

Student name: _____

CE93 -- Engineering Data Analysis
Final Examination
Tuesday, May 17, 2005

Work on all four problems. Write clearly and state any assumptions you make. The exam is closed books and notes, except for three sheets of paper and the table of Matlab® distributions.

The problems have the following weights:

Problem 1 (30 points) _____
Problem 2 (20 points) _____
Problem 3 (20 points) _____
Problem 4 (30 points) _____

Exam grade (100 points) _____

Problem 1. (10+10+10 = 30 points)

Timely completion of a construction project depends on the weather and on the working conditions of machinery used on the job. The probabilities of delay for different combinations of weather condition (rain or no rain) and machinery breakdown are listed in the following table:

Weather	Machinery	Probability of delay in completion of the project
No rain	No breakdown	0.05
Rain	No breakdown	0.20
No rain	Breakdown	0.30
Rain	Breakdown	0.60

Weather forecast indicates 30% chance of rain. From previous experience, we know that the probability that the machinery will breakdown depends on the weather conditions (due to the added force required to move wet soil). If it rains the probability of breakdown is 0.25, whereas if it does not rain the probability of breakdown is only 0.10.

- Determine the probability that the construction job will be completed on time.
- If the job is known to have been completed on time, what is the probability that it rained?
- If the job is not completed on time, what is the probability that machinery breakdown occurred?

Solution

Define R = rain, B = Breakdown, D = Delay.

$$a) \quad P(D) = P(D | RB)P(RB) + P(D | \bar{R}B)P(\bar{R}B) + P(D | R\bar{B})P(R\bar{B}) + P(D | \bar{R}\bar{B})P(\bar{R}\bar{B})$$

$$P(RB) = P(B | R)P(R) = 0.25 \times 0.30 = 0.075$$

$$P(\bar{R}B) = P(B | \bar{R})P(\bar{R}) = 0.10 \times (1 - 0.30) = 0.070$$

$$P(R\bar{B}) = P(\bar{B} | R)P(R) = 0.75 \times 0.30 = 0.225$$

$$P(\bar{R}\bar{B}) = P(\bar{B} | \bar{R})P(\bar{R}) = (1 - 0.10) \times (1 - 0.30) = 0.630$$

$$P(D) = 0.60 \times 0.075 + 0.30 \times 0.070 + 0.20 \times 0.225 + 0.05 \times 0.630 = 0.143$$

$$1 - P(D) = 1 - 0.143 = 0.857 \quad \text{Ans.}$$

$$b) \quad P(R | \bar{D}) = \frac{P(\bar{D} | R)P(R)}{P(\bar{D})}$$

Student name: _____

$$\begin{aligned}P(D | R) &= P(D | RB)P(B | R) + P(D | R\bar{B})P(\bar{B} | R) \\ &= 0.60 \times 0.25 + 0.20 \times (1 - 0.25) \\ &= 0.30\end{aligned}$$

$$P(R | \bar{D}) = \frac{1 - 0.30}{0.857} \times 0.30 = 0.245 \quad \text{Ans.}$$

c)
$$P(B | D) = \frac{P(DB)}{P(D)}$$

$$\begin{aligned}P(DB) &= P(DB | R)P(R) + P(DB | \bar{R})P(\bar{R}) \\ &= P(D | BR)P(B | R)P(R) + P(D | B\bar{R})P(B | \bar{R})P(\bar{R}) \\ &= 0.6 \times 0.25 \times 0.30 + 0.3 \times 0.10 \times (1 - 0.30) \\ &= 0.0660\end{aligned}$$

$$P(B | D) = \frac{0.0660}{0.143} = 0.462 \quad \text{Ans.}$$

Problem 2. (10+10 = 20 points)

The following two equations describe the monthly water available to cities A and B. City A gets 40 million gallons through a contract with the water distribution system of a neighboring city, 40% of the water available from reservoir 1 and 70% of the water available from reservoir 2. City B gets 20 million gallons on contract basis, plus the remainder of the available water from reservoirs 1 and 2.

$$W_A = 40 + 0.4X_1 + 0.7X_2$$

$$W_B = 20 + 0.6X_1 + 0.3X_2$$

Let the available water from reservoir 1, X_1 , be a random variable with mean 120 millions of gallon and a c.o.v. of 30%. Also, let the available water from reservoir 2, X_2 , be a random variable with mean 80 million gallons and a standard deviation of 25 million gallons. Further, assume that the two variables have a correlation coefficient of 0.30.

