde{y}|y)$
- 3) Evaluating the fit of the model and the implications of the posterior distribution.

Notation

θ → Unobservable vector quantities or model parameters.

y → Denotes the observed data.

\tilde{y} → Unknown but potentially observed quantities
Generally

- Greek letter will denote the model parameter
- Lower case Roman letters will denote the observed or observable scalars or vectors
- Upper case Roman letters will denote the observed or observable matrix
- Vectors would be represented as a column vector, $u^T v$ would be a scalar and $u v^T$ is a matrix

Exchangeability

Ex: H, T, H, H, H, T ...

Suppose θ is known

$$P(H, T, H, H, H, T) = P(H) P(T) \cdot P(H) \dots$$

Exchangeable If the sequence is permuted the probability will not change

$$= \theta^{\# \text{heads}} (1-\theta)^{\# \text{tails}}$$

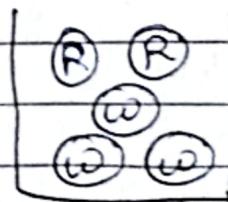
Consider the size of tumors observed over 1 year interval, $S_1, S_2, S_3 \dots$

$$\begin{aligned} P(S_1, S_2, S_3, S_4) & \quad ? \text{ Non Exchangeable Model} \\ P(S_1, S_3, S_4, S_2) \end{aligned}$$

Note Exchangeability is related with iid:
 $iid \Rightarrow$ Exchangeability.
 BUT

Exchangeability $\not\Rightarrow$ iid.

\Rightarrow Polya's Urn



Taking out a ball without replacement

$$P(R, W, R, W, W) =$$

$$\frac{2}{5} \times \frac{3}{4} \times \frac{1}{3} \times 1 \times 1 = \frac{1}{10}$$

$$P(W, R, W, W, R) = \frac{3}{5} \times \frac{2}{4} \times \frac{2}{3} \times \frac{1}{2} \times 1 = \frac{1}{10}$$

This is exchangeable but not iid as the next turn depends on previous one.

Explanatory Variables

The pair (x, y) ; is exchangeable.

Bayes Rule of Inference

$$\theta = \frac{\# \text{ Red Marbles}}{\text{Total Marbles}}$$

$$P(y, \theta) = P(\theta|y) \cdot P(y) = P(y|\theta) \cdot P(\theta)$$

$$P(\theta|y) = \frac{P(y|\theta) \cdot P(\theta)}{P(y)} \rightarrow \text{Prior (Not dependent on past data)}$$

Marginal $\sum_{\theta} P(y|\theta)P(\theta)$

$$P(\theta|y) \propto P(y|\theta)P(\theta)$$

Predictive Inference

To make prediction about an unknown observable.

$$p(y) = \int_0 p(y|\theta) d\theta = \int_0 p(y|\theta) p(\theta) d\theta$$

Prior predictive
Distribution.

→ Prior because no past data is used
Predictive because its making a prediction
on unknown variable.

$P(\tilde{y}|y) \rightarrow$ Suppose an object is weighed on a scale n times, y_1, y_2, \dots, y_n

Posterior Predictive

Distribution. Let the actual weight be M and the variance in the scale be σ^2

Let the unknown parameters denoted by θ

$$P(\tilde{y}|y) = \int_{\theta} P(\tilde{y}, \theta|y)$$

$$= \int_{\theta} P(\tilde{y}|\theta, y) P(\theta|y) \cdot d\theta$$

$$= \int_{\theta} P(\tilde{y}|\theta) P(\theta|y) \cdot d\theta$$

Posterior Odd Ratio

$$\frac{P(\theta_1|y)}{P(\theta_2|y)} = \frac{P(y|\theta_1) \cdot P(\theta_1)}{P(y|\theta_2) \cdot P(\theta_2)} = \frac{P(\theta_1)}{P(\theta_2)} \times \frac{P(y|\theta_1)}{P(y|\theta_2)}$$

↓
Prior Odds Likelihood Ratio.

If θ changes a little $P(y|\theta_1) \cdot P(\theta)$ changes largely.

Problems related to these.

1) Prior may not be available all the time

2) Modeling the likelihood with a distribution

3) $P(y)$. θ can have exponential space and could have multiple θ_A .

Case Study : Autocorrect and Autocompletion

If someone writes "gradom"

The possibilities are : \rightarrow "random"

\rightarrow "gradon"

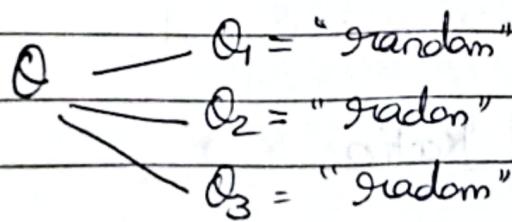
\rightarrow "gradom"

What is the probability that the person has used the word correctly.

Let y be the word the person has written

Here, $y = \text{gradom}$

Let θ be the word that the person is actually supposed to write.



$$P(\theta|y) \propto P(y|\theta) P(\theta)$$

Prior Information : Finding the frequency of θ in search keywords

θ	$P(\theta)$
random	7.6×10^{-5}
gradon	6.5×10^{-6}
gradom	3.12×10^{-7}

Likelihood: $P(y|\theta)$ Can be obtained from spell check

θ	$P(y \theta) = p(\text{random} \theta)$
random	0.00193
gradan	0.000143
gradam	0.975

Posterior

θ	$P(\theta y)P(\theta)$	$P(\theta y)$
random	1.47×10^{-7}	0.325
gradan	8.6×10^{-10}	0.002
gradam	3.04×10^{-7}	0.674

We can include the context words and model it like $P(\theta|y, y)$

We can also add sentiments to improve results.

Useful results from Probability Theory

$f(u, v)$ → Joint probability distribution

$f(u|v)$ → Conditional probability distribution

$f(u)$ → Marginal probability distribution ($\int f(u, v)dv$)

$$f(u, v, w) = f(u, v|w) f(w) = f(u|v, w) f(v|w) f(w)$$

Ex: An urn contains n_r red balls, n_g green balls and n_b blue balls. We pick three balls without replacement. What is the probability that they are of different colors.

Mean

$$E(u) = \int u f(u) du \text{ or } \sum_u u f(u)$$

Variance

$$\text{Var}(u) = \int_u (u - E(u))^2 f(u) du = E[(u - E(u))^2]$$

Covariance

$$\text{Cov}(u) = \int_u (u - E(u))(u - E(u))^T du, \text{ } u \text{ is a column vector}$$

Modeling using conditional probability

- Probability models often express the distribution conditionally or hierarchically
- The height, y , of a student selected randomly in a class

Avg. height of boys = 170

Avg. height of girls = 160

If $f(\text{boy}) = f(\text{girl}) = 1/2$

$$f(y) = f(y|\text{boy}) f(\text{boy}) + f(y|\text{girl}) f(\text{girl})$$

Marginal probability with bimodal distribution.

Conditional Mean

$$E(u) = \int_u u f(u) du = \iint_{uv} u f(u,v) du dv$$

$$E(u) = \iint_{u,v} u f(u|v) f(v) du dv$$

$$= \iint_{u,v} u \underbrace{f(u|v) du}_{\text{conditional mean}} f(v) dv$$

$$= \int_v E(u|v) f(v) dv$$

$$= E[E(u|v)]$$

Conditional Variance

$$\text{Var}(u) = E(\text{Var}(u|v)) + \text{Var}(E(u|v))$$

$$E(\text{Var}(u|v)) + \text{Var}(E(u|v))$$

$$= E(E(u^2|v) - (E(u|v))^2) + E[(E(u|v))^2] - (E(E(u|v)))^2$$

$$= E(E(u^2|v)) - E[(E(u|v))^2] + E[(E(u|v))^2] - (E(E(u|v)))^2$$

$$= E(E(u^2|v)) - (E(E(u|v)))^2$$

$$= E(u^2) - (E(u))^2$$

$$= \text{Var}(u)$$

Q Let X be a discrete random variable with support $S = \{0, 1\}$ and let Y be a discrete Random variable with support $S_2 = \{0, 1, 2\}$. The joint distribution $f_{X,Y}(x,y)$ is given as follows:

		P _{X,Y} (x,y)			P _X (x)
		0	1	2	
0	0	1/8	2/8	1/8	1/2
	1	2/8	1/8	1/8	1/2

i) Find the conditional mean of $Y|X$.

$$\text{Sol: } E(Y|X) = \frac{(1/8 + 2/8 + 1/8) \cdot 0 + (2/8 + 1/8 + 1/8) \cdot 1}{1/2} = \frac{1}{2}$$

$$E(Y|X) = \frac{1}{2} \sum_y y f(y|x)$$

$$E(Y|X=0) = \sum_y y f(y|x=0)$$

$$= 0 \cdot f(y=0|x=0) + 1 \cdot f(y=1|x=0) + 2 \cdot f(y=2|x=0)$$

$$f_{X,Y}(x,y) = f_x(x|y) f(y)$$

~~$$f_x(x) f(y|x) \cdot f(x) \text{ or } f(y|x) = \frac{f_{X,Y}(x,y)}{f(x)}$$~~

$$E(Y|X=0) = 1 \cdot \frac{2/8}{1/2} + 2 \cdot \frac{1/8}{1/2} = 1$$

$$E(Y|X=1) = 0 \cdot \frac{1/8}{1/2} + 1 \cdot \frac{2/8}{1/2} + 2 \cdot \frac{1/8}{1/2} = \frac{3}{4}$$

2) Find the conditional variance of $Y|X$

Sol $\text{Var}_c(Y|X) = E(Y^2|X) - (E(Y|X))^2$

$$\sum y^2 f(y|x)$$

$$0 f(Y|X=0) + 1 f(Y|X=0) + 4 f(Y|X=0)$$

$$\frac{1}{2} +$$

Transformation of Variables

Suppose u is a random variable

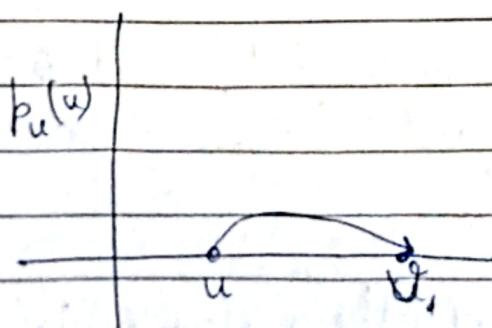
$f_u(u) \rightarrow$ Distribution of the Random Variable u .

$$v = f(u)$$

$f_v(v) \rightarrow$ Distribution of v we want to find.

Case-1 $f_u(u)$ is discrete and v is a one to one mapping of u , $v = f(u)$

$$f_v(v) = f_u(f^{-1}(v))$$



Q Suppose X is a geometric distribution with $f_x(x)$

$$f_x(x) = \frac{3}{4} \left(\frac{1}{4}\right)^{x-1}, x \in \mathbb{Z}$$

Find the probability distribution of the R.V; $Y = X^2$

Sol

$$Y = X^2$$

$$X = \sqrt{Y}$$

$$f_Y(y) = \begin{cases} \frac{3}{4} \left(\frac{1}{4}\right)^{\frac{y-1}{2}} & \text{for } Y=1, 4, 9 \\ 0 & \text{otherwise} \end{cases} \quad \text{for } Y=1, 4, 9$$

Case 1 (a)

Suppose X_1 and X_2 are 2 discrete random variables with a joint probability distribution $f_{X_1, X_2}(x_1, x_2)$ and Let Y_1 be a function $Y_1 = U_1(x_1, x_2)$ and $Y_2 = U_2(x_1, x_2)$

One to One Transformation

So we can solve x_1 and x_2 in terms of Y_1 and Y_2

$$x_1 = w_1(Y_1, Y_2)$$

$$x_2 = w_2(Y_1, Y_2)$$

Find distribution $f_{Y_1, Y_2}(Y_1, Y_2)$

$$f_{Y_1, Y_2}(Y_1, Y_2) = f_{X_1, X_2}(w_1(Y_1, Y_2), w_2(Y_1, Y_2))$$

Q

Suppose X_1 and X_2 be two random independent random variable with Poisson distribution with mean μ_1, μ_2

$$f(x) = \frac{\lambda^n e^{-\lambda}}{x!}$$

Find the distribution of the random variable

$$Y = X_1 + X_2$$

Sol

$$Y = X_1 + X_2$$

$$V = X_1 + X_2$$

~~$$X_1 = Y - V, \quad X_2 = Y - V$$~~

$$f_{X_1, X_2}(x_1, x_2) = \frac{\lambda^{x_1+x_2} e^{-2\lambda}}{x_1! x_2!}$$

$$f_{X_1, X_2}(x_1, x_2) = \frac{\lambda^{x_1} e^{-\lambda}}{x_1!} \cdot \frac{\lambda^{x_2} e^{-\lambda}}{x_2!}$$

$$f_{Y, V}(y, v) = \frac{\lambda^y e^{-\lambda}}{y!} \frac{\lambda^{v-y} e^{-\lambda}}{(v-y)!}$$

$$f_Y(y) = \sum_v \frac{\lambda^y}{v!} \frac{\lambda^{v-y}}{(v-y)!} e^{-2\lambda}$$

$$\frac{e^{-(\lambda_1+\lambda_2)}}{\lambda_1!} \lambda_2^y \sum_v \left(\frac{\lambda_1}{\lambda_2}\right)^v \frac{y!}{v! (y-v)!}$$

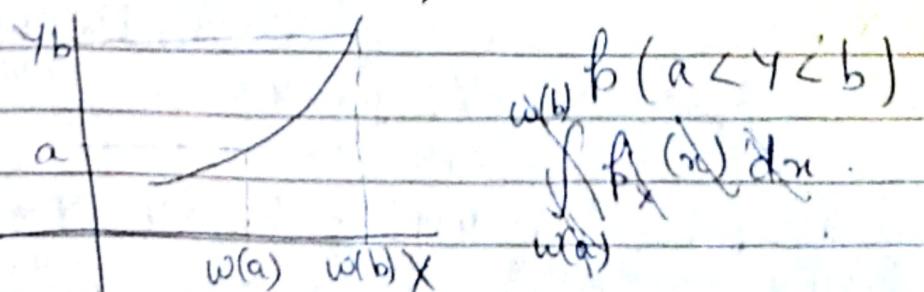
$$\frac{e^{-(\lambda_1+\lambda_2)}}{\lambda_1!} \lambda_2^y \sum_v v! C_v \left(\frac{\lambda_1}{\lambda_2}\right)^v$$

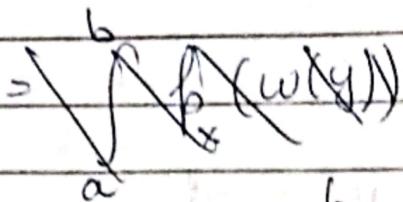
$$\frac{e^{-(\lambda_1+\lambda_2)}}{\lambda_2!} \lambda_2^y \sum_v v! C_v \lambda_1^{y-v} \lambda_2^{-v}$$

$$\frac{e^{-(\lambda_1+\lambda_2)}}{\lambda_1!} (\lambda_1+\lambda_2)^y \sim f_{Y, V}(y)$$

Case 2 Suppose X is the continuous random variable, $f_X(x)$ being its distribution and let $Y = \omega(x)$ define a one to one mapping from $X \rightarrow Y$. So $X = \omega^{-1}(Y)$. Then the probability distribution $f_Y(y) = f_X(\omega(y)) |\omega'(y)|$, where $J = \omega'(y)$, Jacobian.

