# Biostatistics 140.623
# Third Term, 2017-2018
# Problem Set 2

### Survival in Primary Biliary Cirrhosis

**Learning Objectives:**
Students who successfully complete this section will be able to:
- To evaluate whether the drug DPCA prolongs life in patients.
- To identify baseline characteristics of patients which predict longer survival.
- Analyze the survival time data (without grouping) by the Kaplan-Meier estimate of the survival function, the log- rank statistic, and Cox proportional hazards model.
- Check the estimated model for its consistency with the observed data; in particular, check the proportional hazards assumption using the complementary log-log plot of the estimated survival function.
- Summarize the findings for public health readers and document and archive the steps of the statistical analysis by creating a `Stata` do-file.

**Data Set**:
Between January 1974 and May 1984, a double-blinded randomized trial on patients with primary biliary cirrhosis (PBC) of the liver was conducted at the Mayo clinic.  A total of 312 patients were randomized to either receive the drug D-penicillin (DPCA) or a placebo.  Patients were followed until they died from PBC or until censoring, either because of administrative censoring (withdrawn alive at end of study), death not attributable to PBC, liver transplantation, or loss to follow-up.  At baseline, a large number of clinical, biochemical, serological and histologic measurements were recorded on each patient.  This data set is a subset of the original data, and includes information on each patient's time to death or censoring, treatment, age, gender, serum bilirubin, and histologic disease stage (1-4).

The variables included in this dataset include:
**case**: unique patient ID number
**sex**: 0 = male, 1 = female
**drug**: 0 = placebo, 1 = DPCA
**bil** : serum bilirubin in mg/dl
**survyr**:  time (in years) to death or censoring
**death**:  indicator = 1 if patient died, 0 if censored
**ageyr**: age in years [continuous variable]
**histo**: histologic disease stage (1 – 4) [categorical variable]
Also included in the data set for your possible use are the following indicator (dummy) variables:

Age Indicators:
       **agegr_2**:  1 if patient is 45-55 years old, 0 otherwise
       **agegr_3**:  1 if patient is >= 55 years old, 0 otherwise

Histologic Stage Indicators:
>    **hstage2**: 1 if patient is in Stage 2, 0 otherwise
>    **hstage3**: 1 if patient is in Stage 3, 0 otherwise
>    **hstage4**: 1 if patient is in Stage 4, 0 otherwise

The data are stored in the `Stata` data set *pbctrial.dta*, which may be downloaded from the course website using Internet Explorer. This data set is already in `stset` format. If you need more help with the `Stata` survival data commands, use the Help menu for the command `sts`.

**Methods**:
Use the data set described above and the appropriate statistical analyses to address the specific learning objectives listed on the first page.

**Hints**:
(In the following list of commands, the word *varname* indicates that the command must be followed by the name of at least one variable stored in the dataset: ie: if you type `tab` you will get an error, but if you type `tab sex` you will get output)

a.  Explore the data using descriptive statistics:
```
tab varname
summarize
```

b.  Explore differences in time to death by different baseline variables using graphs and complementary log-log plots.

```
sts list
sts graph
sts test varname
stphplot, by(varname)
```

c.  Fit several Cox proportional hazards regression models to the ungrouped survival data:
```
stcox varname   etc.
```

d.  Using your review window or `Stata` log file, create a `Stata` do-file that documents and archives the steps of your statistical analysis. This file will make your analysis "reproducible."

e.  **Summarize your findings in a brief report** (less than two pages with at most one table and one figure) as if for a biomedical/public health journal.

>   A **suggested format** is:
>   - Introduction – a few sentences about the research question(s)
>   - Data description – simple tabulations describing patient characteristics
>   - Results from multiple models that address question(s) (e.g., bivariate and multivariable)
>   - Graphical display that presents evidence in the data relevant to your scientific question.