# IE8990: Adv. Data Analytics for Complex Systems

- Lab 1
  - Getting started with R and Rstudio
  - Use R to fit regression model and perform diagnostics

# R and Rstudio

- R Download
  - https://cloud.r-project.org/

- Rstudio Download
  - https://www.rstudio.com/products/rstudio/download/

# Rstudio interface

# R notebook and markdown

- The structure and functions of R notebook and R markdown are almost the same.

- Coding is exactly the same for these two.

- One major difference is that R notebook is able to provide a preview of the generate report.

# Example 1: Iris data

- Load the data
- Data visualization
- Fit a linear regression model using Petal. Length as response and Petal.Width as predictor
- Evaluate the performance of this model
- Find confidence interval of the parameters
- Diagnostics
  - Residual plots
  - Statistical tests

**How do we interpret the results?**

MISSISSIPPI STATE
UNIVERSITY™

# Data visualization

- Response variable
  - Sepal.Length

**Iris Data**

# Homework 1: Question 1 – Iris data cont.

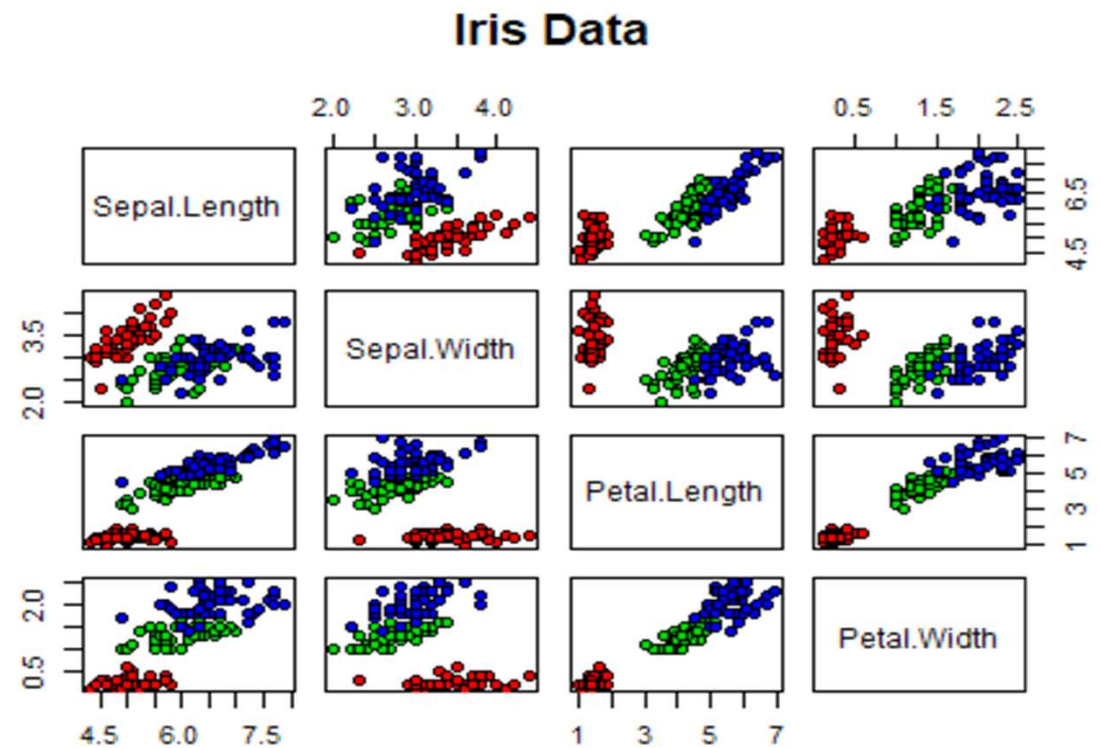1. Fit Model 0 using Sepal.Length as response and Sepal.Width as the only one predictor
   a. Find all the parameters, write down their interpretations
   b. Evaluate Model 0 using adj $R^2$ as the criteria
   c. Is Model 0 good for prediction? If no, how do you want to improve the model?

2. Based on your answer in 1c, fit a new model (Model 1) using Sepal.Length as response, Sepal.Width and Species as predictors
   a. Find all the parameters, write down their interpretations
   b. Is Model 1 a good model for prediction?

3. Can you come up with another model that have better adj $R^2$ than Model 1?

We will discuss what models you have fitted next Thursday.