Proof





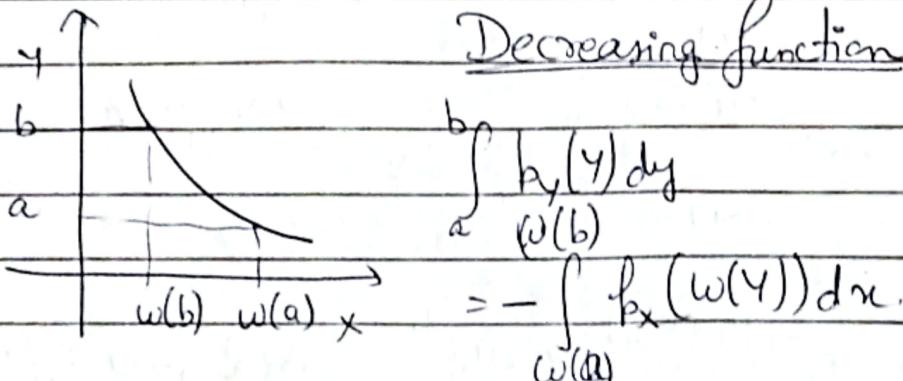
$$P_Y(a < Y < b) = \int b_y(y) dy$$

$$= \int_a^{w(b)} f_x(w(y)) dx \quad x = w(y) \\ dx = w'(y) dy$$

$$= \int_a^b f_x(w(y)) w'(y) dy$$

$$f_x(y) = f_x(w(y)) |w'(y)|.$$

Decreasing function.



$$= - \int_{w(a)}^{w(b)} f_x(w(y)) \underline{w'(y)} dy.$$

≤ 0

makes (+) ve.

$$f_x(y) = f_x(w(y)) |w'(y)|.$$

Q

Let X be a continuous random variable with prob.

distⁿ
$$f_x(x) = \begin{cases} x/12 & \text{for } 0 < x < 5 \\ 0 & \text{otherwise.} \end{cases}$$

Find the probability distⁿ of the R.V. $Y = 2x - 3$

Solⁿ

$$f_y(y) = \begin{cases} \frac{1+3}{24} \cdot \frac{1}{2} & \text{for } 1 < \frac{y+3}{2} < 5 \\ 0 & \text{otherwise.} \end{cases}$$

$$f_y(y) = \begin{cases} \frac{1}{48} & -1 \leq y \leq 7 \\ 0 & \text{Otherwise} \end{cases}$$

Case-3 Suppose X is a vector of continuous random variables and $Y = u(X)$ with one to one mapping $X = w(Y)$

$$f_y(y) = f_x(w(y)) |J|$$

where J is a matrix

Note $\frac{\partial Y}{\partial X}, Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$$\frac{\partial Y}{\partial X} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \frac{\partial y_n}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_n} \end{bmatrix} = J$$

Case-4 If X is a continuous random variable with probability distribution given as $f_x(x)$ and $Y = U(X)$ defines a transformation between the value of X and Y that is not one-to-one. If the interval over which X is defined can be partitioned into k mutually disjoint sets such that each of the inverse functions $x_1 = w_1(y)$, $x_2 = w_2(y) \dots x_n = w_n(y)$ that defines a one to one correspondence then the $f_y(y)$ can be defined as

$$f_y(y) = \sum_{i=1}^k f_x(w_i(y)) |J_i| \quad \text{where } |J_i| = |w_i'(y)|$$

$i = 1, 2, 3 \dots k$

Q

Goat	Car	Nothing
A	B	C

"Let's Make a Deal" / "The Monte-Hall" Problem

When the participant chooses a door some other non-prize door is opened then he is being given a chance to switch his choice.

A is chosen by participant

$$P(A) = \cancel{\frac{1}{3}}$$

Likelihood if A contains prize food of host to open B.

$$P(B=\text{opened} | A) = \frac{1}{2}$$

$$P(B) = \cancel{\frac{1}{3}}$$

$$P(B=\text{opened} | B) = 0$$

$$P(C) = \cancel{\frac{1}{3}}$$

A is chosen and C has the prize so opening of B is certain.

$$P(\text{prize in A} | B=\text{opened}) = P(B=\text{open} | A) \cdot P(A)$$

$$\sum_{A, B, C} P(B=\text{open} | A) \cdot P(A)$$

$$= \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + 0 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} = \frac{1}{3}$$

$$P(B | B) = 0$$

$$P(C | B=\text{opened}) = \frac{1 \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + 0 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} = \frac{2}{3}$$

Probability is in C given B is opened by the host

Distributions

$$P(Y|\theta, n) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

Likelihood

→ Binomial Distribution.

$$\text{Posterior: } P(\theta|y, n) = \frac{P(y|\theta) \cdot P(\theta)}{\int_{\theta=0}^1 P(y|\theta) \cdot P(\theta) d\theta}$$

$$\propto P(y|\theta) \cdot P(\theta)$$

$$\propto P(y|\theta)$$

$$= \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

$$\boxed{\beta(\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}}$$

$$P(\theta|y, n) \propto \beta(y+1, n-y+1, \theta)$$

$$\boxed{\Gamma(n) = (n-1) \Gamma(n-1) = (n-1)!}$$

$$\Gamma(1) = 1$$

$$P(\theta|y, n) = \frac{n!}{y!(n-y)!} \theta^y (1-\theta)^{n-y}$$

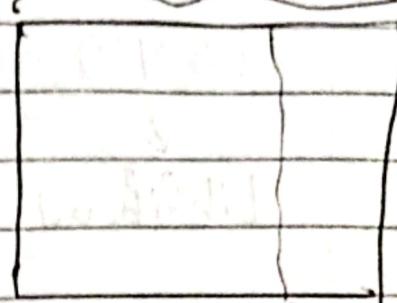
$$= \frac{n(n-1) \dots (n-y+1)}{y!(n-y)!} \theta^y (1-\theta)^{n-y}$$

$$= \frac{n \Gamma(n)}{y \Gamma(y) (n-y) \Gamma(n-y)} \theta^y (1-\theta)^{n-y}$$

$$\propto \frac{n}{y(n-y)} \beta(\theta, y+1, n-y+1)$$

Bayes' Experiment

- 1) A ball w was thrown randomly
 θ is the fraction of the width
 on which it stopped.



- 2) Another ball O was dropped n number of times. The number of success Y in the count of the times the ball landed on the right of θ .

$$P(Y|\theta, n) = \binom{n}{Y} \theta^Y (1-\theta)^{n-Y}$$

What he tried to observe $P(\theta \in (\theta_1, \theta_2))|y$

$$= \frac{\int_{\theta_1}^{\theta_2} P(Y|\theta) P(\theta) d\theta}{\int_{\theta_1}^1 P(Y|\theta) P(\theta) d\theta}$$

Assumptions:

- θ is uniformly distributed in the range $[0, 1]$

→ The denominator

$$\int P(Y|\theta) P(\theta) d\theta$$

$$= \int \binom{n}{Y} \theta^Y (1-\theta)^{n-Y} d\theta$$

$\therefore P(\theta)$ is uniform for all θ

$$= \binom{n}{y} \int_0^y \theta^y (1-\theta)^{n-y} d\theta = I_n \text{ (say)}$$

$$\int u v dx = u \int v dx + \int [u' \int v dx] dx.$$

$$I_n = \binom{n}{y} \left[\int_0^y \int (1-\theta)^{n-y} d\theta + \int (y \theta^{y-1} \int (1-\theta)^{n-y} d\theta) d\theta \right]$$

$$\binom{n}{y} \left[-\theta^y \frac{(1-\theta)^{n-y+1}}{n-y+1} \Big|_0^1 + \int y \theta^{y-1} \frac{(1-\theta)^{n-y+1}}{n-y+1} d\theta \right]$$

Applying limits.

$$\binom{n}{y} \int y \frac{\theta^{y-1} (1-\theta)^{n-y+1}}{n-y+1} d\theta$$

$$\binom{n}{y} \left(\frac{y}{n-y+1} \right) \left[\theta^{y-1} \int (1-\theta)^{n-y+1} d\theta + \int \frac{(y-1) \theta^{y-2} (1-\theta)^{n-y+2}}{n-y+2} d\theta \right]$$

$$\binom{n}{y} \left(\frac{y}{n-y+1} \right) \left[-\theta^{y-1} \frac{(1-\theta)^{n-y+2}}{n-y+2} \Big|_0^1 + \int \frac{y-1}{n-y+2} \theta^{y-2} (1-\theta)^{n-y+2} d\theta \right]$$

$$+ \binom{n}{y} \left(\frac{y}{n-y+1} \right) \left(\frac{y-1}{n-y+2} \right) \cdot \left(\frac{1}{(n+1)} \right) \int_0^1 (1-\theta)^n d\theta$$

$$\binom{n}{y} \frac{y!}{(n-y+1)(n-y+2) \dots n} \left[\frac{-(1-\theta)^{n+1}}{n+1} \right]_0^1$$

$$\binom{n}{y} \frac{y!}{(n-y+1)(n-y+2) \dots n(n+1)}$$

$$\frac{n!}{y! (n-y)!} \frac{y!}{(n-y+1)(n-y+2) \dots n(n+1)} = \frac{n!}{(n+1)!} = \frac{1}{n+1}$$

Q What is the probability of getting a success if he already got y successes.

$$\text{Sol} \quad P(\tilde{Y}=1|y) = \int_0^1 P(\tilde{Y}=1, \theta|y) d\theta$$

$$\int_0^1 P(\tilde{Y}=1|\theta, y) \cdot P(\theta|y) d\theta$$

$$\int_0^1 \underbrace{P(\tilde{Y}=1|\theta)}_{\text{Posterior}} \cdot P(\theta|y) d\theta$$

$$\int_0^1 \theta P(\tilde{Y}=1|\theta) \cdot P(\theta|y) d\theta = E[\theta|y]$$

$$\int_0^1 \theta P(y|\theta) \cdot P(\theta) d\theta$$

Posterior Predictive Distribution

$$P(\tilde{Y}=1|y) = \int_0^1 P(\tilde{Y}=1, \theta|y) d\theta$$

$$= \int_0^1 P(\tilde{Y}=1|\theta, y) \cdot P(\theta|y) d\theta$$

$$= \int_0^1 \theta P(\tilde{Y}=1|\theta) \cdot P(\theta|y) d\theta$$

$$= \int_0^1 \theta P(\theta|y) d\theta$$

$$= \frac{\int_0^1 \theta P(\theta|y) d\theta}{P(y)}$$

$$\propto \int_0^1 \theta \frac{P(\theta|y) d\theta}{\binom{n+y}{n}}$$

$$= (n+1) \int_0^n \theta^y (1-\theta)^{n-y}$$

$$= (n+1) \binom{n}{y} \int_0^n \theta^{y+1} (1-\theta)^{n-y}$$

$$I_{n-y, y+1} = \int_0^n (1-\theta)^{n-y} \theta^{y+1} d\theta$$

$$= (1-\theta)^{n-y} \int_0^n \theta^{y+1} d\theta + \int_0^n (n-y)(1-\theta)^{n-y-1} \int_0^n \theta^{y+1} d\theta$$

$$\left[\frac{(1-\theta)^{n-y} \theta^{y+2}}{y+2} \right] + \int_0^n \frac{n-y}{y+2} (1-\theta)^{n-y-1} \theta^{y+2} d\theta$$

$$= \frac{n-y}{y+2} \int_0^n (1-\theta)^{n-y-1} \theta^{y+2} d\theta$$

$$= \frac{n-y}{y+2} I_{n-y-1, y+2}$$

$$= \frac{(n-y)}{(y+2)} \left(\frac{n-y-1}{y+3} \right) \left(\frac{n-y-2}{y+4} \right) \dots \frac{1}{n+1} \int_0^n \theta^{n+1} d\theta$$

$$= \frac{(n-y)! (y+1)!}{(n+2)!}$$

$$(n+1) \binom{n}{y} I_{n-y, y+1} = (n+1) \frac{n!}{y!(n-y)!} \cdot \frac{(n-y)! (y+1)!}{(n+2)!}$$

$$= \frac{y+1}{n+2}$$

When $y=0$ in n trials

$$P(\tilde{Y}=1|y) = \frac{y+1}{n+2}$$

When $y=n$

$$P(\tilde{Y}=1|y) = \frac{n+1}{n+2}$$

The posterior is actually a compromise between the prior and the data information.

$$E(u) = \int u P(u) du$$

$$= \iint_{uv} u P(u, v) du dv$$

$$\int_u \int_v u P(u|v) \cdot P(v) du dv$$

$$\int_v P(v) \int_u u P(u|v) du$$

$$= \int_v E(u|v) \cdot P(v) dv$$

$$= E(E(u|v))$$

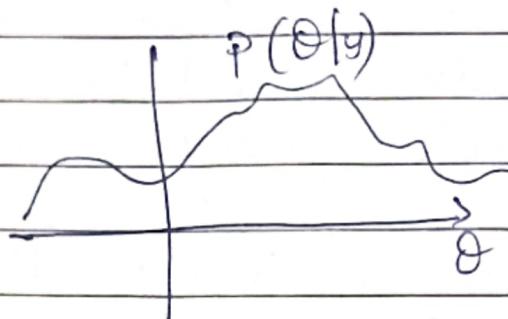
Similarly, $E(\theta) = E(E(\theta|y))$

$$Var(\theta) = E(Var(\theta|y)) + Var(E(\theta|y))$$

Expected variance of posterior Variance of the expected posterior

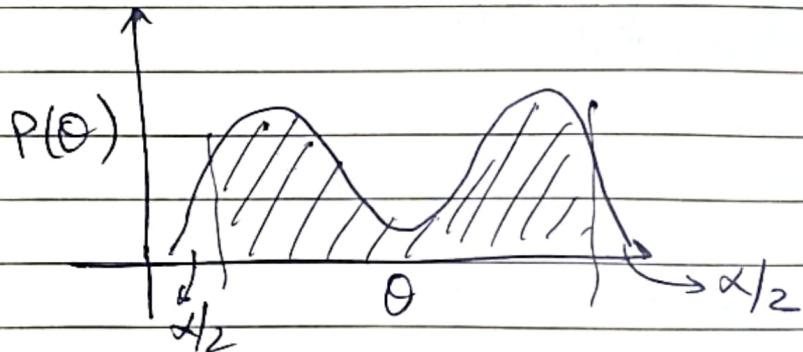
How to summarize the posterior inference

- Posterior probability distribution $P(\theta|y)$
 - Contains all the information about the current parameter θ .
- For Multiple parameters we may use a contour plot.