- Determine the means, the standard deviations and the correlation coefficient of the amounts of water available to cities A and B in one month.
- Suppose the variables X_1 and X_2 are jointly normal. What is the probability that the amount of water available to city B in one month will be less than 80 million gallons?

Solution

$$\text{a) } \mu_{W_A} = 40 + 0.4 \times 120 + 0.7 \times 80 = 144 \text{ million gallons} \quad \text{Ans.}$$

$$\mu_{W_B} = 20 + 0.6 \times 120 + 0.3 \times 80 = 116 \text{ million gallons} \quad \text{Ans.}$$

$$\begin{aligned} \sigma_{W_A}^2 &= 0.4^2 \times (120 \times 0.3)^2 + 0.7^2 \times 25^2 + 2 \times 0.4 \times 0.7 \times 0.30 \times (120 \times 0.3) \times 25 \\ &= 665 \end{aligned}$$

$$\sigma_{W_A} = 25.8 \text{ million gallons} \quad \text{Ans.}$$

$$\begin{aligned} \sigma_{W_B}^2 &= 0.6^2 \times (120 \times 0.3)^2 + 0.3^2 \times 25^2 + 2 \times 0.6 \times 0.3 \times 0.30 \times (120 \times 0.3) \times 25 \\ &= 620 \end{aligned}$$

$$\sigma_{W_B} = 24.9 \text{ million gallons} \quad \text{Ans.}$$

$$\begin{aligned} \text{Cov}[W_A, W_B] &= 0.4 \times 0.6 \times (120 \times 0.3)^2 + 0.7 \times 0.3 \times 25^2 \\ &\quad + (0.4 \times 0.3 + 0.7 \times 0.6) \times 0.30 \times (120 \times 0.3) \times 25 \\ &= 588 \end{aligned}$$

$$\rho_{w_a w_b} = \frac{588}{25.8 \times 24.9} = 0.915$$

- Since W_B is a linear function of jointly normal random variables, it is also normal.

$$P(W_B < 80) = \Phi\left(\frac{80 - 116}{24.9}\right) = \Phi(-1.446) = 0.074 \quad \text{Ans.}$$

Problem 3. (5+5+5+5 = 20 points)

The annual maximum water level in a river fits a lognormal distribution with mean 4m and a standard deviation of 1m.

- What is the probability that the water level in the river will exceed 6m in a given year?
- Suppose the annual maximum water levels in the river can be considered to be statistically independent. What is the probability that the water level in the river will exceed 6m during the next 5 years?
- What is the probability that the first occurrence of a water level above 6m will occur on the fifth year?
- What is the mean number of years that we must wait until the water level exceeds 6m?

Solution

- $$\zeta = \sqrt{\ln(1 + 0.25^2)} = 0.246$$

$$\lambda = \ln(4) - 0.5 \times 0.246^2 = 1.356$$

$$P(\text{Level} > 6) = 1 - \Phi\left[\frac{\ln(6) - 1.356}{0.246}\right] = 1 - \Phi(1.771) = 1 - 0.962 = 0.0383 \quad \text{Ans.}$$
- $$P(\text{water level exceed 6m in 5 tears}) = 1 - (1 - 0.0383)^5 = 0.177 \quad \text{Ans.}$$
- $$P(\text{first occurrence of water level above 6m in 5}^{\text{th}} \text{ year}) = (1 - 0.0383)^4 \times 0.0383$$

$$= 0.0328 \quad \text{Ans.}$$
- $$1/0.0383 = 26.1 \text{ years (average return period)} \quad \text{Ans.}$$

