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

## 1.017 Computing and Data Analysis for Environmental Applications

**Quiz 3 (Solutions provided at the end of each problem)**
Tuesday, December 10, 2002

Please answer all questions on a separate piece(s) of paper with your name clearly identified:

**Problem 1 ( 15 points)**

A field survey attempts to relate nitrogen runoff (in 1000 kg/year) in a small watershed to 2 factors: A) pesticide use (none, light, moderate, or heavy) and B) soil type (sandy, sand-silt, silty clay, or clay). Answer the following questions about the two-way ANOVA table given below:

| Source | SS | df | MS=SS/df | $F$ | $p$ |
|---|---|---|---|---|---|
| Factor A | 0.0403 | 3 | 0.0134 | 0.6992 | 0.5661 |
| Factor B | 0.3044 | 3 | 0.1015 | 5.2856 | 0.0100 |
| Interaction AB | 0.3265 | 9 | 0.0363 | 1.8896 | 0.1278 |
| Error | 0.3072 | 16 | 0.0192 | | |
| Total | 0.9783 | 31 | | | |

1. Which, if any, of the two factors contributes significantly to nitrogen runoff?  Why?
2. Are interactions between factors A and B significant?  Why?

**Problem 1 Solution:**

1.) Soil type contributes significantly to nitrogen runoff because it has a small $p$ value (<0.05).
2.) The interactions are not that significant since $p$ is larger than 0.05.  But 0.1278 is not that large, so you can say that there may be some interaction.

**Problem 2 ( 15 points)**

Species diversity indices are frequently used to measure the health of natural ecosystems.  Suppose that you are given the following two sets of unitless diversity indices for a control site and a site affected by human activity:

Control:  2.72  1.98   8.13   0.42   4.17   6.66

Affected: 2.78  3.02   0.93   3.21

Use the normal probability plot provided at the end of this quiz to determine the $p$ value for a two-sided large-sample test of the null hypothesis that the two means are the same. Do you think the control and affected populations are significantly different? Why? Is the large-sample assumption valid here? Why?

**Problem 2 Solution:**

Calculations may be expressed in terms of MATLAB syntax:

```
control = [2.72 1.98 8.13 .42 4.17 6.66]
affected = [2.78 3.02 .93 3.21]
mc = mean(control) = 4.013
ma = mean(affected) = 2.485
vc = var(control) = 1.11
va = var(affected) = 8.54
sc = std(control) = 2.92
sa = std(affected) = 1.05
zo = (mc-ma)/sqrt(vc/6+va/4) = 1.17
p = 2 * (1-normcdf(zo)) = .24
```

The control and affected populations are not significantly different because p is not small. Although we used the unit normal CDF in this problem, the large-sample assumption is probably not valid here since there are only four affected data points.

**Problem 3 ( 25 points)**

Provide specific one-sentence answers to the following questions:

a) Why did we use the transformation $C_T = ln(C+1)$ when carrying out an ANOVA of Boston Harbor coliform data?
b) Explain the difference between independent and dependent variables in a regression analysis.
c) Why is the sum-of-squared differences between the observed and modeled dependent variables a reasonable measure of "goodness-of-fit"? Suggest at least one other measure that could also be reasonable in some applications.
d) Suppose that there is a significant linear relationship between the area of a temperate watershed and the average annual runoff. Is a set of average annual runoff values selected at random from temperate watersheds of different sizes a random sample? Why?

e) Suppose that the sample mean of a set of 5 data points $x_1,...,x_5$ is 6.0 and the sample standard deviation is 2.5. Compare the $p$ values obtained from small and large sample tests of the hypothesis H0: $E[x]=0$ (i.e. which test will give a larger p value?). You do not need to compute the actual p values .... just rank the two possibilities.

**Problem 3 Solution:**

a.) We used the transformation because the original data is not normally distributed, but the log of the data is nearly normal.
b.) The independent variable varies with the dependent variable in a specific way determined by the regression parameters.

c.) The sum of squared differences is a reasonable measure of goodness of fit because it can only be zero when all the deviations are zero. One other possible measure is the absolute value of the deviation.

d.) This is not a random sample because they are not independent.

e.) The p value from the large sample test will be smaller than the p from the small sample test because the normal distribution is narrower than the t-distribution.

## Problem 4 ( 20 points)

Consider the EPA NOx (nitrous oxide) emissions data set attached to this quiz. Organize the data into groups appropriate for a one-way analysis of variance that tests the influence of fuel type on NOx emissions. Use only three replicates for each treatment level. Indicate directly on the data sheet (by writing two numbers at the end of the appropriate row) the treatment level and replicate for each of the sample you select for your ANOVA.