- Numerical summaries like mean, median and mode

- Variation, Interquartile range using Box plot.
- Posterior quantiles and intervals
 - Central interval of a posterior probability. $\rightarrow 100(1-\alpha)\%$ interval.



- Highest probability density

The part where the frequency of inner part must be larger than all the outside region.

Priors

2 basic interpretations of prior distribution.

1) Population Interpretation

- Prior distribution represents a population of possible parameter values from which the current parameter value is drawn

$$E(\beta) = \frac{\alpha}{\alpha+b} \quad \text{Var}(\beta) = \frac{\alpha b}{(\alpha+b)(\alpha+b-1)}$$

Date _____
Page _____

2) Stage of knowledge representation

- There is an actual prior distribution
- Current value represents a random realization from this distribution.

Problem : The knowledge of entire population from which θ be drawn is not possible. Ex:- Industrial failure

⇒ Binomial example with different priors

$$f(y|\theta) \propto \theta^y (1-\theta)^{n-y}$$

If we assume prior also to be of the same form.

$$f(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \propto \text{Beta}(\alpha, \beta)$$

This is similar to the probability of $\alpha-1$ successes in $\alpha+\beta-2$ trials.

α, β are hyperparameters that are assumed.

$$\begin{aligned} \text{Posterior: } P(\theta|y) &\propto P(y|\theta) \cdot P(\theta) = \theta^{y+\alpha-1} (1-\theta)^{n-y+\beta-1} \\ &= \text{Beta}(y+\alpha, n-y+\beta) \end{aligned}$$

→ When prior & posterior are of the same form they are called Conjugate prior

$$E(\theta|y) = \frac{y+\alpha}{n+\alpha+\beta}$$

$$\begin{aligned} V(\theta|y) &= \frac{(\alpha+y)(\beta+n-y)}{(\alpha+\beta+n)^2 (\alpha+\beta+n-1)} \\ &= \frac{E(\theta|y)(1-E(\theta|y))}{\alpha+\beta+n-1} \end{aligned}$$

When y, n tends to be larger for fixed α, β

$$E(\theta|y) \rightarrow \frac{\alpha}{n}$$

$$\text{Var}(\theta|y) \rightarrow \frac{1}{n} \left(\frac{\alpha}{n} \right) \left(1 - \frac{\alpha}{n} \right) \rightarrow 0 \text{ as } n \rightarrow \infty$$

Formal distribution of a conjugate prior

If F is a class of sampling distribution $f(y|\theta)$ and P is a class of prior distribution θ then class P is conjugate for F if

$$f(\theta|y) \in P \quad \text{if } f(\theta|y) \in P \text{ and } f(\theta) \in P$$

We want P to be of same form as F Then it is called natural conjugate prior.

Exponential family of distribution

Class F is a family of exponential distribution if all its members have the form

$$f(y_i|\theta) = f(y_i) g(\theta) e^{\phi(\theta) u(y_i)}$$

— Not necessary that θ is single parameter.

If θ is multparameter then,

$$f(y_i|\theta) = f(y_i) g(\theta) e^{\phi(\theta)^T u(y_i)}.$$

- $\phi(\theta)$ & $u(y)$ must be of same dimension.
- $\phi(\theta)$ is called natural parameter of F

Suppose a sequence of observation is made $Y = \{y_1, y_2, \dots, y_n\}$

$$L = P(Y|\theta) = \left[\prod_{i=1}^n f(y_i) \right] (g(\theta))^n \exp \left(\phi(\theta)^T \sum_{i=1}^n u(y_i) \right)$$

$$\propto (g(\theta))^n \exp \left(\phi(\theta)^T \sum_{i=1}^n u(y_i) \right)$$

$T(u) = \sum_{i=1}^n u(y_i)$ is called "sufficient statistic"

When y_n tends to

$$f(\theta|y) \propto g(\theta)^n e^{\phi^T(\theta) T}$$

This is the posterior.

$$f(\theta|y) \propto g(\theta)^{n+m} e^{\phi^T(\theta) (T + T(y))}$$

which is of same form as $f(\theta)$.

Non-Informative Prior

When prior distribution has no population basis. They can be difficult to construct. Hence it is necessary to derive prior that play very less role in the posterior distribution.

Such priors are called "Reference prior distribution". The prior density is non informative/flat or diffuse.

Idea is, let the data speak for itself.

Jeffreys' Invariance- Principle.

- Used to define non-informative prior distribution
- Based on one-to-one transformation of the parameters

$$\phi = h(\theta) \Rightarrow \theta = h^{-1}(\phi)$$

- By transformation the prior density is equivalent in terms of expressing the same beliefs to the prior density on ϕ .

$$f(\phi) = f(\theta) \left| \frac{d\theta}{d\phi} \right| = f(\theta) \left| h^{-1}(\phi) \right|$$

- Any rule for determining the prior density $f(\theta)$ should yield equivalent result if applied to the transformed variable, i.e. $f(\phi)$ which is computed by determining $f(\theta)$

$$f(y|\phi) = f(y|\theta) f(\theta)$$

- Considering this principle the non-uniform prior density is given as $f(\theta) \propto |J(\theta)|^{1/2}$, where $J(\theta)$ is the Fisher information for θ .

$$J(\theta) = \mathbb{E} \left(\left(\frac{d \log P(y|\theta)}{d\theta} \right)^2 \mid \theta \right)$$

- Gives a measure of the amount of the information a parameter carries about the likelihood function.

- Helps quantify the sensitivity of the model parameters to changes in data distribution.

$$J(\theta) = -E \left(\frac{d^2 \log P(y|\theta)}{d\theta^2} \mid \theta \right)$$

$$E \left(\left(\frac{d \log P(y|\theta)}{d\theta} \right) \mid \theta \right) = \int \left[\frac{\partial}{\partial \theta} \log P(y|\theta) \right] f(y|\theta) dy$$

$$y \int \frac{1}{p(y|\theta)} \frac{\partial}{\partial \theta} f(y|\theta) \cdot f(y|\theta) dy$$

$$\int y \frac{\partial}{\partial \theta} f(y|\theta) dy = \underbrace{\frac{\partial}{\partial \theta} \int f(y|\theta) dy}_1 = 0$$

$$\frac{d^2 \log f(y|\theta)}{d\theta^2} = \frac{d}{d\theta} \left(\frac{1}{p(y|\theta)} \right) \frac{\partial}{\partial \theta} \left(\frac{1}{p(y|\theta)} \right) = \frac{1}{p(y|\theta)} \frac{d^2 p(y|\theta)}{d\theta^2}$$

$$= \frac{1}{p(y|\theta)} \frac{d^2 p(y|\theta)}{d\theta^2} - \left(\frac{1}{p(y|\theta)} \cdot \frac{d}{d\theta} p(y|\theta) \right)^2$$

$$= \frac{1}{p(y|\theta)} \frac{d^2 p(y|\theta)}{d\theta^2} - \left(\frac{d}{d\theta} \log p(y|\theta) \right)^2$$

Taking expectation on both side.

$$E \left(\frac{d^2 \log p(y|\theta)}{d\theta^2} \right) = E \left(\frac{1}{p(y|\theta)} \frac{d^2 p(y|\theta)}{d\theta^2} \right) - E \left(\frac{d}{d\theta} \log p(y|\theta) \right)^2$$

$$\int \frac{1}{p(y|\theta)} \frac{d^2 p(y|\theta)}{d\theta^2} p(y|\theta) dy$$

$$\frac{d^2}{d\theta^2} \int p(y|\theta) dy \cdot \frac{\partial^2 p}{\partial \theta^2} = 0.$$

$$E\left(\frac{d^2 \log f(y|\theta)}{d\theta^2}\right) = -E\left(\left(\frac{d}{d\theta} \log f(y|\theta)\right)^2\right| \theta)$$

Q Prove that Jeffreys prior is invariant to parametral

$$p(\theta) = \{J(\theta)\}^{1/2}$$

$$\hookrightarrow J(\phi) \text{ at } \theta = h^{-1}(\phi)$$

$$\text{show } p(\phi) = f(\theta) \left| \frac{d\theta}{d\phi} \right|$$

$$J(\phi) = -E\left(\frac{d^2 \log f(y|\theta)}{d\phi^2}\right| \phi)$$

$$= -E\left(\frac{d^2}{d\theta^2} \log f(y|\theta = h^{-1}(\phi))\right| \theta) \left(\frac{d\theta}{d\phi}\right)^2$$

$$= J(\theta) \left| \frac{d\theta}{d\phi} \right|^2$$

$$(J(\phi))^{1/2} = (J(\theta))^{1/2} \left| \frac{d\theta}{d\phi} \right|$$

Q Take a binomial distribution and find Jeffreys prior

$$J(\theta) = -E\left(\frac{d^2}{d\theta^2} \log f(y|\theta)\right| \theta)$$

$$-E\left(\frac{d^2}{d\theta^2} \left[\binom{n}{y} \theta^y (1-\theta)^{n-y} \right]\right| \theta)$$

$$-E\left(\binom{n}{y} \frac{d^2}{d\theta^2} \left[\theta^y (1-\theta)^{n-y} \right]\right| \theta)$$

$$\binom{n}{y} \frac{d}{d\theta} \left(y\theta^y (1-\theta)^{n-y} + \theta^y (n-y)(1-\theta)^{n-y-1} \right) \bigg| \theta$$

$$-E(\theta) \frac{\partial}{\partial \theta} (y\theta^{y-1}(1-\theta)^{n-y} - \theta^y (n-y)(1-\theta)^{n-y-1}) | \theta$$

$$-E(\theta) \frac{\partial}{\partial \theta} \left[y(y-1)\theta^{y-2}(1-\theta)^{n-y} - y\theta^{y-1}(n-y)(1-\theta)^{n-y-1} \right] | \theta$$

$$J(\theta) = -E \left[\log \left(\frac{y}{\theta} \right) + y \log \theta + (n-y) \log (1-\theta) \right] | \theta$$

$$-E \left(\frac{\partial}{\partial \theta} \left(\frac{y}{\theta} - \frac{(n-y)}{1-\theta} \right) | \theta \right)$$

$$-E \left(\frac{\partial}{\partial \theta} \left(\frac{y}{\theta} \right) + \frac{(n-y)}{(1-\theta)^2} | \theta \right)$$

$$- \left(-\frac{1}{\theta^2} E(y) - \frac{1}{(1-\theta)^2} E(n-y) \right)$$

$$- \left(-\frac{n\theta}{\theta^2} - \frac{1}{(1-\theta)^2} (n-n\theta) \right)$$

$$\frac{n}{\theta} + \frac{n}{1-\theta} = \frac{n}{\theta(1-\theta)}$$

$$f(\theta) \propto \sqrt{\frac{n}{\theta(1-\theta)}} = n^{1/2} \theta^{-1/2} (1-\theta)^{-1/2}$$

$$\propto \text{Beta}(1/2, 1/2)$$

For two cases all principles (for generating non-informative prior) will give similar results.

i) If $f(y)$ is such that $f(y-\theta | \theta)$ is free of θ &

Then $y-\theta$ is called as a pivotal quantity

θ is called a location parameter

b) If $f(y)$ is such that

$$f(y-\theta|y) \propto f(\theta) f(y-\theta|\theta)$$

↳ constant.

b) If $f(y)$ is such that $f(y|\theta)$ is free of θ and y then
 $\frac{y}{\theta} = u$ is also a pivotal quantity
 θ in this case is a scale parameter.

Problem with non informative prior: It is difficult to find some prior which is flat. The cdf may not be equal to 1 in all case.

Weakly informative prior: Proper prior with some general information.

Ex Money in the wallet \rightarrow May range from 100 to 3000
 So we can consider a normal distribution with mean 500

⇒ CASE STUDY

A study was conducted to know what is the prob. of girl birth and boy birth in case the mother is suffering from Placenta Previa.

Proportion of girl birth in the population is 0.485
 For sample of Placenta Previa out of 980 samples,
 437 were girls. $P(\text{girl}) = 0.485$

Sol $P(\text{girl} \mid \text{Mother has placenta previa}) = \frac{437}{980}$

θ = prob. of a girl birth in patients with PP
 To find $E(\theta)$ and $\text{Var}(\theta)$

$$P(\theta|y, n) \rightarrow \text{Posterior}$$

$$\hookrightarrow P(\theta) \cdot P(y|\theta, n)$$

$$P(y|\theta, n) \propto \theta^y (1-\theta)^{n-y}$$

$$P(\theta) \propto \text{Beta}(1, 1)$$

$$P(\theta) \cdot P(y|\theta, n) \propto \theta^y (1-\theta)^{n-y}$$

$$= \text{Beta}(y+1, n-y+1)$$

$$E(\text{Beta}(y+1, n-y+1)) = \frac{\alpha(y+1)}{n+2} = 0.446$$

$$\text{Var}(\text{Beta}(y+1, n-y+1)) = \frac{\alpha(y+1)(n-y+1)}{(n+2)^2(n+3)}$$

$$= 251 \times 10^{-6}$$

$$S.D \approx 0.016$$

Q Suppose you have $\text{Beta}(4, 4)$ prior distribution on the probability θ that a coin will yield a head when spun in a specified manner. The coin is independently spun 10 times and head appears fewer than 3 times. You are not told how many heads were seen only the number is less than 3. Calculate the posterior density (exactly) upto a proportionality constant.

Sol $P(\theta) \sim \text{Beta}(4, 4)$

Unobserved parameter $\rightarrow \theta$

Observation $\rightarrow y < 3$

\hookrightarrow no. of heads in 10 trials

$P(\theta | y < 3) \rightarrow$ To find posterior.

$$P(\theta | y < 3) \propto P(y < 3 | \theta) \cdot P(\theta).$$

$$= P(y=0|\theta) P(\theta) + P(y=1|\theta) \cdot P(\theta) + P(y=2|\theta) P(\theta).$$

$$= \left[\binom{n}{0} \theta^0 (1-\theta)^n + \binom{n}{1} \theta^1 (1-\theta)^{n-1} + \binom{n}{2} \theta^2 (1-\theta)^{n-2} \right] P(\theta)$$

$$= \left[(1-\theta)^n + n \theta (1-\theta)^{n-1} + \frac{n(n-1)}{2} \theta^2 (1-\theta)^{n-2} \right] P(\theta)$$

$$= (1-\theta)^n + 10 \theta (1-\theta)^9 + 45 \theta^2 (1-\theta)^8 \theta^3 (1-\theta)^3$$

$$= \theta^3 (1-\theta)^n ((1-\theta)^2 + 10 \theta (1-\theta) + 45 \theta^2).$$

$$= \theta^3 (1-\theta)^n (1 + \theta^2 - 2\theta + 10\theta - 10\theta^2 + 45\theta^2)$$

$$= \theta^3 (1-\theta)^n (36\theta^2 + 8\theta + 1)$$

Q

Mind Game

One person is guessing random numbers from 1 to 10 20 times and based on that the other person will guess the 10 samples whether they are greater than 5 or less than equal to 5.