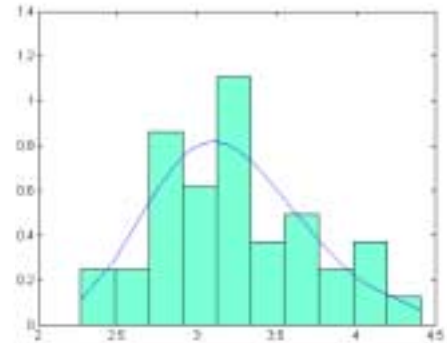
Problem 4. (10+10+10 = 30 points)

In the following Matlab® sessions, x is a 38x3 matrix. The first column contains the biochemical oxygen demand (BOD) and the second and third columns contain the corresponding values of nitrates (NO₃-N) and ammonia (NH₃-N) in a river in units of milligrams per liter. Carefully review the Matlab® commands and calculation results in each session and answer the questions listed below the session. Note that the results from each session are carried over to the next session.

Session 1

```
>> x = load('E6_2.dat');
>> size(x)
ans =
    38     3
mean(x)
ans =
    3.2176    8.5332    0.7811
>> std(x)
ans =
    0.5025    1.6012    0.2157
>> corrcoef(x)
ans =
    1.0000    0.6495    0.5153
    0.6495    1.0000    0.4154
    0.5153    0.4154    1.0000
>> [par, pari] = lognfit(x(:,1),0.05)
par =
    1.1569    0.1554
pari =
    1.1058    0.1267
    1.2080    0.2011
>> normfreq(x(:,1))
>> hold on
>> plot(sort(x(:,1)),lognpdf(sort(x(:,1)),par(1),par(2)))
>> hold off
```

Figure 1



Questions for Session 1

- What is the sample mean, sample standard deviation and sample c.o.v. of BOD data?
 - Sample mean = 3.2176 milligrams per liter
 - Sample stdev = 0.5025 milligrams per liter
 - Sample c.o.v. = $0.5025/3.2175 = 0.156$ or 15.6%
- What is the sample correlation coefficient between BOD and NO₃-N data?
 - $\rho_{BOD,NO_3-N} = 0.6495$
- What are the unbiased point estimates and 95% confidence intervals of the parameters λ and ζ of a lognormal distribution that is fitted to the BOD data?
 - $\hat{\lambda} = 1.1569, \hat{\zeta} = 0.1554$ unbiased estimates
 - $\langle \lambda \rangle_{0.95} = (1.1058, 1.2080), \langle \zeta \rangle_{0.95} = (0.1267, 0.2011)$ 95% confidence intervals
- Explain what is shown in Figure 1.

Shown in Figure 1 are the normalized frequency diagram of the BOD data together with the PDF of the fitted lognormal distribution with the unbiased parameter estimates $\hat{\lambda}$ and $\hat{\zeta}$.

Session 2

```
>> sx = sort(x(:,1));
>> [h,p,kstat,cv] = kstest(sx,[sx,normcdf(sx,mean(sx),std(sx))])
h =
    0
p =
    0.6822
kstat =
    0.1138
cv =
    0.2155
>> [h,p,kstat,cv] = kstest(sx,[sx,logncdf(sx,par(1),par(2))])
h =
    0
p =
    0.9436
kstat =
    0.0837
cv =
    0.2155
```

Questions for Session 2

- Would you reject the hypothesis of a normal distribution fit to the BOD data? Explain your reasoning.

I would not reject the hypothesis of a normal distribution fit to the BOD data because $h=0$ after the first kstest command. Furthermore, the K-S statistic $D_n = 0.1138$ is smaller than the critical value $D_n^\alpha = 0.2155$ at the 5% significance level.

- Would you reject the hypothesis of a lognormal distribution fit to the BOD data? Explain your reasoning.

I would not reject the hypothesis of a lognormal distribution fit to the BOD data because $h=0$ after the second kstest command. Furthermore, the K-S statistic $D_n = 0.0837$ is smaller than the critical value $D_n^\alpha = 0.2155$ at the 5% significance level.

- Which of the two distributions (normal or lognormal) better fits the data? Explain your reasoning.