Next use your annotated data sheet to read off the values of the input data array required by the MATLAB function `anova1` (documentation attached). Please read the documentation carefully to make sure that you define the array properly.

Finally, construct a one-way ANOVA table with the degree of freedom values filled in for all table rows. Identify the quantity (e.g. sum-of-squares, etc.) that goes in each of the remaining table entries but do not compute numerical values for any quantities other than the number of degrees of freedom.

### Problem 4 Solution:

```
% Factor = fuel type, 4 treatments
PNG = [53.7,18.8,0]
C = [1032,927.1,2411.6]
Oil = [136.6,1407.9 8.5]
DSL = [7.6,2.3,15.1]
matrix = [PNG',C',oil',DSL']
anova1(matrix)
```

## Problem 5 ( 25 points)

Suppose that you are given the following data describing the concentration (in mg/L) of decaying organic carbon remaining in a treatment tank after the indicated time has passed:

| Time | 10 | 12 | 15 | 17 | 22 | 24 |
|------|------|------|------|------|------|------|
| concentration | 1.77 | 0.88 | 1.08 | 0.23 | 0.43 | 0.34 |

Also, suppose that you model the decaying organic carbon by the following exponential function:

$$C(t) = C(0)e^{-rt}$$

This equation can be expressed as a linear regression model (with dependent variable $y(t) = log_eC(t)$ and unknown regression parameters $a_1 = log_eC(0)$ and $a_2 = -r$ ) if you take the $log_e$ of each side and add a random measurement error to the result. You can use a regression approach to estimate the unknown regression parameters from the tabulated data and then to check the model's ability to explain observed temporal changes in concentration. The results of this regression are summarized in the ANOVA table provided below.

| Source | SS | df | MS=SS/df | $F$ | $p$ |
|--------|------|----|----------|------|-------|
| Regression | 1.09 | 1 | 1.09 | 7.15 | 0.056 |
| Error | 0.61 | 4 | 0.15 | | |
| Total | 1.70 | 5 | | | |

Carry out the following tasks:

1) Define the arrays needed to carry out the regression analysis with the MATLAB function `regress` (documentation attached). Provide specific numerical values inferred from the data.

2) Use the $F$ and $p$ values in the table to determine whether the regression for this problem is significant (i.e. does the exponential model provide a good explanation of observed temporal variability).

3) Calculate the $R^2$ value for this regression. Does it suggest that the exponential fit is good or bad?
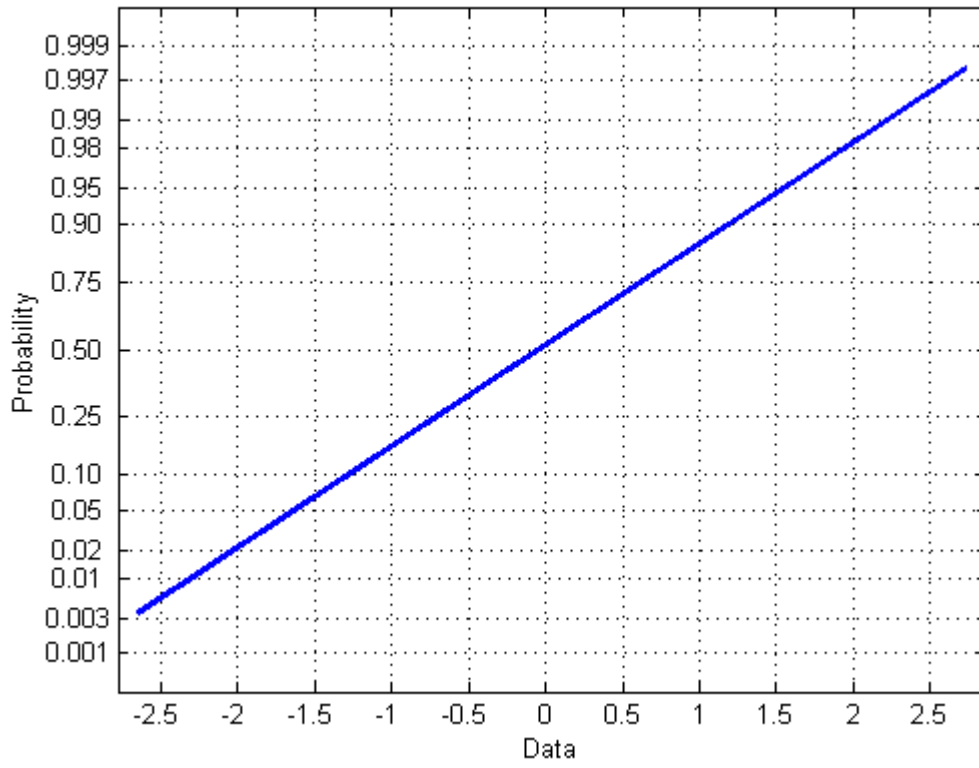
**Problem 5 Solution:**

```
1.)
time = [10 12 15 17 22 24]
conc = [1.77 .88 1.08 .23 .43 .34]
lnconc = log(conc);
x=[ones(6,1), time'];
[B,BINT,R,RINT,STATS] =regress(lnconc',x)
```