Sol'

$$\text{Observation} = Y = \{y_1, y_2, \dots, y_{20}\}$$

where $y_i = 0$ if count < 5

1 otherwise.

$$P(\tilde{Y} | y) = \int_0^1 P(\tilde{Y}, \theta | y)$$

$$= \int_0^1 P(\tilde{Y} | \theta) \cdot P(\theta | y)$$

$$= \int_0^1 \theta P(\theta | y)$$

$$= \int_0^1 \theta \cdot \frac{P(y|\theta) \cdot P(\theta)}{P(y)} = \frac{\theta + 1}{n + 2}$$

→ What we did till now:

→ Assumed a likelihood distribution

- Binomial
- Normal
- Poisson
- Exponential

→ Priors : Ininformative \rightsquigarrow Conjugate .
 Non-Informative
 Weak Informative .

→ Posterior

→ Prior predictive distribution

→ Posterior prediction .

~~Multi-parameter Modeling~~

NORMAL DISTRIBUTION

$$\text{Likelihood } P(y | N(\mu, \sigma^2)) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$$

Assume Known variance but unknown mean μ .

$$\text{Prior (Conjugate)} \rightarrow e^{A\theta^2 + B\theta + C} \sim N(\mu_0, \tau_0^2)$$

Since the likelihood is exponential
 $(\mu_0 \text{ & } \tau_0 \text{ are hyperparameters})$ so the prior will also be like this .

$$f(\theta) = \frac{1}{\sqrt{2\pi}\tau_0} e^{-\frac{1}{2\tau_0^2}(\theta - \mu_0)^2}$$

→ Prior distⁿ

Posterior

$$f(\theta|y) \propto f(y|\theta) \cdot f(\theta)$$

where $\theta = \mu$

$$= \frac{1}{2\pi\sigma\tau_0} e^{-\frac{1}{2\sigma^2}(y-\theta)^2 - \frac{1}{2\tau_0^2}(\theta - \mu_0)^2}$$

$$\exp\left(-\frac{1}{2}\left[\frac{(y-\theta)^2}{\sigma^2} + \frac{(\theta - \mu_0)^2}{\tau_0^2}\right]\right)$$

$$= \frac{1}{2\sigma^2\tau_0^2} \left[(\tau_0 y - \tau_0 \theta)^2 + (\theta - \mu_0)^2 \right]$$

$$= \frac{1}{2\sigma^2\tau_0^2} \left(\tau_0^2 y^2 + \tau_0^2 \theta^2 - 2\tau_0^2 y\theta + \sigma^2 \theta^2 + \sigma^2 \tau_0^2 - 2\sigma^2 \theta \tau_0 \right)$$

$$= \frac{1}{2\sigma^2\tau_0^2} \left[\tau_0^2 y^2 + \tau_0^2 \theta^2 - 2\tau_0^2 y\theta + \sigma^2 \theta^2 + \sigma^2 \mu_0^2 - 2\sigma^2 \theta \tau_0 \right]$$

$$= \frac{1}{2\sigma^2\tau_0^2} \left[(\tau_0^2 + \sigma^2) \theta^2 - 2(\tau_0^2 y + \sigma^2 \theta) \theta + (\tau_0^2 y^2 + \sigma^2 \mu_0^2) \right]$$

$$= \frac{1}{2\sigma^2\tau_0^2(\tau_0^2 + \sigma^2)} \left(\theta^2 - 2(\tau_0^2 y + \sigma^2 \theta) \theta + \frac{\tau_0^2 y^2 + \sigma^2 \mu_0^2}{\tau_0^2 + \sigma^2} \right)$$

let $\frac{\tau_0^2 y + \sigma^2 \theta}{\tau_0^2 + \sigma^2}$ as M_1 and $\frac{\tau_0^2 y^2 + \sigma^2 \mu_0^2}{\tau_0^2 + \sigma^2}$ as Z

$$= \frac{1}{2\sigma^2\tau_0^2(\tau_0^2 + \sigma^2)} (\theta^2 - 2\theta M_1 + M_1^2 - M_1^2 + Z)$$

let $\frac{\sigma^2 \tau_0^2}{\tau_0^2 + \sigma^2}$ as T

$$\exp \left(-\frac{1}{2\sigma^2} \left[\frac{(\theta - \mu_1)^2}{\sigma^2} + \frac{\mu_1^2 - z}{\sigma^2} \right] \right)$$

$$\propto \exp \left(-\frac{1}{2} \frac{(\theta - \mu_1)^2}{\sigma^2} \right)$$

Posterior: $\sim N(\mu_1, \sigma^2)$

MULTI PARAMETER MODEL

\Rightarrow Concept of "Nuisance parameter"

- Unknown ~~or~~ and unobserved quantities common in practical problem.
- Out of these one or few of these quantity is relevant.

- Obtain a marginal probability distribution
 - First obtain the joint probability distribution over all unknowns
 - Integrate the distribution over the unknowns that are not of interest.
 - These unknowns are the nuisance parameters

\Rightarrow Averaging over the nuisance parameters

$f(\theta_1, \theta_2 | y) \rightarrow$ Given
Suppose, We want to find $f(\theta_1 | y)$

$$f(\theta_1 | y) = \int_{\theta_2} f(\theta_1, \theta_2 | y) d\theta_2$$

$$= \int_{\Omega_2} f(y|\Omega_1, \Omega_2) \cdot f(\Omega_1, \Omega_2) d\Omega_2$$

Q

Consider 2 coins C_1 and C_2 with the following characteristic
 $P(\text{Heads} | C_1) = 0.6$ and $P(\text{Heads} | C_2) = 0.4$. Choose one of the coin as random and assume spinning it repeatedly. Now we observe that first two spin results in "Tails". Now what is expectation of the number of additional spin until head shows up?

Q.21

Expected no of trial given coin 1 = $1/p$.

$$P(C_1 | TT)$$

$$\mathbb{E}(\mathbb{E}(N | C, TT))$$

~~Normal~~ \Rightarrow Normal Data with non-informative priors.

- Consider a vector y of n independent observations from a univariate normal distribution $N(\mu, \sigma^2)$
- Consider a non-informative prior $f(\mu, \sigma^2) \propto (\sigma^2)^{-1}$
- Finding the joint posterior density

$$f(\mu, \sigma^2 | y) \propto f(y | \mu, \sigma^2) \cdot f(\mu, \sigma^2)$$

$$\propto \left[\prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - \mu)^2} \right) \right] \left(\frac{1}{\sigma^2} \right)$$

$$\propto (\sigma)^{-n-2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right)$$

$$\propto (\sigma)^{-n-2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - \mu)^2\right)$$

$$\propto (\sigma)^{-n-2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y})^2 + (\bar{y} - \mu)^2 + 2(y_i - \bar{y})(\bar{y} - \mu)\right)$$

$$= (\sigma)^{-n-2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\bar{y} - \mu)^2 + 2 \sum_{i=1}^n y_i \bar{y} - \sum_{i=1}^n y_i \mu\right)$$

$$= (\sigma)^{-n-2} \exp\left(-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \mu)^2 \right]\right)$$

$$\text{Let } S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$= (\sigma)^{n-2} \exp\left(-\frac{1}{2\sigma^2} \left((n-1)S^2 + n(\bar{y} - \mu)^2 \right)\right) \quad \textcircled{1}$$

Sufficient Statistic : The statistic that we need to sufficiently define the distribution in this case
 S^2 and \bar{y} are sufficient statistic.

$f(\mu | y, \sigma^2)$, assuming σ^2 is known.

$$\begin{aligned} f(\mu | y, \sigma^2) &\propto f(y | \mu, \sigma^2) \cdot f(\mu | \sigma^2) \\ &\propto f(y | \mu, \sigma^2) \cdot f(\sigma^2 | \mu) \cdot f(\mu). \end{aligned}$$

$$f(\mu | y, \sigma^2) \propto f(y | \mu, \sigma^2) f(\mu | \sigma^2)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} (n(\bar{y} - \mu)^2)\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} (\mu - \bar{y})^2\right) \sim N(\bar{y}, \frac{\sigma^2}{n}) \quad \textcircled{2}$$

$$f(\sigma^2 | y) = \int f(\mu | \sigma^2, y) d\mu$$

Marginal posterior of

distribution of the Variance.

$$(\sigma)^{-n-2} e^{-\frac{1}{2\sigma^2} (n-1)S^2} \int_{\mu} \exp\left(-\frac{1}{2\sigma^2} n(\bar{y} - \mu)^2\right) d\mu$$

~~Exp~~ $\propto \sigma^{(n-1)} \cdot (2\sigma^2)^{-\frac{n}{2}} \cdot (2\sigma^2)^{\frac{n}{2}} \cdot (2\sigma^2)^{-\frac{n}{2}}$.

$$\propto \int e^{\frac{n}{2\sigma^2}(\bar{y}^2 - 2\bar{y}\mu + \mu^2)} d\mu$$

$$\propto e^{-\frac{n\bar{y}^2}{2\sigma^2}} \int e^{\frac{n\mu^2}{2\sigma^2}} d\mu$$

$$e^{-\frac{n\bar{y}^2}{2\sigma^2}}$$

$$\propto \int_{\mu} \exp\left(-\frac{1}{2\sigma^2}(\mu - \bar{y})^2\right) d\mu = \frac{\sqrt{2\pi\sigma^2}}{\sqrt{n}} \int N(\bar{y}; \frac{1}{n}) d\mu$$

$$f(\sigma^2 | y) \propto (\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}(n-1)s^2\right) \sqrt{\frac{2\pi\sigma^2}{n}}$$

$$\propto (\sigma^2)^{-\frac{n+1}{2}} \exp\left(-\frac{1}{2\sigma^2}(n-1)s^2\right) \quad \text{--- (3)}$$

Scaled Inverse χ^2 density.

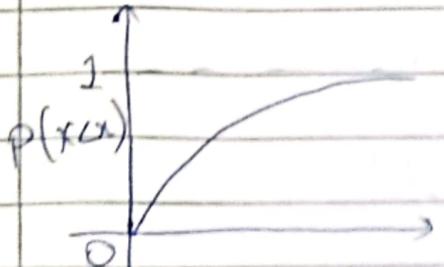
Q Find $f(\mu | y)$.

?? How to sample from joint posterior density?

② First draw σ^2 from ③

Note \Rightarrow Techniques for generating sample that follow a distribution

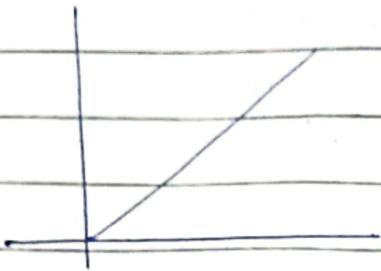
1) Inverse Transform



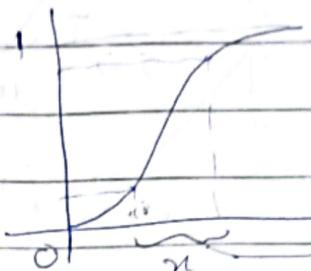
$$P(x < u) = \sum_{x=-\infty}^u P(x)$$

Take a number b/w 0 & 1 and find the corresponding value of x .

For Uniform



For Normal

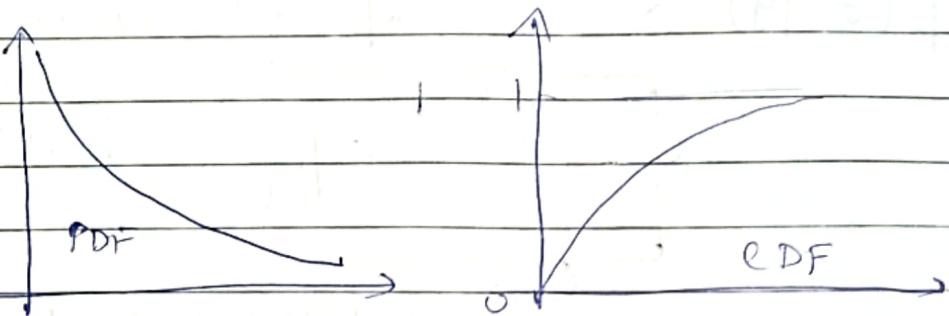


For Exponential distribution

For Exponential distribution

$$P(X=x) = \lambda e^{-\lambda x}$$

$$P(X \leq x) = 1 - e^{-\lambda x}$$

Generate a random number $U \in [0, 1]$

$$U = 1 - e^{-\lambda x}$$

~~$$U = e^{-\lambda x}$$~~

$$e^{-\lambda x} = 1 - U$$

$$-\lambda x = \ln(1 - U)$$

$$x = -\frac{1}{\lambda} \ln(1 - U)$$

In this technique we will take a random no. and for the CDF we find x corresponding to y by finding inverse transform.

2) Acceptance / Rejection Technique

- Consider an alternate distribution(A) whose inverse CDF is solvable.
- Generate samplesⁿ based on this alternate distribution.

$\frac{f_0(x)}{f_A(x)}$ if this ratio is greater than a threshold
 then accept otherwise reject.

(b) For a given σ^2 , draw μ from ②

Now we have $\mu \sim \sigma^2$ which will follow ①

$$P(\tilde{y}|y) \propto \iint_{\sigma^2 \mu} P(\tilde{y}, \mu, \sigma^2 | y) d(\sigma^2) d\mu$$

$$= \iint_{\sigma^2 \mu} P(\tilde{y}|\mu, \sigma^2, y) \cdot P(\mu, \sigma^2 | y) d(\sigma^2) d\mu.$$

$$\iint_{\sigma^2 \mu} N(\tilde{y}, \mu, \sigma^2) \quad \text{derived using ①}$$

The integration is difficult to calculate so we will use the Sampling techniques.

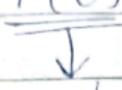
- first draw samples from the joint posterior distribution $P(\mu, \sigma^2 | y)$
- Given μ and σ^2 ~~draw samples from~~ find $N(\tilde{y}, \mu, \sigma^2)$
 ie $\tilde{y} \sim N(\mu, \sigma^2)$

Multinomial model for categorical Data

- ↓
- The number of
- Data for which each observation has k possible outcomes.
 - If γ is the vector of the counts of number of observations of each outcome
- $$f(\gamma | \theta) \propto \prod_{j=1}^k \theta_j^{y_j} ; \sum_{j=1}^k \theta_j = 1, \sum_{j=1}^k y_j = n$$

constraints.

$$P(O|y) \propto P(y|O) \cdot P(O)$$



for binomial conjugate prior

For Multinomial conjugate prior is

Dirichlet Distribution.