The lognormal distribution fits better because the maximum deviation between the empirical and theoretical CDF's, the K-S statistic, is smaller for it ($D_n = 0.0837$) than for the normal distribution.

Session 3

```
>> [b,bi,r,ri,stats] = regress(x(:,1),[ones(38,1),x(:,2)]);
>> b
    1.4784
```

```

0.2038
>> stats
0.4219 26.2738 0.0000 0.1500
>> [b,bi,r,ri,stats] = regress(log(x(:,1)),[ones(38,1),x(:,2)]);
>> b
0.5964
0.0657
>> stats
0.4579 30.4080 0.0000 0.0135
>> [b,bi,r,ri,stats] = regress(log(x(:,1)),[ones(38,1),x(:,2),x(:,3)]);
>> b
0.5224
0.0520
0.2438
>> stats
0.5527 21.6230 0.0000 0.0114

```

Questions for Session 3

- We have regressed BOD against NO₃-N and *ln*(BOD) against NO₃-N. Which model in your opinion is better? Explain your reasoning.

The regression of *ln*(BOD) against NO₃-N is a better model than the regression of BOD against NO₃-N. This is because the R² value of the *ln*(BOD) regression (=0.4579) is greater than the R² value of the BOD regression (=0.4219).

- What fraction of the variability in *ln*(BOD) is being “explained” by NO₃-N?
45.79 percent of the variance in *ln*(BOD) is being “explained” by NO₃-N.
- Sketch the mean and mean ± one standard deviation plots of the two regression models on the scatter diagram shown on the next page.

For regression of BOD against NO₃-N:

$$E[\text{BOD} \mid \text{NO}_3\text{-N}] = 1.4784 + 0.2038(\text{NO}_3\text{-N})$$

$$\sigma_{\text{BOD} \mid \text{NO}_3\text{-N}} = \sqrt{0.1500} = 0.387$$

For regression of *ln*(BOD) against NO₃-N:

$$E[\ln(\text{BOD}) \mid \text{NO}_3\text{-N}] = 0.5964 + 0.0657(\text{NO}_3\text{-N})$$

$$\sigma_{\ln(\text{BOD}) \mid \text{NO}_3\text{-N}} = \sqrt{0.0135} = 0.116$$

The mean and ± one standard deviation curves are sketched in the figure. Note that the curves for *ln*(BOD) are obtained by taking exponential of the computed values and, therefore, are nonlinear.

- In regressing *ln*(BOD) against both NO₃-N and NH₃-N, how much of the variability in *ln*(BOD) is being explained by the two variables? Is this a better prediction model than the previous two? Explain your reasoning.

Student name: _____

When regressing $\ln(\text{BOD})$ against both $\text{NO}_3\text{-N}$ and $\text{NH}_3\text{-N}$, 55.27% of the variability in $\ln(\text{BOD})$ is being explained, which is greater than the 45.79% mentioned above. Thus, the regression on both $\text{NO}_3\text{-N}$ and $\text{NH}_3\text{-N}$ is a better model.

- For $\text{NO}_3\text{-N} = 10$ milligram per liter and $\text{NH}_3\text{-N} = 0.8$ milligram per liter, what is the probability that BOD will exceed 4 milligrams per liter?

$$\begin{aligned} E[\ln(\text{BOD}) \mid \text{NO}_3\text{-N}, \text{NH}_3\text{-N}] &= 0.5224 + 0.0520(\text{NO}_3\text{-N}) + 0.2438(\text{NH}_3\text{-N}) \\ &= 0.5224 + 0.0520 \times 10 + 0.2438 \times 0.8 = 1.237 \end{aligned}$$

$$\sigma_{\ln(\text{BOD}) \mid \text{NO}_3\text{-N}, \text{NH}_3\text{-N}} = \sqrt{0.0114} = 0.107.$$

$$P[\text{BOD} > 4] = P[\ln(\text{BOD}) > \ln(4)]$$

$$= 1 - \Phi\left(\frac{\ln(4) - 1.237}{0.107}\right) = 1 - \Phi(1.391) = 1 - 0.918 = 0.082$$

