# Problem set 1

Your answers are due on September 24, 2025 at 1pm. You should submit them through Gradescope.

1. A health company claims that a new vitamin supplement increases concentration levels in adults. To test this claim, a researcher recruits two independent groups of participants. One group (the treatment group) takes the supplement daily for a month. The second group (the control group) receives a placebo.

   At the end of the month, both groups are administered the same standardized concentration test, scored out of 100 points. The scores are assumed to be normally distributed with equal but unknown variances.

   The average score in the **treatment group** ($n = 35$) is 78, with a sample standard deviation of 12. In the **control group** ($n = 30$), the average score is 72, with a sample standard deviation of 10. A statistician sets up the hypotheses as follows:

   $$H_0 : \mu_{\text{treat}} = \mu_{\text{control}} \quad \text{vs.} \quad H_1 : \mu_{\text{treat}} > \mu_{\text{control}}.$$

   (a) Sketch the sampling distribution of the difference in means under the null and alternative hypotheses.

   (b) What is the standard error of the difference in sample means? What is the distribution of the test statistic under the null?

   (c) Compute the $t$-statistic for the observed data.

   (d) What is the critical value for a one-sided test at the 5% significance level? Should the researcher reject the null?

   (e) What is the $p$-value corresponding to the observed $t$-statistic?

   (f) Suppose the actual difference in population means is 8 points. What is the power of the test in that case? What steps could increase the power?

2. In the following two problems you will investigate whether money influences elections, and learn how to work in STATA.

   The file VOTE1.dta (available on Canvas) contains data on election outcomes and campaign expenditures for 173 two-party races for the U.S. House of Representatives

in 1988. There are two candidates in each race, A and B. Let *voteA* be the percentage of the votes received by Candidate A and *shareA* be the percentage of total campaign expenditures accounted for by Candidate A. Many factors other than *shareA* affect the election outcome (including the quality of the candidates and possibly the dollar amount spent by A and B). Nevertheless, we can estimate a simple regression model to find out whether spending more relative to one's challenger implies a higher percentage of the vote. Denote the variable X to be *shareA* and the variable Y to be *voteA*.

Do the following tasks using STATA:

(a) Construct variable $X$ from the data available to you;

(b) Calculate the sample means of X and Y;

(c) Calculate the sample standard deviations of X and Y and the sample correlation coefficient between X and Y;

(d) Produce the OLS estimated regression coefficients from the regression $Y_i = \beta_0 + \beta_1 X_i + u_i$;

(e) Create a new variable $\hat{Y}_i$ $i = 1, ..., n$, containing the predicted vote for each election based on the regression above;

(f) Create a new variable $\hat{u}_i$ containing the OLS residual for each election;

(g) Calculate the sum of $\hat{u}_i$ and explain why it should be zero;

(h) Graph the scatterplot of the data points and the regression line.

3. We continue to work with VOTE1 data set. Estimate a regression of votes on spending share, using the "`robust`" option.

(a) What is the estimated slope? Explain in words what it means. Is the estimated effect of spending on share large or small? Explain what you mean by "large" or "small".

(b) Report the 95% confidence interval for $\beta_1$, the slope of the population regression line.

(c) Does spending explain a large fraction of the variance in vote? Explain.

(d) Look at the correlation coefficient between share and vote computed in the previous problem, and compare its square to the $R^2$. How are they related? Provide a simple mathematical derivation of this fact.

(e) What is the root mean squared error of the regression? What does this mean?

(f) Based on your graph from 2(h), does the error term appear to be homoskedastic or heteroskedastic?

(g) Run the regression again without the "robust" option. Compare the results to what you obtained with the "robust" option. What is the same and what is different?

STATA HINTS. Note that STATA has on-line help. The following commands will be useful:

| | |
|---|---|
| `list` | lists the data |
| `summarize` | computes sample means and standard deviations (the option ",detail" gives additional statistics, including the sample variance) |
| `correlate` | produces correlation coefficients (with the option ", covariance" this command produces covariances) |
| `regress` | estimates regression by OLS |
| `predict` | computes OLS predicted values and residuals |
| `generate` | creates a new variable |
| `scatter` | creates a scatter-plot. There is an option to connect the dots by adding the option ",connect()" |