Suppose that we are interested in understanding the effect of $X_1$ ~~change~~ on $Y$. The model is

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

(a) Would you include the variable $X_1 \times X_3$ as a regressor in the model? Explain this in detail.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 (X_1 \times X_3) + u$$

Interaction terms are often added to models to capture potential interaction effects

(b) Suppose that in the model in (1) all the assumptions for the unbiasedness of OLS hold, but the errors are heteroscedastic. What would you do? Make sure that you mention the reasons for your decision.

If errors are heteroscedastic, the ordinary least squares (OLS) estimates remain unbiased, but they are no longer the best linear unbiased estimates (BLUE) because their variance is not minimized.

(c) The model is far from complete. Propose 2 other variables which should have been included?
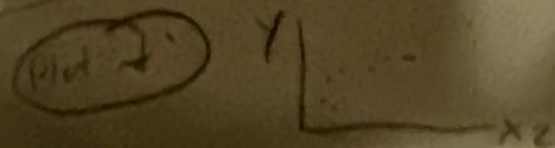
- Interest rates, inflation          Income/wealth

- time                               Geographic region

(d) Do you believe it is appropriate for the model to be linear in the variable $X_2$? Discuss this matter in detail.

(Plot 2) $Y$ [plot] $X_2$    + relation seems straight, (ii) look
or check $R^2$

(e) suppose that $\hat{\beta}_1 = 0.04$ and $SE(\hat{\beta}) = 0.015$. A 5% significance level, can we say that $x_1$ affects $y$? why:

$H_0: \beta_1 = 0$ (no effect)
$H_1: \beta_1 \neq 0$ (no affect)

$t = \dfrac{0.04}{0.015} = 2.67$

Significance level = 1.96

Since $|t| = 2.67 > 1.96$, we reject $H_0$. $x_1$ has significant effect on $y$.

(f) Based on the hypothesis test from (e), how does the result impact policy? Which recommendation do you have?

Since $x_1$ has significant effect on $y$, suggest a relevant change may matter but won't solve societal issues

(g) How would you test the hypothesis that the variable $x_2$ is superfluous in the model.

$H_0: \beta_2 = 0$     perform F-test
$H_1: \beta_2 \neq 0$

(i) How would you reply to the following statement: "Statistics about one or more coefficients" in a scientific name?

-To test hypothesis test about coefficients
• State Null and alternative hypothesis
-Use test statistics
-Evaluate p-values against significance levels and draw conclusion
Interpret results with consideration of assumption.

(J) Suppose that the true model is actually

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_3 + u$$

but you run the model in (1) what is the consequence of this?

leads to misspecification bias which leads to omitted variable bias. The estimated coefficients for the remaining variables may be bias and inconsistent b/c the model will not capture the curve between $y$ and $x_1$

(K) What is the results of $R^2$ od a regression? How do you expect $R^2$ to change from the model in (2) to the one in (1)?

$R^2$ represents the proportion of variance in $y$ and $x$

$R^2$ will decrease b/c here omitting $x_1$ in model 2 which lead to more variation in $y$ the model (1).

(L) Assume the estimated coefficients for $\beta_1$ and $\beta_2$ in (2) are 0.04 and 0.001 respectively and $x_1 = 1$, keeping everything else constant, what is the predicted change in $y$ when we increase $x_1$ by 1 unit. What about when $x_1 = 5$

$$x_1 = 1: \quad 0.04(1) + 0.001(1)^2 = 0.04 + 0.001 = 0.041$$

$$x_1 = 5: \quad 0.04(5) + 0.001(5)^2 = 0.20 + 0.025 = 0.225$$

(M) Suppose that our data shows $y$ is self reported. what problem we may have? What is the consequence of this problem?

Lead to measurement error and inaccuracies in model, produce biased and inconsistent parameter in regression model.

(a) Suppose that $X_1$ is Nully for 5% of the observations. What problems do you fear? How would you handle this? We will face selection Bias, to handle this we perform Stratified sampling.