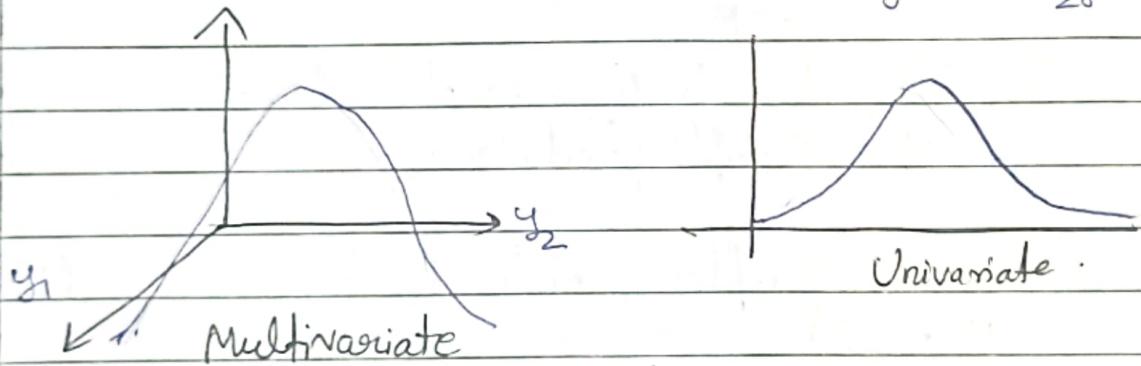
$$D(\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_k) \propto \prod_{j=1}^k \theta_j^{\alpha_j-1}$$

$$P(O|y) \propto \prod_{j=1}^k \theta_j^{\alpha_j + y_j - 1} \sim D(\alpha_1 + y_1, \alpha_2 + y_2, \dots, \alpha_k + y_k)$$

Multivariate Normal Model

Likelihood

$$\text{For univariate case: } f(y; \mu, \sigma^2) \propto \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} (y - \mu)^2\right)$$



Suppose y is an observable quantity with d components. $y \in \mathbb{R}^d$

$$y | \mu, \Sigma \sim N(\mu, \Sigma)$$

$$\propto \left| \Sigma \right|^{1/2} \exp\left(-\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu)\right)$$

$$E((y_i - \mu_i)(y_j - \mu_j))$$

$$\Sigma = \begin{bmatrix} \text{Cov}(y_1, y_1) & \text{Cov}(y_1, y_2) & \dots & \text{Cov}(y_1, y_d) \\ \text{Cov}(y_2, y_1) & \text{Cov}(y_2, y_2) & \dots & \text{Cov}(y_2, y_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(y_d, y_1) & \text{Cov}(y_d, y_2) & \dots & \text{Cov}(y_d, y_d) \end{bmatrix}$$

$$\sigma^2$$

\Leftrightarrow Symmetric

\Leftrightarrow Positive definite matrix (all eigen values)

$$x^T \Sigma x = 0 \Rightarrow \text{The soln to this}$$

is unique and it becomes quadratic while solution

Suppose you are given n observations, y_1, y_2, \dots, y_n

Likelihood: $f(y_1, y_2, \dots, y_n | \mu, \Sigma) \propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \mu)^T \Sigma^{-1} (y_i - \mu)\right)$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}^T \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1d} \\ \vdots & \vdots & & \vdots \\ \Sigma_{d1} & \Sigma_{d2} & \dots & \Sigma_{dd} \end{bmatrix} \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_d \end{bmatrix}$$

$$\begin{aligned} & \hat{y}_1 (\hat{y}_1 \Sigma_{11} + \hat{y}_2 \Sigma_{21} + \dots + \hat{y}_d \Sigma_{d1}) \\ & + \hat{y}_2 (\hat{y}_1 \Sigma_{12} + \hat{y}_2 \Sigma_{22} + \dots + \hat{y}_d \Sigma_{d2}) \\ & + \dots \\ & + \hat{y}_d (\hat{y}_1 \Sigma_{d1} + \hat{y}_2 \Sigma_{d2} + \dots + \hat{y}_d \Sigma_{dd}) \end{aligned}$$

$$= \text{Tr}(\Sigma^{-1} S_0), \quad \text{where } S_0 = \sum_{i=1}^n (y_i - \mu)(y_i - \mu)^T$$

$$f(y_1, y_2, \dots, y_n | \mu, \Sigma) \propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \text{Tr}(\Sigma^{-1} S_0)\right)$$

Conjugate Analysis

→ Assume we do not know μ, Σ is known

→ Prior for μ is assumed to be normal $\mu \sim N(\mu_0, \Lambda_0) \in \mathbb{R}^d$

$$\begin{aligned} \text{posterior } f(\mu | y, \Sigma) & \propto f(y | \mu, \Sigma) f(\mu | \Sigma) \\ & \propto N(y, \mu, \sigma^2 \Sigma) N(\mu, \mu_0, \Lambda_0) \end{aligned}$$

$$\propto \exp \left(-\frac{1}{2} \left[(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda}_0^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) + \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \right)$$

$$(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda}_0^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)$$

$$\underbrace{\boldsymbol{\mu}^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}}_{\text{Quadratic}} - \boldsymbol{\mu}^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0$$

$$\quad \quad \quad \downarrow \text{Scalar quantities}$$

$$- \boldsymbol{\mu}^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0$$

$$= \underbrace{\boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}}_{\substack{\text{Quadratic} \\ \text{in terms of } \boldsymbol{\mu}}} - 2 \underbrace{\boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0}_{\substack{\downarrow \\ \text{Linear in terms} \\ \text{of } \boldsymbol{\mu}}} + \underbrace{\boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0}_{\substack{\downarrow \\ \text{Constant in terms} \\ \text{of } \boldsymbol{\mu}}} \quad - \textcircled{1}$$

$$\sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) = \textcircled{2}$$

$$= \sum_{i=1}^n \left[\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - 2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}_i + \mathbf{y}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{y}_i \right]$$

$$= n \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - 2n \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{y}} + \text{constant} \quad - \textcircled{11}$$

~~cancel~~ from $\textcircled{1} + \textcircled{11}$

$$\boldsymbol{\mu}^T \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu} + n \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - 2 \boldsymbol{\mu}^T (\boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + n \boldsymbol{\Sigma}^{-1} \bar{\mathbf{y}}) + \text{constant}$$

$$\boldsymbol{\mu}^T (\boldsymbol{\Lambda}_0^{-1} + n \boldsymbol{\Sigma}^{-1}) \boldsymbol{\mu} - 2 \boldsymbol{\mu}^T (\boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + n \boldsymbol{\Sigma}^{-1} \bar{\mathbf{y}}) + \text{const.}$$

$$\boldsymbol{\Lambda}_n^{-1} = \boldsymbol{\Lambda}_0^{-1} + n \boldsymbol{\Sigma}^{-1}$$

$$\boldsymbol{\Lambda}_n^{-1} \boldsymbol{\mu}_n = (\boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + n \boldsymbol{\Sigma}^{-1} \bar{\mathbf{y}})$$

$$\boldsymbol{\mu}_n = \boldsymbol{\Lambda}_n (\boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + n \boldsymbol{\Sigma}^{-1} \bar{\mathbf{y}})$$

$$= (\boldsymbol{\Lambda}_0^{-1} + n \boldsymbol{\Sigma}^{-1})^{-1} (\boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + n \boldsymbol{\Sigma}^{-1} \bar{\mathbf{y}})$$

$$\sim N(\mu; \mu_n, \Lambda_n)$$

Posterior distribution of μ in a multivariate normal distribution with known Σ . Assuming prior $\mu \sim N(\mu_0, \Lambda_0)$

$$\begin{aligned} f(\mu | y, \Sigma) &\sim N(\mu_n, \Lambda_n) \\ \mu_n &= (\Lambda_0^{-1} + n \Sigma^{-1}) (\Lambda_0^{-1} \mu_0 + n \Sigma^{-1} \bar{y}) \\ \Lambda_n &= (\Lambda_0^{-1} + n \Sigma^{-1})^{-1} \end{aligned}$$

Posterior, conditionals and marginal distribution of a subvectors of μ with known Σ

We know, $y \in \mathbb{R}^d$, $\mu \in \mathbb{R}^d$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{bmatrix} \begin{matrix} \rightarrow \text{known} \\ \rightarrow \text{known} \\ \rightarrow \text{unknown} \\ \rightarrow \text{unknown} \end{matrix}$$

we want to find their distribution

- Marginal posterior of a subset of parameter (say $\mu^{(1)}$)
- Given a known subset ($\mu^{(2)}$)
- A dataset y
- Q We want $\mu^{(1)} | \mu^{(2)}, y$

Result is a multivariate normal distribution

$$\mu^{(1)} | \mu^{(2)}, y \sim N(\mu_n^{(1)} + \beta^{1/2} (\mu^{(2)} - \mu_n^{(2)}), \Lambda^{1/2})$$

$\mu_n^{(1)} \rightarrow$ Posterior mean vector

$\mu_n^{(1)} \rightarrow$ Appropriate subvector of μ_n with similar grouping as $\mu^{(1)}$

$\mu_n^{(2)} \rightarrow$ Appropriate ... of μ_n with ... $\mu^{(2)}$

$\Lambda_n \rightarrow$ Appropriate Submatrix of the given n variance matrix

~~Block~~

$$\beta'^{1/2} = \Lambda_n^{(1,2)} (\Lambda_n^{(2,2)})^{-1}$$

$$\Lambda'^{1/2} = \Lambda_n^{(1,1)} - \Lambda_n^{(1,2)} (\Lambda_n^{(2,2)})^{-1} \Lambda_n^{(2,1)}$$

$$\Lambda' = \begin{bmatrix} \Lambda_n^{(1,1)} & \Lambda_n^{(1,2)} \\ \Lambda_n^{(2,1)} & \Lambda_n^{(2,2)} \end{bmatrix}$$

Posterior Predictive Distribution of a new data

Suppose we want to find $\tilde{Y}|y$ for known variance Σ

$$f(\tilde{Y}|y) = \int f(\tilde{Y}, \mu|y, \Sigma) d\mu$$

$$\mu = \int \underbrace{f(\tilde{Y}|\mu, y, \Sigma)}_{\text{likelihood}} \cdot f(\mu|y, \Sigma) d\mu$$

$$= \int \underbrace{f(\tilde{Y}|\mu, \Sigma)}_{\mu} \cdot \underbrace{f(\mu|y, \Sigma)}_{\text{Posterior of } \mu \text{ with known } \Sigma} d\mu$$

$$\int_{\mu} N(\mu, \Sigma) \cdot N(\mu_n, \Lambda_n) d\mu$$

which is normal distribution

Rather than integrating we find the mean and the variance of posterior predictive distribution.

$$\text{We know } E(u) = E(E(u|v))$$

$$\begin{aligned} E(\tilde{Y}|y) &= E(E(\tilde{Y}|y, \mu)|y) \\ &= E(\mu|y) = \mu_n \end{aligned}$$

$$\begin{aligned}\text{Var}(\tilde{y}|y) &= \mathbb{E}(\text{var}(\tilde{y}|y, \mu))|y) + \text{Var}(\mathbb{E}(\tilde{y}|y, \mu)|y) \\ &= \mathbb{E}(\sum|y) + \text{Var}(\mu|y) \\ &= \sum + \Lambda_0\end{aligned}$$

Non-informative Prior

$$P(\mu) \propto \text{constant}.$$

The distribution generated for a non-informative prior is proper only when $n \geq d$

↑
no. of samples ↑
dimension of the vector (y)
multivariate.

Multivariate normal with unknown mean and variance

Chi square distribution when extended to multivariate it becomes wishart distribution.

$$\Sigma \sim \text{Inv-Wishart}(\Lambda_0^{-1})$$

$$\mu|\Sigma \sim N(\mu_0, \Sigma/k_0) ; k_0 \text{ is a scale parameter}$$

Joint prior $= p(\mu, \Sigma) = \text{Inv-Wishart} \rightarrow \# \text{ of prior measurements on the } \Sigma \text{ scale.}$

The joint posterior density is also a ~~is~~ Inv-Wishart distribution

Case Study - Bioassay Experiment

Scientific Problem:

- Administer various levels of dose of a drug to some batch of test animals.

- The animal response is characterized by a dichotomous outcome. (alive or dead)
- Data is of the following form
 (x_i, n_i, y_i)
 \downarrow
 i^{th} dose level No. of animals to which i^{th} dose level is given How many of the n_i animals died.
- Sample data for tests conducted on 20 animals

Dose (x_i in $\text{log}(\text{gm/ml})$)	Number of animals (n_i)	Number of deaths (y_i)
0.86	5	0
0.3	5	1
0.05	5	3
0.73	5	5

We want to establish a dose-response relation.

8.1

Modeling the dose-response relation

- Assume that the outcomes of 5 animals within each group is exchangeable.
- Reasonable to model them as independent with equal probabilities.
- Let θ_i be that probability for group i
 $\text{So } y_i | \theta_i \sim \text{Bin}(n_i, \theta_i) - ①$
 \hookrightarrow prob. of death given dose x_i
- Assume that the outcome probability of each group is independent of each other.
- Simple assumption is $\theta_1, \theta_2, \theta_3, \theta_4$ are exchangeable.

$$f(\theta_1, \theta_2, \theta_3, \theta_4) \propto 1$$

Assuming $\theta_i \sim \text{Beta}$ posterior.

Assuming, $\theta_i = \alpha + \beta x_i$

α, β are parameters of the model

- Problem is with $x_i \rightarrow \infty, \theta_i \rightarrow \infty$; But $\theta_i \in [0, 1]$
- So let $\text{logit}(\theta_i) = \alpha + \beta x_i$

$$\text{logit}(\theta_i) = \log\left(\frac{\theta_i}{1-\theta_i}\right) = \alpha + \beta x_i$$

$$\frac{\theta_i}{1-\theta_i} = e^{\alpha + \beta x_i}$$

$$\frac{1-\theta_i}{\theta_i} = e^{-(\alpha + \beta x_i)}$$

$$\text{or } \theta_i = \frac{1}{1 + e^{-(\alpha + \beta x_i)}}$$

Likelihood

for group i

$$f(y_i | \alpha, \beta, n_i, x_i) \propto \theta_i^{y_i} (1-\theta_i)^{n_i - y_i}$$

$$\text{where } \theta_i = \frac{1}{1 + e^{-(\alpha + \beta x_i)}} = \text{logit}^{-1}(\alpha + \beta x_i)$$

$$\rightarrow f(y_i | \alpha, \beta, n_i, x_i) \propto [\text{logit}^{-1}(\alpha + \beta x_i)]^{y_i} [1 - \text{logit}^{-1}(\alpha + \beta x_i)]^{n_i - y_i}$$

Joint posterior of α, β

$$f(\alpha, \beta | y_i, n_i, x_i) \propto f(y_i | \alpha, \beta, n_i, x_i) \cdot f(\alpha, \beta)$$

$$f(\alpha, \beta | y, n, x) \propto f(\alpha, \beta) f(y | \alpha, \beta, n, x)$$

$$f(\alpha, \beta | y, n, x) \propto f(\alpha, \beta) \prod_{i=1}^n f(y_i | \alpha, \beta, n, x_i)$$

As likelihood is complex so we have to rely on computation

Assuming $f(\alpha, \beta) \propto 1$ ie. constant.

