# Pol Sci 630: Problem Set 9 Heteroskedasticity

Prepared by: Anh Le (anh.le@duke.edu)

Due Date: Nov 2 (Beginning of Class)
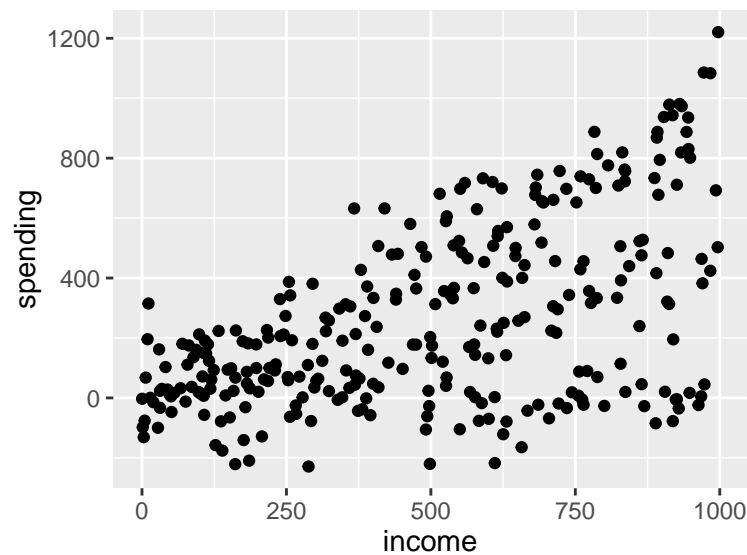
```
rm(list = ls())
library(ggplot2)
```

## 1   Heteroskedasticity

This exercise nudges you to think about heteroskedasticity as a theoretical / social science problem, not a mechanical / statistical issue to be blindly fixed.
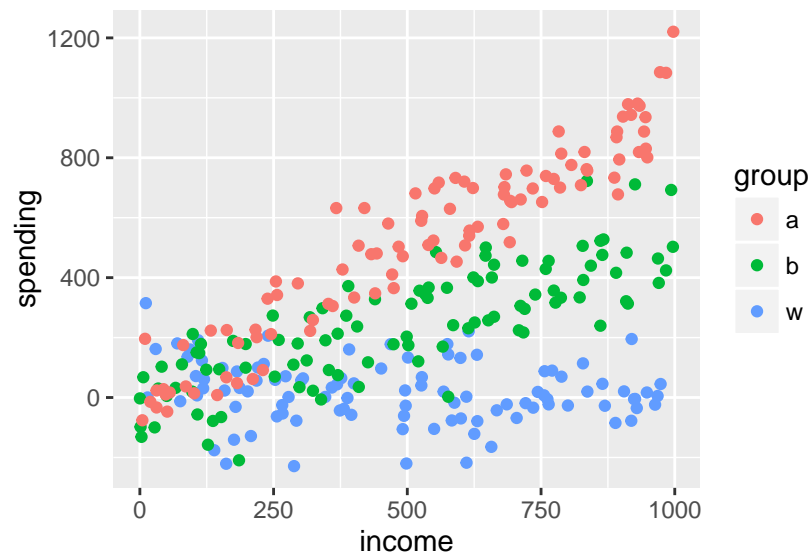
One common cause of heteroskedasticity is that our model does not take into account heterogenous effect across sub-populations. For example, we have a model of spending (dependent var) as a function of income (independent var), and the propensity to spend differs across ethnic groups. Formally,

$$spending = \beta_{ethnic} income + \epsilon \tag{1}$$

where $\beta_{ethnic}$ takes a different value for white, black, and asian. If we don't know about this heterogeneity of propensity to spend across ethnic groups, the graph will show heteroskedasticity:

Buf if we are smart researcher, we'll realize the underlying cause of the heterogeneity, as shown in the following plot:



The take-home point is that heteroskedasticity could be a signal of underlying model specification, and we should think hard about the cause of heteroskedasticity instead of applying a quick fix.

## 1.1 Simulating

Simulate the spending and income pattern for three ethnic groups as described above. (Try to) Re-create the two plots above (doesn't have to be ggplot2). The numbers don't have to be the same – just make sure that your data has heteroskedasticity due to underlying heterogenous effect across ethnic groups as described in the example above. Note: Don't look at my code.

## 1.2 Diagnostics: Visual

Using the simulated data above, regress spending on income, plot the residual against the predicted value.

## 1.3 Diagonistics: Hypothesis test

Conduct BP test and White test. Why do the tests reach the same conclusion here, unlike in the lab tutorial?

## 1.4 Diagnostics: Repeat the White's test manually

Here's the instruction. Compare the result you get doing it by hand vs using R.
   *White test (Wooldridge "Introductory", Testing for heteroskedasticity)*
   1. Estimate the model `y ~ x_1 + x_2 + ... + x_k` by OLS, as usual. Obtain the OLS residual $\hat{u}$ and the fitted values $\hat{y}$. Compute $\hat{u}^2$ and $\hat{y}^2$.
   2. Run the regression $\hat{u}^2 = \delta_0 + \delta_1\hat{y} + \delta_2\hat{y}^2$. Keep the R square.
   3. **I want you to use the LM for this problem** Form either the F or LM statistic and compute the p-value (using the $F_{2,n-3}$ distribution in the former case and the $\chi_2^2$ distribution in the latter case).

## 1.5 Fixing: robust standard error

Run hypothesis test without and with robust standard error. What's the conclusion?

## 1.6 Fixing: calculate robust standard error by hand

Show that it's the same as given by R.

## 1.7 Fixing: Provide a correct model

Specify a regression model that takes into account heterogenous effect of income on spending across ethnic groups. Show that there's no longer heteroskedasticity.