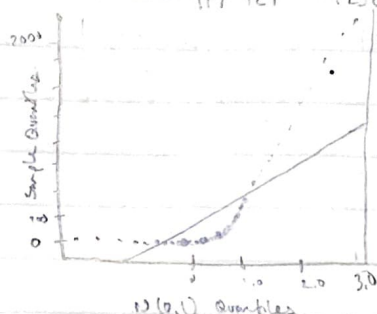


See Handouts 8.9, 10, Chp 1, 6, 14, sections 4.6, 5.3-5.6, 7.2-7.4 in Devore book.

(H.O.4) 1) 250 iid observations y_1, y_2, \dots, y_{250} from a process yield the following normal probability plot.



Using the plot given, describe the process distribution relative to a normal plot w/ respect to tail weight and symmetry.

From the Q-Q plot shown we can tell that the process distribution has a heavier right tail and lighter left tail than the $N(0,1)$ distribution. Furthermore, the process distribution is skewed heavily to the right.

Need to redo because $\hat{E}_0 < 1$. Need to combine last two groups.

(H.O.9) 2) Nylon Bars were tested for brittleness. Each of the 500 bars was modeled under similar conditions and was tested by placing a specified stress at 5 locations on the bar. Assuming each bar has uniform composition, the number of breaks on a given bar should be binomially distributed w/ an unknown probability p of breaking. The following table summarizes the outcome of the experiment:

Breaks/Bar	0	1	2	3	4	5	Total
Frequency	140	197	115	41	5	2	500

Use a GOF test to evaluate whether the data appear to be from a binomial model.

① Find MLE of θ ; $\hat{\theta}$ = Sample probability of a location on a bar breaking.

$$\hat{\theta} = \frac{0(140) + 1(197) + 2(115) + 3(41) + 4(5) + 5(2)}{500(5)} = 0.232 = \hat{\theta}$$

② Let P_{i+1} be the probability that a randomly selected bar has i breaks for $i \in [0, \dots, 5]$. That is $P_{i+1} = P[X=i]$ where $X \sim \text{Binomial}(5, \theta)$ distribution

$$\hat{P}_{i+1} = P(X=i) = f_i(i, \hat{\theta}) = \binom{5}{i} (0.232)^i (0.768)^{5-i} \text{ for } i \in [0, \dots, 5]$$

$$E_i = 500 \hat{P}_i$$

$$\hat{Q} = \frac{\sum_{i=1}^5 (O_i - E_i)^2}{E_i}$$

③ Follow table was calculated using R.

i	1	2	3	4	5	6	Total
\hat{P}_i	0.267813	0.43521	0.243814	0.073682	0.01112457	0.0006721092	1.00
E_i	133.94066	217.756	121.90727	36.82615	5.562184	0.33605466	500
O_i	140	197	115	41	5	2	500
$\hat{Q}_i = \frac{(O_i - E_i)^2}{E_i}$	0.3075	0.1131	0.3913	0.4730	0.05684	8.23888	9.58077

NOTE: $\hat{E}_0 < 1 \Rightarrow$ need to recalculate w/ 5 and 6 combined.

2) (contd)

$$Q_{new} = Q_{old} - E_{old} - E_{old} + E_{y, new} = 9.380771 - 5.542284 - 0.126054 + \frac{(7 - 5.898229)^2}{5.898229}$$

$$Q_{new} = 9.490813 \quad df_{new} = 5 - 1 - 1 = 3$$

$$P(Q \geq 9.490813) = 1 - P(Q \leq 9.490813) = 1 - pchisq(9.490813, 3)$$

$$P(Q \geq 9.490813) = 0.6043223 \Rightarrow \text{w/ a p-value that large, we would conclude a normal distribution is a good fit of random.}$$

(4.0.8) 3) A random sample of 500 data values are selected from 4 separate processes having cdfs F_1, F_2, F_3, F_4 . The plot of the sample quantiles vs a standard normal quantile for each of the four samples is given below. For each of these plots, select one of the following distributions to describe the pdf which generated the data. (Hint: make sure to take into consideration the size of values associated w/ each distribution).

Plot 1: N

Plot 2: H

Plot 3: D (w/ x_2 goes from 0 to 1)

Plot 4: Gr

- 4.) An experiment was conducted to investigate if the impact of the carcinogen DMBA could be delayed by treatment w/ a potential beta-blocker. 50 mature rats of the same general health were given the beta-blocker then injected w/ DMBA. The time in days, after exposure, at which the carcinoma was diagnosed for the rats are given below. From past studies, 150 days after exposure, a carcinoma was detected in untreated rats.

Data given in table

- Does a Weibull distribution appear to provide an adequate fit to the data? Justify your answer using both a GOF and a graphical plot.