For Sampling

We need to have some initial estimate of the parameter

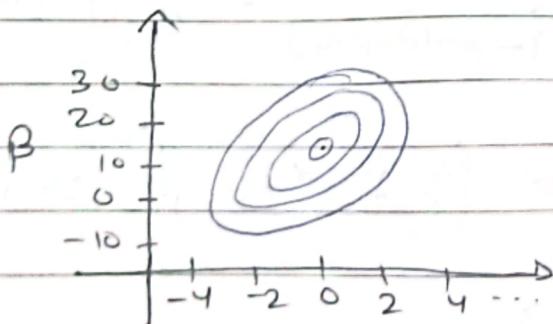
Get an estimate of α and β by fitting a logistic regression using MLE.

$$\text{Estimated } (\hat{\alpha}, \hat{\beta}) = (0.8, 7.7)$$

$$S.E(\hat{\alpha}) = 1.02$$

$$S.E(\hat{\beta}) = 4.9$$

- Consider the range $(\alpha, \beta) \in [-5, 10], [-10, 40]$
Take small values for incrementing and get the value of posterior for and get the contour plot with these grid values -



Steps to find $f(\theta)$

- Start with an initial estimate of $(\hat{\alpha}, \hat{\beta})$ by fitting with a LR using MLE
 - ↳ Logistic Regression

$$y = \frac{1}{1 + e^{\alpha + \beta x}}$$

$$\theta_i = \frac{\text{No. of deaths}(y_i)}{\text{Total number of animals}(n_i)}$$

$$f_p(\alpha, \beta | m, n, y) = \cancel{p(\alpha, \beta)} \cdot p(y | \alpha, \beta, n, n) \\ \propto \cancel{\prod_{i=1}^K p(y_i | \alpha, \beta, n_i, \theta_i)} \\ \propto \prod_{i=1}^K \theta_i^{y_i} (1 - \theta_i)^{n_i - y_i}$$

$$\alpha_{\text{new}} = \alpha_{\text{old}} + \frac{\delta \alpha}{\text{grid size}}$$

If grid size is very large then we will get fine info. about the distribution.

2) The distribution that we get is unnormalized so there is a need of normalization. Sampling should be used.

Sampling from the joint posterior

$$f_p(\alpha | y_i, n_i, \theta_i)$$

2) Compute the marginal posterior distribution of α by numerically summing over β .
 → in the discrete distribution computed over the grid

$$f_p(\alpha, \beta | y, n, \theta) \approx \sum_{\beta=1}^S f_p(\alpha, \beta | y, n, \theta)$$

2) For $S = 1 \text{ to } 1000$

Draw α^S directly from the directly computed $f_p(\alpha, \beta | y, n, \theta)$

b) Draw β^s from the discrete conditional distribution.
 $P(\beta | \alpha, y)$

c) For each sampled α^s and β^s add a ~~gray~~ random jitter ϵ
 $E(\epsilon) = 0$ and it is uniform
 $\epsilon \sim \mathcal{U}[-u, +u]$



Introduction To Bayesian Computation Techniques

Revolves around 2 steps:

- Computation of the posterior $f(\theta | y)$
- Computation of the posterior predictive distribution $f(y | \theta)$

- Standardized distributions are easier to deal with
 like normal, poisson, Exponential, Gamma

Normalized & Unnormalized

The distribution (which may be multivariate) that needs to be computed (simulated) is called Target distribution
 Let it be denoted as $f(\theta | y)$

Assume $f(\theta | y)$ to be ^{easily} computed for any value of θ
 upto a factor only involving y

We assume there is some easily computable function
 $g(\theta | y)$ which is easily computed and is unnormalized
 density of $f(\theta | y)$
 proposal distribution

$\frac{q(\theta|y)}{p(\theta|y)}$ is a constant that only depends on y

and is proportional to the posterior density.

#

Numerical Integration / Quadrature

- It refers to a method in which the integration over the continuous function is evaluated
- By computing the value of the function at finite set of points.

Two methods for numerical integration

- Simulation based (stochastic)
- Deterministic method.

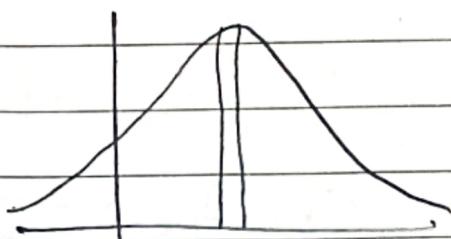
~~Simulation
Based~~

$$E(h(\theta|y)) = \int h(\theta|y) p(\theta|y) d\theta = \int h(\theta) p(\theta|y) d\theta$$

Draw posterior draws θ^s from $p(\theta|y)$

Estimate the integral

$$\frac{1}{S} \sum_{i=1}^S h(\theta^i)$$



When we integrate we need to have the sum the value of the $h(\theta)$ at that i as when we sample the $h(\theta)$ will be higher at the peak and less at the tails.

- Samples drawn are independent of each other

Deterministic Method

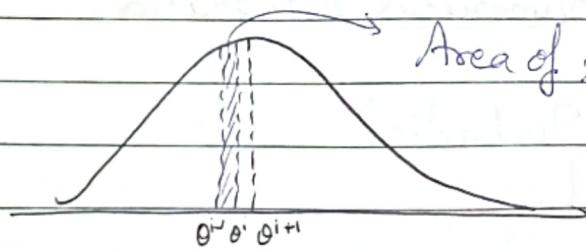
Evaluate the integrand $\int h(\theta) f(\theta|y) d\theta$ at selected points θ^s .

$$E(h(\theta|y)) = \int h(\theta) f(\theta|y) = \sum \theta^s$$

$$= \frac{1}{S} \sum_{i=1}^S \omega^{(i)} h(\theta^{(i)}) f(\theta^{(i)}|y)$$

weight given to a particular integral.

$\omega^{(i)}$ is the weight corresponds to volume of space at $\theta^{(i)}$



$$\text{Area of trapezium} = \frac{1}{2} (h(\theta^{(i-1)}) + h(\theta^{(i)})) \Delta \theta^{(i)}$$

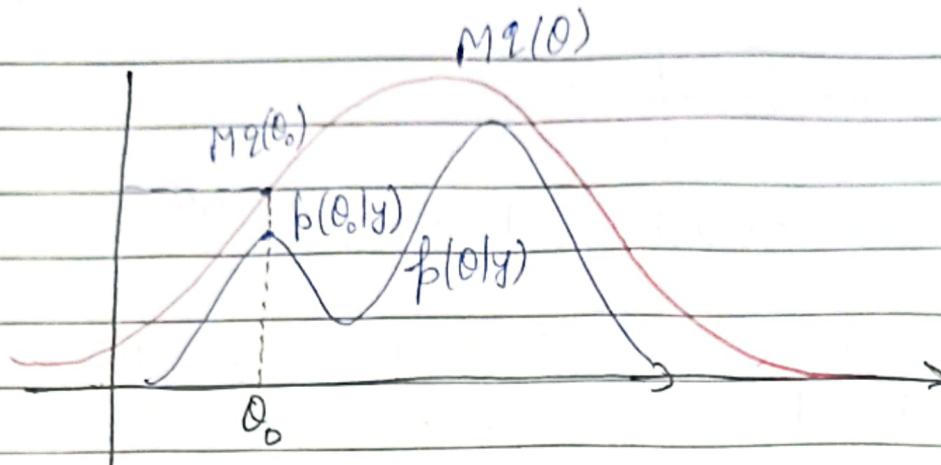
$$E(h(\theta|y)) = \frac{1}{2} \sum_{i=1}^S (h(\theta^{(i-1)}) + h(\theta^{(i)})) \Delta \theta^{(i)} f(\theta^{(i)}|y)$$

Rejection Sampling

Objective : To estimate $f(\theta|y)$

- Consider a proposal distribution $q(\theta)$ where inverse of $q(\theta)$ is easy to find
- $q(\theta)$ is defined for all θ for which $f(\theta|y)$ is defined.
- $q(\theta)$ when multiplied by a constant M is always greater than $f(\theta|y)$ $\forall \theta$ for which $f(\theta|y)$ is defined

$$M q(\theta) \geq f(\theta|y)$$



- Use $q(\theta)$ to sample a point say θ_0
- Generate a random sample u uniformly over $[0, 1]$
- Accept θ_0 with a probability $\frac{f(\theta_0|y)}{M q(\theta_0)} \geq u$
- Here θ_0 gets accepted with a higher prob. if $f(\theta_0|y)$ is closer to $M q(\theta_0)$

Problem with Rejection Sampling:

- If there is a multivariate distribution then we need to have exponentially high values ~~for~~ for θ

Advantage

- Generation of normalized distribution.

If we have $\alpha_1, \alpha_2, \dots, \alpha_n$ are generated Sample

Normalized procedure :
$$\frac{f(\alpha_i|y)}{\sum_{j=1}^n f(\alpha_j|y)}$$

Importance Sampling

- Suppose we are interested in $E(h(\theta|y))$
- Importance Sampling finds the point estimate of the distribution directly without the use of envelope fn.

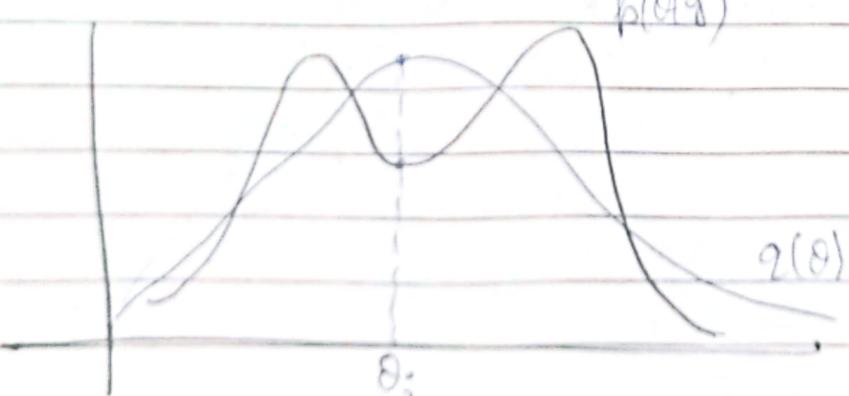
$$\begin{aligned} E(h(\theta|y)) &= \int h(\theta) f(\theta|y) d\theta \\ &= \int h(\theta) \frac{f(\theta|y)}{g(\theta)} \cdot g(\theta) d\theta \end{aligned}$$

$g(\theta)$ is a distribution whose samples are easily generated. and given θ it is easy to determine the density $p(\theta|y) \propto g(\theta|y)$.

- For each sample θ_i generated from $g(\theta)$
 - Find the ratio $\frac{f(\theta_i|y)}{g(\theta_i)} = w_i$

Generate S sample of θ from $g(\theta)$

$$E(h(\theta|y)) = \frac{1}{S} \sum_{i=1}^S h(\theta_i|y) w_i$$



For θ_i $g(\theta_i)$ is large but $w_i = \frac{f(\theta_i|y)}{g(\theta_i)}$ will be < 1 and acts as correction factor.

For the case when $f(\theta|y)$ is unnormalized.

$$p(\theta|y) = \frac{g(\theta|y)}{\int g(\theta|y) d\theta}$$

$$E(h(\theta|y)) = \frac{\int h(\theta) g(\theta|y) d\theta}{\int g(\theta|y) d\theta}$$

Suppose proposal distribution $q(\theta|y)$

$\frac{g(\theta|y)}{\int g(\theta|y) d\theta}$ is estimated using proposal distribution

$$\frac{q(\theta|y)}{\int q(\theta|y) d\theta}$$

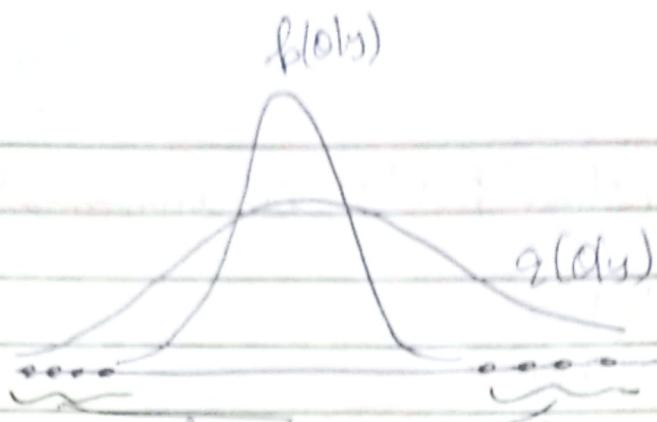
$$E(h(\theta|y)) = \frac{\int h(\theta) \left(\frac{g(\theta|y)}{q(\theta|y)} \right) q(\theta|y) d\theta}{\int \frac{g(\theta|y)}{q(\theta|y)} q(\theta|y) d\theta}$$

$$\text{Let } w_i = \frac{g(\theta|y)}{q(\theta|y)}$$

Given S samples from $q(\theta|y)$

$$\frac{\frac{1}{S} \sum_{i=1}^S h(\theta_i) w_i}{\frac{1}{S} \sum_{i=1}^S w_i} = \frac{\sum_{i=1}^S h(\theta_i) w_i}{\sum_{i=1}^S w_i}$$

Problem in determining distribution instead of point estimate
 There are many points which are sampled from $q(\theta|y)$ which are not the part of the $f(\theta|y)$ line.



These points are not belonging of form $f(O|y)$

II. Sample Importance Resampling.

$$g_i(O|y) = \frac{g_i(O|y)}{g(O|y)} g(O|y)$$

Suppose you have generated samples from $g(O|y)$

$$g(O|y) \rightarrow \{O^1, O^2, O^3, \dots, O^N\}$$

N is quite large.

Generate another set of samples

$$\tilde{g} \rightarrow \{\tilde{O}^1, \tilde{O}^2, \tilde{O}^3, \dots, \tilde{O}^N\}$$

Accept O^i with probability $\frac{w_i}{\sum w_i}$ to generate \tilde{O}^i

Steps

Take sample O^i from g

Generate a random number $\gamma \in [0, 1]$

If $\gamma < \frac{w_i}{\sum w_i}$ accept O^i to \tilde{O}^i

Else repeat it for O^i , γ

In general we have to sample without replacement. But analytically it is difficult to do this as

Most of the samples are from 1st peak so the data will be biased if we do with replacement.

Markov Chain Monte Carlo (MCMC)

- Focus the posterior density evaluations to the part of the parameter space.

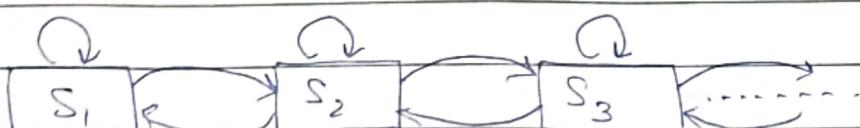
$$E_{p(\theta|y)} f(\theta) = \int f(\theta) p(\theta|y) d\theta$$

We can use the unnormalized posterior $q(\theta|y) = p(y|\theta) \cdot h(\theta)$

Grid $E_{p(\theta|y)} f(\theta) \approx \sum_{s=1}^S f(\theta^s) \frac{q(\theta^s|y)}{\sum_{s=1}^S q(\theta^s|y)}$

Works for few dimensions. As the no. of dimensions increases the no. of grids increases exponentially and most of the grid points have posterior as 0. so it is computation expensive.