2.) $F = 7.15$, $p = .056$:
The F and p values are borderline, but the model provides a fair explanation of the observed temporal variability.

3.) $R^2$ = SSR/SST = 1.09/1.70 = 0.641 -- the exponential fit is pretty good but not great.

Normal Probability Plot

# Data for Problem 4

U.S.EPA - EMISSIONS TRACKING SYSTEM (ETS)
Preliminary Ozone Season Values for May to September, 2002
Report for Massachusetts

| | Unit/ Stack ID | Boiler Type | Primary Fuel | NOx Monitoring Methodology | NOx Controls | Operating Time (hours) | Heat Input (mmBtu) | NOx Emissions (tons) |
|---|---|---|---|---|---|---|---|---|
| ANP Bellingham Energy Proj. | 1 | CC | PNG | CEM | DLNB | 759 | 993533 | 53.7 |
| | 2 | CC | PNG | CEM | SCR | 175 | 70315 | 25.9 |
| ANP Blackstone Energy Co. | 1 | CC | PNG | CEM | DLNB | 2610 | 3949647 | 18.8 |
| | 2 | CC | PNG | CEM | DLNB | 2733 | 4123868 | 20.3 |
| Bellingham | CP1 | | | | | 3672 | 9220740 | 0 |
| | CS1 | | | | | 3672 | 9220740 | 394.7 |
| | 1 | CC | PNG | CEM | STM | 3668 | 4724269 | 0 |
| | 2 | CC | PNG | CEM | STM | 3672 | 4496471 | 0 |
| Berkshire Power | 1 | CC | PNG | CEM | H2O | 2412 | 3764556 | 19.5 |
| Blackstone | CP2 | | | | | 30 | 2746 | 0 |
| | CS2 | | | | | 1324 | 71668 | 6.3 |
| | 11 | DB | PNG | CEM | CM | 825 | 43473 | 0 |
| | 12 | DB | PNG | CEM | CM | 619 | 28195 | 0 |
| Brayton Point | 1 | T | C | CEM | LNC3 | 3593 | 7215909 | 1032 |
| | 2 | T | C | CEM | LNC3 | 3522 | 6927859 | 927.1 |
| | 3 | DB | C | CEM | LNBO | 2438 | 11878591 | 2411.6 |
| | 4 | DB | OIL | CEM | LNB | 843 | 1364761 | 136.6 |
| Canal Station | 1 | DB | OIL | CEM | LNBO | 3348 | 11782280 | 1407.9 |
| | 2 | DB | OIL | CEM | LNBO | 2296 | 5761147 | 654.5 |
| Cleary Flood | 8 | DB | OIL | CEM | LNB | 255 | 63200 | 8.5 |
| | 9 | OB | PNG | CEM | LNBO | 1362 | 1258594 | 87.4 |
| Dartmouth Power | 1 | CC | PNG | CEM | H2O | 2602 | 1364167 | 20.6 |
| Deer Island Treatment | S42 | CT | DSL | AE | H2O | 103 | 20330 | 7.6 |
| | S43 | CT | DSL | AE | H2O | 26 | 3275 | 2.3 |
| Dighton | 1 | CC | PNG | CEM | SCR | 3518 | 4139225 | 23.9 |
| Doreen | 10 | CT | DSL | NOXG | | 14 | 3416 | 2 |
| Framingham Station | FJ-1 | CT | DSL | NOXU | | 41 | 7293 | 2.3 |
| | FJ-2 | CT | DSL | NOXU | | 62 | 8718 | 3.5 |
| | FJ-3 | CT | DSL | NOXG | | 150 | 25155 | 15.1 |
| GE Aircraft Engines Lynn | 3 | DB | PNG | CEM | LNBO | 3595 | 875932 | 68.6 |
| | 5 | CC | PNG | CEM | STM | 684 | 136049 | 9.4 |
| Indeck Pepperell | CC1 | CC | PNG | CEM | STM | 1389 | 348406 | 22.8 |