

# Unit 9 Homework: Large-Sample Regression Theory

w203: Statistics for Data Science

## What Makes a Successful Video Game?

The file `video_games.csv` contains data on 1212 video games that were on sold in 2011. It was compiled by Joe Cox, an economist at the University of Portsmouth.

Three key variables are as follows:

Variable	Meaning
Metrics.Sales	The total sales, measured in millions of dollars.
Metrics.Review.Score	Metacritic review score, an indicator of quality, out of 100.
Length.Completionists.Average	The mean time that players reported completing everything in the game, in hours.

You can find an explanation of other variables at [https://think.cs.vt.edu/corgis/csv/video\\_games/](https://think.cs.vt.edu/corgis/csv/video_games/).

You want to fit a regression predicting `Metrics.Sales`, with `Metrics.Review.Score` and `Length.Completionists.Average` as predictors.

0. Rename the variables that you are going to use to something sensible – variable names that have both periods and capital letters are not sensible. :fire: Better would be, for example changing `Metrics.Sales` to just `sales`.
1. Examining the data, and using your background knowledge, evaluate the assumptions of the large-sample linear model.
2. Whether you consider the large-sample linear model sufficiently valid or not, proceed to fit the linear model using `lm()`.
3. Examine the coefficient for `Metrics.Review.Score` and give an interpretation of what it means.
4. Perform a hypothesis test to assess whether video game quality has a relationship with total sales. Please use `vcovHC` from the `sandwich` package with the default options (“HC3”) to compute robust standard errors. To conduct the test, use `coeftest` from the `lmtest` package.
5. How many more sales does your model predict for a game one standard-deviation higher than the mean review, vs. a game one standard-deviation lower than the mean review, holding all else equal? Answer this in two different ways:
  - (a) Compute the standard deviation of the review score, and multiply the appropriate model coefficient by two-times this standard deviation.
  - (b) Use the `predict` function with the model that you have estimated. You can read the documentation for `predict.lm` which is the predict method for linear model objects (the type that you have fit here). Include a data frame (that has the same variable names as the data frame that you fitted the model against) in the `newdata` argument to `predict`. This data frame should have two rows and two columns. The column for the reviews should change from  $\mu - \sigma$  to  $\mu + \sigma$ ; the column for the play time should be set to a constant, sensible level (perhaps the  $\mu$  of this variable).

5. **Optional:** Open the attached paper by Joe Cox, and read section 3. Which assumption did the author focus on, and why do you think that is?

*Note: Maximum score on any homework is 100%*