Markov Chain



$$p(S_7 | S_1, S_2, S_3, \dots, S_{7-1}) = p(S_7 | S_{7-1})$$

Stationarity Property: If we run the chain multiple times, then the fraction of time we are on state S_i is same for all i .

Gibbi's Sampling

If $y \sim N(\mu, \sigma^2)$ and $\mu \& \sigma^2$ are both unknown so we can't apply markov chain directly so we will use gibbi sampling.

We try to generate samples based on 1D conditional distribution $f(\mu | \sigma^2)$ and $f(\sigma^2 | \mu)$

Algo

Sample θ_j^t from $f(\theta_j | \theta_{-j}^{t-1}, y)$

where $\theta_{-j}^{t-1} = (\theta_1^t, \dots, \theta_{j-1}^t, \theta_{j+1}^{t-1}, \dots, \theta_n^{t-1})$

If $\mu \& \sigma^2$ are highly correlated then the algo will be very slow.

For Joint pdf

For Conditional the correlation creates problem of slow convergence.

Metropolis Algorithm

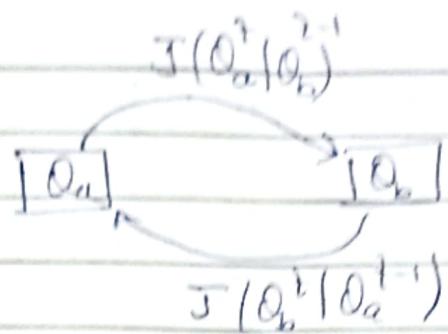
1) Starting point θ^0

2) $t = 1, 2, \dots$

a) pick a proposal θ^* from the proposal dist $J_1(\theta^* | \theta^{t-1})$

Proposal dist has to be symmetric ie.

$$J_2(\theta_a | \theta_b) = J_2(\theta_b | \theta_a) \forall \theta_a, \theta_b$$



b) Calculate acceptance ratio

$$\alpha = \frac{p(O^* | y)}{p(O^{**} | y)}$$

c) set

$$O^t = \begin{cases} O^* & \text{with prob min}(a, 1) \\ O^{**} & \text{otherwise.} \end{cases}$$

→ No need of normalization coeff. because they will cancel out while calculating α .

1 Prove that simulated series is a Markov chain which has unique stationary distribution.

→ a) Irreducible : the prob. of eventually reaching any state from any other state.

We have continuous distⁿ so from any state we can go to any other state

b) Aperiodic : return times are not periodic
holds for a random walk on any proper distⁿ.

$\pi_k = p(O^t | O^{t-u}, y) \rightarrow$ If for any k the prob. is 1
since every state has some time prob. for jumps in random walk so the prob is never 1

→ Recurrent/Non transient : prob to return to a state is 1.

- holds for a random walk on any proper distⁿ

2. Prove that the stationary distⁿ is the target distⁿ $f(\theta|y)$

Consider $\theta^{t-1} \sim p(\theta|y)$

Consider any two such points θ_a and θ_b drawn from $p(\theta|y)$ and labeled so that $f(\theta_b|y) \geq f(\theta_a|y)$

$$f(\theta^{t-1} = \theta_a, \theta^t = \theta_b) = f(\theta_a|y) \cdot J(\theta_b|\theta_a)$$

Due to this the
gr is not 1

$$\rightarrow f(\theta^{t-1} = \theta_b, \theta^t = \theta_a)$$

$$= f(\theta_b|y) J(\theta_a|\theta_b) \cdot \frac{f(\theta^t = \theta_a)}{f(\theta^{t-1} = \theta_b)}$$

$$= f(\theta_a|y) \cdot J(\theta_a|\theta_b)$$

The marginal to be in θ_a at time t & $t-1$ is same.

Metropolis-Hastings Algorithm

Generalization of metropolis algo with non symmetric proposal distⁿ.

$$\alpha = \frac{f(\theta^t|y) / J_p(\theta^*|\theta^{t-1})}{f(\theta^{t-1}|y) / J_p(\theta^t|\theta^*)}$$

$$L(\theta|y) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}}$$

Ex:

Using Gibbs sampler to sample from the posterior density of μ and σ^2 in an 1D normal distribution

Start
and
End
and
End

Posterior distribution of μ, σ^2 $f(\mu, \sigma^2 | y)$

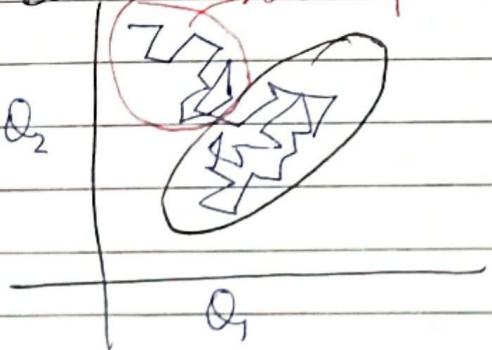
Posterior distribution of μ

Problems with Metropolis Algorithm:

Usually doesn't scale well to high dimension as there is exponentially growth of search space.

↑ In the higher dimension it is difficult to get the correct proposal distribution due to which a most of the sampled points will be rejected so the efficiency of the algorithm reduces.

(Warm-up): Remove draws from the beginning of the chain
or
Burn-in ~~warmup~~



Using several chains (To cover the distribution well)

We have to remove warmup and run the chain until they get well mixed and undistinguishable.

Sometimes in visual representation due to shrink in scale we can get false idea about convergence.

\hat{R} : Comparison of within and between variances of the chains

M chain each having N draws

Within chain variance (W)

$$W = \frac{1}{M} \sum_{m=1}^M S_m^2, \quad S_m^2 = \frac{1}{N-1} \sum_{n=1}^N (\theta_{nm} - \bar{\theta}_{..})^2$$

Between chain variance (B)

$$B = \frac{N}{M-1} \sum_{m=1}^M (\bar{\theta}_{..m} - \bar{\theta}_{..})^2$$

$$\bar{\theta}_{..m} = \frac{1}{N} \sum_{n=1}^N \theta_{nm}, \quad \bar{\theta}_{..} = \frac{1}{M} \sum_{m=1}^M \bar{\theta}_{..m}$$

$$\widehat{\text{Var}}^+(\theta|y) = \frac{N-1}{N} W + \frac{1}{N} B$$

This overestimates marginal posterior variance if the starting points are overdispersed

$$\hat{R} = \frac{\widehat{\text{Var}}^+}{W}$$

For finite N, W underestimates marginal variance

It should be very near to 1 (0.95-0.99)

if $R > 1.01$, keep sampling.

If \hat{R} close to 1, it is still possible that chain have not converged:

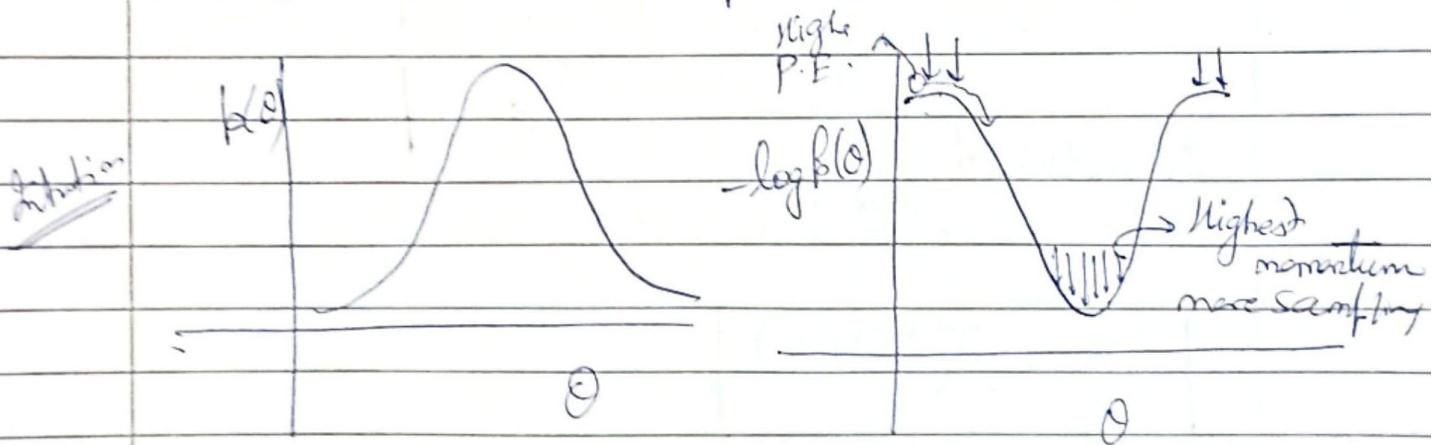
- If starting points are not overdispersed
- dist far from mean (inf variance)
- just by chance

SP62-R

We remove the warm-up then divide the remaining chain into two halves and then we test whether the both the half have similar dist" (mean & variance). This makes the assurance that the series has become stationary.

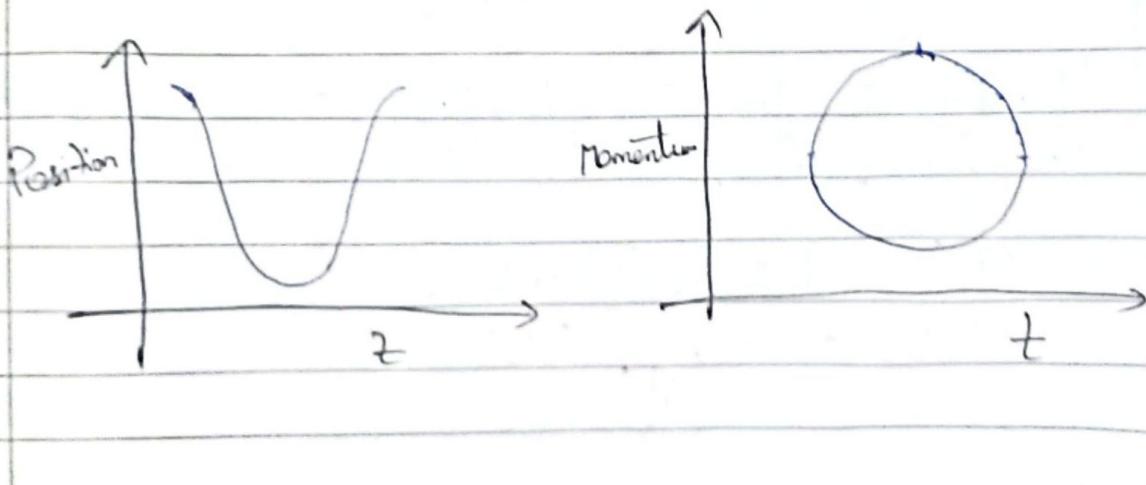
Hamiltonian Monte Carlo

Faster and efficient than metropolis algorithm.
Scale well with multiple dimension.



- Hamiltonian Equation

- operates on a d -dimensional position vector θ and a d -dimensional momentum vector ϕ
 $\theta \in \mathbb{R}^d$ and $\phi \in \mathbb{R}^d$



So the full state given as $(\theta, \phi) \in \mathbb{R}^{2d}$
 Described by a function $H(\theta, \phi)$

Equations of the motion

↳ How θ and ϕ changes over time

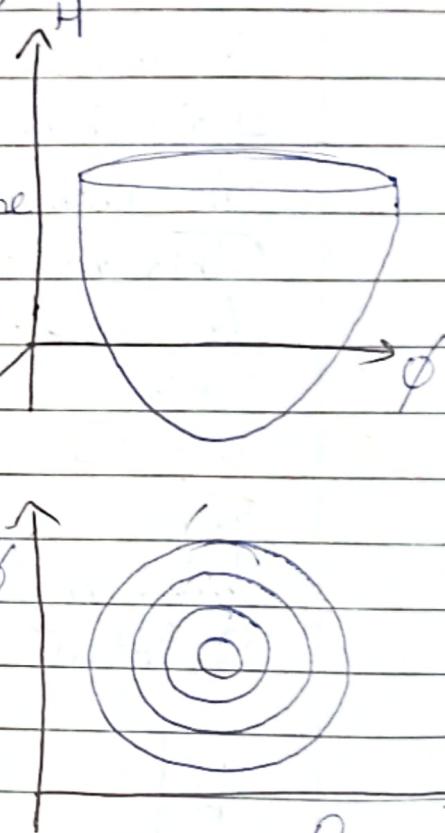
Given by hamiltonian equations

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \phi} \quad \frac{d\phi}{dt} = -\frac{\partial H}{\partial \theta}$$

for d -dimension

$$\frac{d\theta_i}{dt} = \frac{\partial H}{\partial \phi_i} \quad \frac{d\phi_i}{dt} = -\frac{\partial H}{\partial \theta_i}$$

for $i = 1, 2, \dots, d$



In matrix form-

Combine θ and ϕ in a vector $z = (\theta, \phi) \in \mathbb{R}^{2d}$

$$\frac{\partial z}{\partial t} = J \nabla H(z)$$

→ $\nabla H(z)$ is the gradient of H i.e.

$$[\nabla H(z)]_k = \frac{\partial H}{\partial z_k}$$

$$J = \begin{bmatrix} 0_{d \times d} & I_{d \times d} \\ I_{d \times d} & 0_{d \times d} \end{bmatrix} \in \mathbb{R}^{2d \times 2d}$$

Hamiltonian function gives the total energy which is the sum of the potential & kinetic energy.

$$H(\theta, \phi) = V(\theta) + K(\phi)$$

$$V(\theta) = -\log f(\theta)$$

↳ Potential Energy.

$$K(\phi) = \frac{1}{2} \phi^T M^{-1} \phi$$

$$\left[\text{Analogy} \quad \text{KE} = \frac{p^2}{2m} \right]$$

- M is the mass matrix

- Symmetric, Positive Definite

- Typically M is considered to be diagonal with a multiple of the identity matrix I

$$\left[\begin{matrix} c & 0 \\ 0 & c \end{matrix} \right] \text{ this type}$$

$$\rightarrow K(\phi) = -\log p(\phi)$$

where $p(\phi) = \mathcal{N}(0, M)$

Thus equations ① and ⑪

$$\begin{aligned} \frac{d\theta_i}{dt} &= \frac{\partial H}{\partial \phi_i} = \frac{\partial K(\phi)}{\partial \phi_i} = M^{-1} \phi \\ \frac{\partial K(\phi)}{\partial \phi_i} &= \frac{1}{2} (M^{-1} \phi) \frac{\partial \phi}{\partial \phi_i} + (M^{-1} \phi) \frac{\partial \phi}{\partial \phi_i} \\ &= \frac{2}{2} (M^{-1} \phi) = M^{-1} \phi \end{aligned}$$

$$\frac{\partial \phi_i}{\partial t} = \frac{\partial H}{\partial \theta_i} = -\frac{\partial U(\theta)}{\partial \theta_i} = -\frac{\partial (-\log p(\theta, y))}{\partial \theta_i}$$

Gradient of the posterior

If we will solve the above differential Eq's -

$$\theta_i(t+s) = f(\theta_i(t))$$

$$\phi_i(t+s) = g(\phi_i(t))$$

Now the sampling will become deterministic and that's the benefit of this method.