(b) $X_4$ is likely to affect $Y$. However, we do not have access to data about it. Do you expect the bias of $\hat{\beta}_1$ to be positive or negative?

If $X_4$ affects $Y$ and $X_1$, and is omitted, the bias of $\hat{\beta}_1$ depends on the correlation.

- if $X_4$ is positively correlated with $Y$ and $X_1$, omitting it leads to upward bias.

(c) $Y$ should depend on $X_5$. List at least one control variable that we could use as proxy for $X_5$. Explain your variable of choice.

Proxy - when $Y$ depends on $X_5$, but $X_5$ data is unavailable

Ex) $X_5 =$ interest rate, proxy is raw GDP

(d) Suppose that your data was collected only on certain type of area. Do you fear some form of sample selection? What would be the implication on the estimated coefficient of $X_1$? Elaborate on this issue.

data collected in a certain area is not representative of whole population leads to sample selection bias

Leads to biased estimates for $\hat{\beta}_1$ as it does not reflect entire population.

(V) Suppose that you have a lot of variables which are connected to the topic. It may be related to the topic substantially or tangentially. Explain your decision process about including those variables or not.

- We need to consider how relevant is the variables.
- We need also need to exclude highly collinear variables to avoid inflating standard errors.
- And recognize that too much variables can lead to overfilling.

# Essay Question

Prompt:

With the growth of industrialized food, high level of sugar consumption and large use of unneeded medications, American policy makers have become more and more concerned with the health of its population. Poor health jeopardizes people's welfare and generates high monetary costs for society.

In economics the link between health and income seems pretty uncontroversial. After all, healthier people can work longer and harder. Healthier children are likely to stay in school longer and learn more, earn more. Even across countries the relationship seems clear: those with better health are generally richer, and those that improve their citizen's health grow faster.

So, consider the following question: what is the effect of physical activity on type 2 diabetes? Your task is to describe how you might go about answering this question. Construct a model, and be specific about how you define your variables. Justify your choices and assumptions you need. Explain how you might go about finding the data you need. Comment shortcomings of your approach, and provide solutions you would implement when possible.

The goal of this essay is to develop an empirical study in order to understand the effect of physical exercise on the likelihood of developing Type 2 diabetes.

We will need two groups of people. Group 1 is the treatment group that will be assigned to do physical exercise for a certain amount of time per week. Group 2 will be the control, where they will live their regular lives. The allocation of people for this study will be done randomly. It will also be possible to divide the treatment group into subgroups of people based on hours of physical exercise per week / or frequency of exercise per week, and whether more exercise or frequency of exercise matters with the likelihood of being diagnosed with Type 2 Diabetes.

In our dataset needed for this experiment, we will need 1 column of hours or frequency of exercise done by individual, and whether the individual has been diagnosed with Type 2 diabetes.

$$Diabetes_{it+1} = \beta_1 Exer + \beta X + \varepsilon$$

Our dependent variable will be called Diabetes, and this variable will take value 1 if the individual is diagnosed with diabetes and 0 otherwise. For our continuous variables, Exer represents the average number of exercise done in a week.

Additionally, we would include controls for individuals personal characteristics— such as family history, race/ethnicity. This will mitigate the omitted variable bias to be correlated with exercise an affect the probability of developing type 2 diabetes.

$\beta_1$ can be interpreted as the change in the probability of being diagnosed with type 2 diabetes by increasing 1 hour of average weekly exercise.

There an required assumption to interpret the coefficient of $\beta_1$ as the causal effect of exercise on the probability of being diagnosed with type 2 diabetes, is this conditional with the other variables, is true nothing else that is correlated with the amount of exercise causing out by t that could also affect the probability of being diagnosed with type 2 diabetes.

One major shortcomings of this approach is the potential existence of variables that change the likelihood of being diagnosed with type 2 diabetes and are also correlated with exercise. For example, alcohol consumption or sugar consumption. If they were present in our dataset, we would include them as controls. An additional problem would be the existence of relatively large measurement error in Exercise. This measurement error would bias down our estimates of the effect of exercising on the probability of being diagnosed with diagnosed with type 2 diabetes.