

**MASSACHUSETTS INSTITUTE OF TECHNOLOGY**  
**Department of Civil and Environmental Engineering**

**1.017/1.010 Computing and Data Analysis for Environmental Applications /  
Uncertainty in Engineering**

---

**Problem Set 6: Estimates and Confidence Intervals**  
**Due: Thursday, Oct. 30, 2003**

---

Please turn in a hard copy of your MATLAB program as well as all printed outputs (tables, plots, etc.) required to solve the problem.

**Problem 1: Comparing Alternative Estimates of Population Properties**

Suppose that you have a sample of 10 observations of a random variable  $x$  which you believe to be exponentially distributed. Your objective is to estimate the 90 percentile value  $x_{90}$  of this variable. This value is the solution of the equation  $F_x(x_{90}) = 0.9$ .

Propose at least two different methods for estimating  $x_{90}$  from the 10 observations.

Compare the performance of these alternative estimators with a stochastic simulation that performs the following steps:

1. Generate many (e.g. 1000) replicates, each consisting of 10 observations drawn from an exponential distribution with parameter  $a = E[x]$  specified by you.
2. For each replicate derive an estimate  $\hat{x}_{90}$  of  $x_{90}$  from each of your two proposed estimators.
3. For each estimator compute the sample mean and variance of  $\hat{x}_{90}$  over all replicates. Also, for each estimator construct an  $\hat{x}_{90}$  histogram and an  $\hat{x}_{90}$  CDF plot (using MATLAB's `normplot` function).
4. Determine whether your estimators are unbiased and consistent (check consistency by plotting the rerunning your simulation for a much larger number of observations).

Which of your estimators is better? Explain your reasoning.

**Problem 2: The Bivariate Normal Distribution**

**Two dependent normally distributed** random variables (parameters  $\mu_x$ ,  $\mu_y$ ,  $\sigma_x$ ,  $\sigma_y$ , and  $\rho$ ):

$$f_{xy}(x, y) = \frac{1}{2\pi|C|^{0.5}} \exp\left\{-\left[\frac{(Z - \mu)'C^{-1}(Z - \mu)}{2}\right]\right\}$$

$Z =$  vector of **random variables**  $= [x \ y]'$

$\mu =$  vector of **means**  $= [E(x) \ E(y)]'$

$$C = \text{covariance matrix} = C = \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}$$

$\sigma_x = \text{Std}(x)$ ,  $\sigma_y = \text{Std}(y)$ ,  $\rho = \text{Correl}(x,y)$

$|C| =$  **determinant** of  $C = \sigma_x^2\sigma_y^2(1-\rho^2)$

$$C^{-1} = \text{inverse of } C = \frac{1}{|C|} \begin{bmatrix} \sigma_y^2 & -\rho\sigma_x\sigma_y \\ -\rho\sigma_x\sigma_y & \sigma_x^2 \end{bmatrix}$$

Note that the ' symbol is used to indicate the vector transpose in the bivariate normal probability density expression. The argument of the exponential in this expression is a scalar.

In this problem you will use the MATLAB function `mvnrnd` to generate scatterplots of correlated bivariate normal samples. This function takes as arguments the means of  $x$  and  $y$  and the covariance matrix defined above (called `SIGMA` in the MATLAB documentation).

Assume  $E[x] = 0$ ,  $E[y] = 0$ ,  $\sigma_x = 1$ ,  $\sigma_y = 0$ . Use `mvnrnd` to generate 100  $(x, y)$  realizations. Use `plot` to plot each of these as a point on the  $(x,y)$  plane (do not connect the points). Vary the correlation coefficient  $\rho$  to examine its effect on the scatter. Consider  $\rho = 0., 0.5, 0.9$ . Use `subplot` to put plots for all 3  $\rho$  values on one page.

### **Problem 3: Effect of Sample Size on Estimate Accuracy**

Reconsider the arsenic data set from Problem Set 4. Estimate the mean of the complete data set (population) from smaller samples of size  $N$ , randomly selected (without replacement) from the complete data set. Compute the sample mean and standard deviation for  $N = 4, 8, 32, 64,$  and  $128$ . Plot the differences between the sample and population means and standard deviations (on two different plots) as functions of  $N$ . Explain your results.

### **Problem 4: Confidence Intervals**

The following random sample was drawn from a continuous probability distribution  $F_X(x)$ . Estimate the mean and standard deviation of this distribution and specify 90%, 95%, and 99% confidence intervals for the mean.

$x = [ 2.6287 \ 7.0923 \ 2.3959 \ 0.4207 \ 2.8124 \ 4.1257 \ 3.1121 \ 0.8913$   
 $1.2885 \ 0.1863 \ 0.5489 \ 2.2652 \ 1.3867 \ 8.5322 \ 1.8364 \ 2.3576$

0.4417 0.4693 2.2507 0.7189 ]