# STATS 202: Data Mining and Analysis
Instructor: Linh Tran

HOMEWORK # 1
*Due date: July 7, 2023*

Stanford University

## Introduction

Homework problems are selected from the course textbook: *An Introduction to Statistical Learning*.

### Problem 1 (4 points)

Chapter 2, Exercise 2 (p. 52).

### Problem 2 (4 points)

Chapter 2, Exercise 3 (p. 53).

### Problem 3 (4 points)

Chapter 2, Exercise 7 (p. 54).

### Problem 4 (4 points)

Chapter 12, Exercise 1 (p. 548).

### Problem 5 (4 points)

Chapter 12, Exercise 2 (p. 548).

### Problem 6 (4 points)

Chapter 12, Exercise 4 (p. 549).

### Problem 7 (4 points)

Chapter 12, Exercise 9 (p. 550).

### Problem 8 (4 points)

Chapter 3, Exercise 4 (p. 122).

### Problem 9 (4 points)

Chapter 3, Exercise 9 (p. 123). In parts (e) and (f), you need only try a few interactions and transformations.

## Problem 10 (4 points)

Chapter 3, Exercise 14 (p. 127).

## Problem 11 (5 points)

Let $x_1, \ldots, x_n$ be a fixed set of input points and $y_i = f(x_i) + \epsilon_i$, where $\epsilon_i \overset{iid}{\sim} P_\epsilon$ with $\mathbb{E}(\epsilon_i) = 0$ and $\text{Var}(\epsilon_i) < \infty$. Prove that the MSE of a regression estimate $\hat{f}$ fit to $(x_1, y_1), \ldots, (x_n, y_n)$ for a <u>random</u> test point $x_0$ or $\mathbb{E}\left(y_0 - \hat{f}(x_0)\right)^2$ decomposes into variance, square bias, and irreducible error components. *Hint: You can apply the bias-variance decomposition proved in class.*

## Problem 12 (5 points)

Consider the regression through the origin model (i.e. with no intercept):

$$y_i = \beta x_i + \epsilon_i \tag{1}$$

(a) *(1 point)* Find the least squares estimate for $\beta$.

(b) *(2 points)* Assume $\epsilon_i \overset{iid}{\sim} P_\epsilon$ such that $\mathbb{E}(\epsilon_i) = 0$ and $\text{Var}(\epsilon_i) = \sigma^2 < \infty$. Find the standard error of the estimate.

(c) *(2 points)* Find conditions that guarantee that the estimator is consistent. *n.b. An estimator $\hat{\beta}_n$ of a parameter $\beta$ is consistent if $\hat{\beta} \overset{p}{\to} \beta$, i.e. if the estimator converges to the parameter value in probability.*