Use R to find MLE of γ , α , β , $\hat{\alpha}$.

$$\hat{\gamma} = 1.3861041, \hat{\alpha} = 201.1786473$$

$$AD = 0.24708605$$

$$ADM = 0.2232267$$

- From Table 5 in H.O. 9 (pg 34) we see that our p-value is > 0.25 . Thus we have a very good fit.
- Similarly, looking at our QQ plot in R we see that the Weibull model is a good fit.

*note corr test

and AD are not matching up.

AD \Rightarrow good fit, corr test \Rightarrow not good.

*Ask about corr test for non-normal data.

$$X \sim N$$

might want to, also I using the Shapiro-Wilk's test instead of AD

A major problem in the Gulf of Mexico is the excessive capture of game fish by shrimpers.

A random sample of the catch of 50 shrimpers yield the following data concerning the catch per unit effort (CPUE) of red snappers, a highly sought game fish. Let C_i be the CPUE of the i^{th} shrimper. The data C_1, \dots, C_{50} is given in the below table.

(1.) CPUE data is often modeled using a lognormal distribution. Does the above data appear to be from a log-Normal distribution? Explain your answer w/ both a normal reference distribution plot and a Q-Q test.

Let $X_i = \ln(C_i)$.

$$\text{Then } X_i \sim N(\mu, \sigma^2)$$

Make I
d this
instead of
AD.

Using MLE estimates from R, we get $\hat{\mu} = 2.590627$, $\hat{\sigma}^2 = 2.070229$

$$AD = 0.4774204, \text{ADP} = 0.4550114.$$

From table 5 (4.0.9 pg 34) we see that our p-value is between (0.10, 0.25) thus the lognormal distribution provides a moderately good fit of the data.

In R: `data = log(data)`, `x = sort(data)`, `shapiro.test(x)`

$W = 0.97731$, $p\text{-value} = 0.4451 \Rightarrow$ the lognormal distribution provides an "excellent fit" (4.0.9 pg 5) to our data.

(2.) Use the Box-cox transformation of the CPUE data to determine the

most appropriate power transformation to transform the CPUE distribution to Normality.

How does the fit from the box-cox transformation compare to the fit for the log transformation?

From R: $\hat{\theta} = 0.074$;

In R: `data = data^0.074`; `shapiro.test(data)`

$$\Rightarrow W = 0.98151, p\text{-value} = 0.6173.$$

The fit from the box-cox transformation is better than the fit for the log transformation.

5) (Contd.)

(3) Use the R program from H.O. 10 to draw 10,000 Bootstrap samples from the CPUE data. From the 10,000 samples, estimate the standard error of the sample mean for the $Y_i = \log(C_i)$ data. Compare this estimate to the usual estimate S_y / \sqrt{n} , where S_y is the sample standard deviation computed from the $n=50$ values of $Y_i = \log(C_i)$.

$$\frac{S_y}{\sqrt{n}} = 0.2055464$$

$$SD(\bar{y}) = 0.2035765$$

The estimated standard error of the sample mean is approximately equal to the sample standard error.

(4) Use your bootstrap samples to estimate the mean and standard deviation of the following sample statistics for $Y_i = \log(C_i)$

(a) Sample median; $\hat{Q}(0.5)$

$$\bar{\hat{Q}}(0.5) = 2.814969, \quad SD(\bar{\hat{Q}}(0.5)) = 0.3721707$$

(b) Sample SD; S_y

$$\bar{S}_y = 1.433990, \quad SD(\bar{S}_y) = 0.1187079$$

(c) The sample MAD; \hat{MAD}_y

$$\bar{MAD} = 1.569568, \quad SD(\bar{MAD}) = 0.2244567$$

(H.O. 10 pg. 18)

(5) Historically, the $\log(\text{CPUE})$ data was modeled as a random sample from a $N(3, (1.5)^2)$ distribution. Compare your bootstrap estimates of the mean and standard deviation of $\hat{Q}(0.5)$ and S_y from part 4 of this problem to the theoretical mean and standard deviation of $\hat{Q}_y(0.5)$ and S_y based on $Y_i = \log(C_i)$ having a $N(3, (1.5)^2)$ distribution.

• Asymptotic Mean of $\hat{Q}(0.5)$ is $\mu = Q(0.5) = 3$

• This is approximately the same as the mean of the sample median

• Asymptotic SD of $\hat{Q}(0.5)$: $\sigma_A = \frac{\sqrt{0.527}}{\sqrt{50} f(Q(0.5))} = \frac{\sqrt{0.527}}{\sqrt{50} (0.1/\sqrt{1.5})} = \frac{1.5\sqrt{\pi/2}}{\sqrt{50}} = 0.2658681$

• The asymptotic SD of $\hat{Q}(0.5)$ is less than our estimated SD of $\hat{Q}(0.5)$.