For a 1D case:

θ and ϕ are scalar

$$H(\theta, \phi) = U(\theta) + K(\phi)$$

$$\text{Let } U(\theta) = \theta^2/2$$

$$U(\theta) = -\log(f(\theta/y)) = \theta^2/2 \quad (\text{standard normal as posterior})$$

$$K(\phi) = \phi^2/2 \quad \left[\begin{array}{l} \text{This is standard normal} \\ \text{which we generally use} \end{array} \right]$$

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \phi} = \phi \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Equations of SHM.}$$

$$\frac{\partial \phi}{\partial t} = -\frac{\partial H}{\partial \theta} = -\theta \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$\theta(t) = A \cos(t) + B \sin(t)$$

$$\phi(t) = -A \sin(t) + B \cos(t)$$

Initial Conditions

$$\text{At } t=0, \theta(0) = \theta_0$$

$$\phi(0) = \phi_0$$

$$\theta_0 = A; \phi_0 = B$$

$$\theta(t) = \theta_0 \cos(t) + \phi_0 \sin(t)$$

$$\phi(t) = -\theta_0 \sin(t) + \phi_0 \cos(t)$$

$$\text{Let } \theta_0 = r_1 \cos a \quad \phi_0 = -r_1 \sin a$$

$$\theta(t) = r_1 \cos(a+t) \quad r_1^2 = \theta_0^2 + \phi_0^2$$

$$\phi(t) = -r_1 \sin(a+t) \quad \tan a = -\frac{\phi_0}{\theta_0}$$

Certain properties of Hamiltonian Monte Carlo dynamics

- a) **Reversible** : Knowledge of $\theta(t)$ gives a one-to-one mapping to $\theta(t+s)$. Hence path can be traced.

b) Conservation of Hamiltonian

$H(\theta, \phi)$ will be conserved for a ϕ and θ . It will follow a fixed contour. So it will efficiently cover the contour.

c) Volume preservation

If we take a set of θ and ϕ to constitute a region. If the θ and ϕ moves then the region will move with constant volume.

d) Symplecticity

Ensure that coverage based on dynamics is good enough to get enough values of θ & ϕ .

Discretizing Hamiltonian Equations

$$\text{Let } K(\phi) = \sum_{i=1}^d \frac{\phi_i^2}{2m_i}$$

- Simple Method
 - Euler method

$$\begin{aligned} \phi_i(t + \epsilon) &= \phi_i(t) + \epsilon \frac{d\phi_i}{dt} \\ &= \phi_i(t) - \epsilon \frac{\partial H}{\partial \theta_i} \Big|_{\theta_i(t)} \\ &= \phi_i(t) - \epsilon \frac{\partial U}{\partial \theta_i} \Big|_{\theta_i(t)} \end{aligned}$$

$$\begin{aligned}
 \theta_i(t+\epsilon) &= \theta_i(t) + \epsilon \frac{\partial \theta_i}{\partial t} \bigg|_{t=t} \\
 &= \phi_i \theta_i(t) + \epsilon \frac{\partial \phi_i}{\partial t} \bigg|_{t=t} \\
 &= \theta_i(t) + \epsilon \frac{\phi_i(t)}{m_i} \bigg|_{t=t}
 \end{aligned}$$

If we take small ϵ then it will converge to the original contour lines but if it is large then it will deviate from the original contours.

The spiraling out is the main problem in this approach.

- Modified Euler Method

$$\phi_i(t+\epsilon) = \phi_i(t) + \epsilon \frac{\partial \phi_i}{\partial t} \bigg|_{t=t} \bigg|_{\theta_i(t) = \theta_i(t)}$$

$$\begin{aligned}
 \theta_i(t+\epsilon) &= \theta_i(t) + \epsilon \frac{\partial \theta_i}{\partial t} \bigg|_{t=t} \bigg|_{\theta_i(t+\epsilon)} \xrightarrow{\text{Taking the gradient of the same step.}} \\
 &= \theta_i(t) + \epsilon \frac{\phi_i(t+\epsilon)}{m_i} \bigg|_{t=t}
 \end{aligned}$$

- Leapfrog Integration

$$\phi_i(t+\epsilon/2) = \phi_i(t) + \epsilon/2 \frac{\partial \phi_i}{\partial t} \bigg|_{t=t} \bigg|_{\theta_i(t) = \theta_i(t)}$$

$$\theta_i(t+\epsilon) = \theta_i(t) + \epsilon \frac{\phi_i(t+\epsilon/2)}{m_i} \bigg|_{t=t}$$

$$\phi_i(t+\epsilon) = \phi_i(t) + \epsilon/2 \frac{\partial \phi_i}{\partial t} \bigg|_{t=t} \bigg|_{\theta_i(t+\epsilon) = \theta_i(t+\epsilon)}$$

~~Problem~~

If there is a peak in a gradient then with larger stepsize

Problem

In case of high divergence leap frog integrator error increases. If small step size \rightarrow very huge time req. If very large step size then ~~jumps to~~ diverges more.

Sol'

Use an adaptive method

- Step size can be high for lower divergence to move faster
- Decrease step size adaptively for high divergence

This is followed in NUTS (No U-turn Sampling)
 \rightarrow Sampling after 1 steps.

After getting $\theta^* \& \theta^{**}$ after leap frog we can apply

metropolis hastings like method to accept it with

$$r = \frac{p(\theta^*|y)}{p(\theta^{**}|y)} \frac{p(\theta^{**})}{p(\theta^*)}$$

$$\theta^2 = \begin{cases} \theta^* & \text{with prob. } \min(r, 1) \\ \theta^{**} & \text{otherwise} \end{cases}$$

#

Hierarchical Model

$p(\theta_j | \tau)$ $\theta_1 \theta_2 \theta_3 \dots \theta_n$ parameters

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

$p(y_i | \theta_j)$ $y_{i1} y_{i2} y_{i3} \dots y_{in}$ observation

Joint posterior

$$f(\theta, \tau | y) \propto f(y | \theta, \tau) \cdot f(\theta, \tau)$$

$$f(y | \theta) \cdot f(\theta | \tau) \cdot f(\tau)$$

$$\{ \theta_j | \alpha, \beta \sim \text{Beta}(\theta_j | \alpha, \beta)$$

$$\text{multiply and divide} \quad y_j | n_j, \theta_j \sim \text{Bin}(y_j | n_j, \theta_j)$$

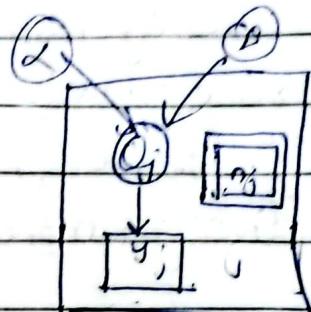


Plate diagram for j^{th} hospital.

Joint posterior

$$f(\theta_1, \theta_2, \dots, \theta_n | \alpha, \beta | y)$$

$$= \prod_{i=1}^n f(\theta_i | \alpha, \beta, y) f(\alpha, \beta | y)$$

$$f(y | \alpha, \beta) f(\alpha, \beta) \rightarrow \text{Apply sampling}$$

$$\rightarrow \alpha \sim \text{Beta}(\alpha + \beta) \sim \frac{1}{2}$$

Hierarchical model has less variance for smaller value of N .

If ϕ denotes the hyperprior parameters & θ denotes the parameters

- Joint posterior distribution $f(\theta, \phi | y)$

Joint prior distribution

$$f(\theta, \phi) = f(\phi) f(\theta | \phi)$$

$$\rightarrow f(\theta, \phi | y) \propto f(\theta, \phi) f(y | \theta, \phi)$$

$$\propto f(\theta, \phi) f(y | \theta)$$

Hierarchical model works in case of independent systems which are exchangeable

Data
Page

- Marginal posterior density $\int f(\theta, \phi | y) d\theta$

$$f(\theta, \phi | y) = f(\theta | \phi, y) f(\phi | y)$$

$$\text{or } f(\phi | y) = \frac{f(\phi, \theta | y)}{f(\theta | \phi, y)}$$

The rat tumor model

- Data from 71 experiment

- Each experiment follows an independent binomial distribution.

$$y_j \sim \text{Bin}(\theta_j; n_j)$$

No. of rats with (+)ve outputs n_j is the total number of rats.

θ_j is assumed to be independent samples from a prior beta distribution.

$$\theta_j \sim \text{Beta}(\alpha, \beta)$$

Joint posterior distribution of $f(\theta, \alpha, \beta | y)$ where

$$\theta = (\theta_1, \theta_2, \theta_3, \dots, \theta_j)$$

$$\propto f(\theta, \alpha, \beta) f(y | \theta, \alpha, \beta)$$

$$\propto \prod f(\theta_j)$$

$$\propto f(\alpha, \beta) \prod f(\theta_j | \alpha, \beta) f(y_j | \theta_j, \alpha, \beta)$$

$$f(\theta_1, \theta_2, \dots, \theta_j | \alpha, \beta) = \prod_{j=1}^J \frac{1}{\Gamma(\alpha + \beta)} \theta_j^{\alpha-1} (1-\theta_j)^{\beta-1}$$

$$f(y | \theta, \alpha, \beta) \propto \prod_{j=1}^J \theta_j^{y_j} (1-\theta_j)^{n_j - y_j}$$

$$\propto f(\alpha, \beta) \prod_{j=1}^J \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \Gamma(\beta)} \theta_j^{\alpha-1} (1-\theta_j)^{\beta-1} \prod_{j=1}^J \theta_j^{y_j} (1-\theta_j)^{n_j-y_j}$$

$$\propto f(\alpha, \beta) \prod_{j=1}^J \theta_j^{\alpha+y_j-1} (1-\theta_j)^{\beta+n_j-y_j-1}$$

Conditional Posterior Distribution $[(\theta|\alpha, \beta), y]$

$$f[(\theta|\alpha, \beta), y] = f(\theta|\alpha, \beta) f(y|\{\theta|\alpha, \beta\})$$

$$\propto \prod_{j=1}^J \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \Gamma(\beta)} \theta_j^{\alpha-1} (1-\theta_j)^{\beta-1} \prod_{j=1}^J \frac{\theta_j^{n_j} (1-\theta_j)^{y_j}}{\Gamma(y_j) \Gamma(n_j-y_j)} \theta_j^{n_j-y_j}$$

$$\prod_{j=1}^J \frac{\Gamma(\alpha+\beta+n_j)}{\Gamma(\alpha+y_j) \Gamma(\beta+n_j-y_j)} \cdot \theta_j^{\alpha+y_j-1} (1-\theta_j)^{\beta+n_j-y_j-1}$$

Marginal posterior distribution $f(\alpha, \beta|y)$

$$f(\alpha, \beta|y) = \frac{f(\alpha, \beta, \theta|y)}{f(\theta|\alpha, \beta, y)} \quad \text{(Dividing the eqn that we got)}$$

$$\propto f(\alpha, \beta) \prod_{j=1}^J \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \Gamma(\beta)} \frac{\prod_{j=1}^J \frac{\Gamma(\alpha+\beta+n_j)}{\Gamma(n_j) \Gamma(\beta+n_j-y_j)}}{\prod_{j=1}^J \frac{\Gamma(\alpha+y_j)}{\Gamma(n_j) \Gamma(\beta+n_j-y_j)}}$$

$$\left\{ \begin{array}{l} f(\theta|y) \propto f(y|\theta) \cdot f(\theta) \\ f(y|\{\theta|\alpha, \beta\}) \cdot f(\theta|\alpha, \beta) \end{array} \right.$$

Model Checking



Internal Validation External Validation

(On the Slides)

External Validation

Suppose observed data or

$$y = y_1, y_2, \dots, y_n$$

$$f(y|\theta) = \prod_{i=1}^n f(y_i|\theta)$$

Measures

Change the θ and
generate diff y_i

1) Mean Square Error: $\frac{1}{n} \sum (y_i - E(y_i|\theta))^2$

2) Weighted MSE: $\frac{1}{n} \sum_{i=1}^n \frac{(y_i - E(y_i|\theta))^2}{\text{Var}(y_i|\theta)}$

Advantages of these techniques:

- Easy to compute.

Disadvantage

Doesn't perform well for models other than normal.

2) Log predictive density or Log Likelihood or (LPD)

$\log(f(y|\theta)) \rightarrow$ For normal dist the lpd is same as MSE.

- When n is large then the lpd converges towards the a value of θ that follows the highest value of $\log(f(y|\theta))$ that follows the highest posterior density.

Information Criterion Based Measures

Akaike Information Criterion

- Find the expected log predictive density based on observed data \hat{elpd} .
- Find the measure of penalty to subtract.

k = no. of parameters of the estimated model.

$$\log f(y|\hat{\theta}_{MLE})$$

$$E(\hat{elpd}) = \int \log f(y|\hat{\theta}_{MLE}) \cdot f(y) dy$$

Original/True dist of y (data)
This is unknown.

Estimating \hat{elpd}

$$\log(f(y|\hat{\theta}_{MLE})) - k$$

$$AIC = -2 \log(f_y(\hat{\theta}_{MLE})) + 2k$$

Disadvantages

- Not much effective for non-linear model. Only work for linear model with flat prior.
- For hierarchical model, it counts the bias as k is too harsh penalty.

Deviance Information Criterion (DIC)

(2 changes than AIC)

- 1) Log predictive Density, uses $\hat{\theta}_{MAP}$ $\log p(\mathbf{y}|\hat{\theta}_{MAP})$
- 2) Rather than using k as penalty term it uses the variance term of $\log p(\mathbf{y}|\theta_j)$

$$DIC = -2 \log p(\mathbf{y}|\hat{\theta}_{MAP}) + 2 P_{DIC}$$

$$P_{DIC} = \sum_{j=1}^s \text{Var} \log p(\mathbf{y}|\theta_j)$$

More var high uncertainty so more penalty.

Watanabe-Akaike Information Criterion (WAIC)

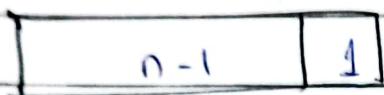
- Use the uncertainty of the measure of θ in the lpd term.
- Uses the variance of each data point.

$$WAIC = -2 \sum_{i=1}^n \log \frac{1}{s} \sum_{j=1}^s p(\mathbf{y}_i|\theta_j) + 2 P_{WAIC}$$

$$P_{WAIC} = \sum_{i=1}^n \sum_{j=1}^s \text{Var} \log (p(\mathbf{y}_i|\theta_j))$$

First we will calculate posterior (we can sample s data points and then calculate the mean of likelihood $\log p(\mathbf{y}_i|\theta_j)$, sum it over i)

Leave-one-out - Cross Validation (Loo-Cv)



Take $n-1$ data point for training and use the rest 1 data as out of sample testing

No need for penalty.

Leave 1.

$$\text{Loo-Cv} : \sum_{i=1}^n \log \left(\frac{1}{S} \sum_{j=1}^S p(y_i | y_{-i}; \theta_j) \right)$$

Which measure is better ??

