

BGN: 2191A

P6

Q.7

a.i. Let L : likelihood $= P(x_1, x_2, \dots, x_n | \alpha)$

$$\log L = \sum_{i=1}^n \log P(x_i | \alpha)$$

$$= \sum_{i=1}^n \log \left(\frac{\alpha}{x_i^{\alpha+1}} \right)$$

$$= \sum_{i=1}^n \log \alpha - (\alpha+1) \log x_i$$

Let $\hat{\alpha}$ maximise $\log L$

$$\therefore \left(\frac{\partial}{\partial \alpha} \log L \right) \Big|_{\alpha=\hat{\alpha}} = \sum_{i=1}^n \frac{1}{\hat{\alpha}} - \log x_i$$

$$= \frac{n}{\hat{\alpha}} - \sum_{i=1}^n \log x_i$$

$$= 0$$

$$\therefore \hat{\alpha} = \frac{n}{\sum_{i=1}^n \log x_i}$$

ii. $P(\text{Pareto}(1, \alpha) \leq q) = 0.99$

$$\therefore q \geq 1 \text{ and } 1 - q^{-\alpha} = 0.99$$

$$\therefore q^{-\alpha} = 0.01 = \frac{1}{100}$$

$$\therefore q^{\alpha} = 100$$

$$\therefore q = 100^{1/\alpha}$$

iii. Calculate $\hat{a} = \frac{n}{\sum_{i=1}^n \log x_i}$ as in part i.

Generate n datapoints independently sampled from Pareto(1, \hat{a}), $\hat{x}_i^{(1)}$. Repeat this a large number, N , times yielding $\hat{x}_i^{(j)}$ values.

For each iteration, calculate a new $\hat{a}^{(j)}$ and $q^{(j)} = 100^{-\hat{a}^{(j)}}$

Sort the values of $q^{(j)}$ into q_{\uparrow}

$[a = q_{\uparrow}^{(\lfloor \frac{s}{2} N \rfloor)}, b = q_{\uparrow}^{(\lfloor (1 - \frac{s}{2}) N \rfloor)}]$ is the ^{two-tailed} confidence interval at significance level s (confidence $1-s$)

Pseudo code:

$$\hat{a} = \frac{n}{\text{sum}(\log(x))}$$

$$\hat{x} = \text{pareto_sample}(1, \hat{a}, \text{shape} = (n, N))$$

$$\hat{a} = \frac{n}{\text{sum}(\log(\hat{x}), \text{axis} = 0)}$$

$$q = 100^{-\hat{a}}$$

$$q = \text{sort}(q)$$

$$a, b = q[\lfloor \frac{s}{2} N \rfloor], q[\lfloor (1 - \frac{s}{2}) N \rfloor]$$

e.g. $N=10000$, $s=5\%$.

b. Generate n data points sampled according to the empirical distribution of $\{x_1, \dots, x_n\}$.
Repeat N times to get N datasets, where N is large

sort each dataset and access position $[0.99n]$ for each of them, q .

If α is the significance level ($1 - \text{confidence}$)
then to get a two-tailed confidence interval, calculate

$$[a = q[\lfloor \frac{\alpha}{2} N \rfloor], b = q[\lfloor (1 - \frac{\alpha}{2}) N \rfloor]]$$

Pseudocode:

$r = \text{random_choice}(x, \text{shape} = (N, n))$

$r = \text{sort}(r, \text{axis} = 0)$

$q = r[\lfloor 0.99n \rfloor, :]$

$q = \text{sort}(q)$

$a, b = q[\lfloor \frac{\alpha}{2} N \rfloor], q[\lfloor (1 - \frac{\alpha}{2}) N \rfloor]$

This would be unreliable when n is small, and so the empirical distribution of x would be unlikely to closely resemble the "true" underlying distribution of each x_i .

c. For some confidence level c , we wish to find a, b such that

$$P(a \leq Q \leq b) = c$$

where Q is the random variable representing the 99th percentile

We must assign a prior (arbitrarily) to X , say $\text{Normal}(0, 1)$

We then calculate the likelihood of the data given a parameter (e.g. μ) and multiply by the prior (and normalise) to give the posterior for X .

Using this, calculate the posterior for Q .

Pseudocode:

$\mu = \text{inspact} - 100, 100, \delta\mu$

prior = $\text{Normal}(0, 1)$

likelihood = $[\text{Normal}(\mu, 1).pdf(x_i) \text{ for } x_i \text{ in } x]$

likelihood = $\text{product}(\text{likelihood}, \text{axis}=0)$

posterior = $\text{likelihood} \times \text{product}(\text{prior}.pdf(x_i) \text{ for } x_i \text{ in } x)$

posterior = $\text{posterior} \div \text{sum}(\text{posterior} \cdot \delta\mu)$

// note: posterior is a function of μ_i ($P(\mu_i | \text{data})$)

percentile = $\text{Normal}(\mu_i, 1).cdf(0.99)$ // range of percentiles

perm = $\text{argsort}(\text{posterior})$

percentile, posterior = $\text{percentile}[\text{perm}], \text{posterior}[\text{perm}]$

$a = \text{cumsum}(\text{posterior})$

$s = 1 - c$ // significance

$i, j = \text{first_above}(\text{posterior}, s/2), \text{last_below}(\text{posterior}, 1 - s/2)$

$a, b = \text{percentile}[i], \text{percentile}[j]$ // $[a, b]$ is the confidence interval.