• Asymptotic Mean of S_y is $\sigma = 1.5$

• The asymptotic mean of S_y is approximately equal to our estimate \bar{S}_y .

• Asymptotic SD of S_y is $\sigma_A = \frac{\sqrt{15.1875 - (1.5)^4}}{20\sqrt{n}}$, $15/\sigma_y = 3\sigma^4 = 3(1.5)^4 = 15.1875$

$$\sigma_A = \sqrt{15.1875 - (1.5)^4} / (2(1.5)\sqrt{50}) = 0.15$$

• The Asymptotic SD of S_y is slightly larger than our estimate

$$SD(\sigma_y)$$

- (6) A company has designed a new battery system for electric powered automobiles. To estimate the lifetime of the system, the design engineers place the batteries in 25 electric powered cars and test them under simulated city driving. Let Y_i be the time to failure of the batteries of the i th car, $i \in [1, \dots, 25]$. The failure times are recorded in units of 20,000 miles. The company wants to know the probability that the sample mean based on 25 observations will estimate the true mean w/in a margin of error of ± 0.2 (40,000 miles), provided that the true mean has a value of 5 (100,000 miles), that is, approximately, $P[-0.2 \leq \bar{Y} - 5 \leq 0.2]$
 $= P[4.8 \leq \bar{Y} \leq 5.2]$

- (1) From past studies, the distribution of the time to failure of the battery system is exponential w/ a mean value of 5 (100,000 miles). Let \bar{Y} be the sample mean time to failure of the 25 cars.

- (40.10.12) (a) What is the exact distribution of \bar{Y} if the exponential model is still valid?

→ *I'm thinking we use the fact that if X_1, \dots, X_n i.i.d. $\text{Exp}(\beta)$, then the distribution
 $n\hat{\beta} = \sum_{i=1}^n X_i \sim \text{Gamma}(n, \hat{\beta})$

Let $X = \sum_{i=1}^{25} Y_i \Rightarrow X \sim \text{Gamma}(25, 5)$

Then $\bar{Y} = (\frac{1}{25})X \Rightarrow X = h^{-1}(\bar{y}) = 25\bar{y} \quad \frac{d}{d\bar{y}} [h^{-1}(\bar{y})] = 25$

$f_{\bar{Y}}(\bar{y}) = \frac{1}{\Gamma(25) 5^{25}} (25\bar{y})^{25-1} e^{-(25\bar{y}/5)} \quad (25)$

$= \frac{5^{50}}{\Gamma(25) 5^{25}} \bar{y}^{25-1} e^{-(5\bar{y})} = \frac{5^{25}}{\Gamma(25)} \bar{y}^{25-1} e^{-5\bar{y}}$

$= \frac{1}{\Gamma(25) (\frac{1}{5})^{25}} \bar{y}^{25-1} e^{-5\bar{y}} \Rightarrow \boxed{\bar{Y} \sim \text{Gamma}(25, 1/5)}$

- (b) What is the mean and SD of \bar{Y} if the exponential model is still valid?

$E[\bar{Y}] = 25(1/5) = 5 = E[\bar{Y}]$

$SD(\bar{Y}) = \sqrt{25(1/5)^2} = 1 = SD(\bar{Y})$

- (2) Simulate 10000 random samples of size 25 from the exponential distribution w/ $\beta = 5$. Compute the sample mean for each of the 10000 samples. Display a normal reference plot for the 10000 sample means. Does the plot suggest that the sampling distribution of \bar{Y} is approximately normal?

• NO, the sampling distribution of \bar{Y} is not approximately normal

(Contd.)

(3) Compute or estimate $P[0.2 \leq \bar{Y} - S \leq 0.2]$ in each of the following ways

(a) Using the exact distribution of \bar{Y}

$$\begin{aligned} \bullet P[0.2 \leq \bar{Y} - S \leq 0.2] &= P[4.8 \leq \bar{Y} \leq 5.2] \\ &= P[\bar{Y} \leq 5.2] - P[\bar{Y} \leq 4.8] \end{aligned}$$

$$P[0.2 \leq \bar{Y} - S \leq 0.2] = 0.1580747$$

(b) Using Central Limit Theorem

$$0.1585194$$

(c) Using your simulated 10000 \bar{Y} 's

$$\hat{P}[0.2 \leq \bar{Y} - S \leq 0.2] = 0.1546$$

(4) What is the level of agreement in the three computations/estimators of $P[0.2 \leq \bar{Y} - S \leq 0.2]$?

- All three estimators agree to the nearest thousandth, the central limit theorem and the exact distribution agree to the nearest thousandth.