

# A Causal Approach to Understanding Student Depression

**Maryam Almahasnah**  
malmahasnah@ucsd.edu

**Kevin Chan**  
kec021@ucsd.edu

**Hunter Brownell**  
hbrownell@ucsd.edu

**Jelena Bradic**  
jbradic@ucsd.edu

**Biwei Huang**  
bih007@ucsd.edu

## Abstract

Our project aims to establish a framework for the root causes of depression, specifically for students. To do this, we utilize various causal mechanisms to investigate and analyze features that could potentially be directly related to depression. While past projects have focused on uncontrollable factors, our goal is to focus on modifiable aspects. The ultimate goal here is to understand exactly what aspects of a potential student's life may be contributing to depression, and set a foundation to make improvements and emphasize the importance of providing targeted resources for students' mental health.

Code: <https://github.com/malmahasnah/capstone-team2-winter-quarter>

1	Introduction . . . . .	2
2	Methods . . . . .	2
3	Results . . . . .	6
4	Discussion . . . . .	10
5	Conclusion . . . . .	10
	References . . . . .	10
6	Contributions . . . . .	A1
	Appendices . . . . .	A1

# 1 Introduction

One of the main difficulties in understanding depression is differentiating its causes and its outcomes. Mental health in the lives of students is not only challenging to understand, but has continuously impacted students' physical health, academic performance, personal growth and relationships.<sup>1</sup> The main objective of our project is to understand the causal relationships between various factors and the role they play in depression. In doing so, we will establish a framework as to what factors are most important to focus on in order to prevent the onset of depression, and provide a road-map for what type of mental health resources should be provided for students.

The data we are using contains a variety of factors about a population of students and their lives. The dataset includes variables such as, but not limited to, their quantity of sleep, financial stress level, academic pressure, dietary habits, family history of mental illness, and a binary variable indicating whether or not they suffer from depression as well as suicidal thoughts. Since our plan is to aid students in their battle against depression, we plan to mainly study aspects of their lives that they have some level of control over. Focusing our efforts on features such as dietary habits and academic pressure can allow us to help students make more informed decisions. Ultimately, this will help better understand how to combat depression by identifying potential risk factors, seeking support systems as appropriate, and, most importantly, taking measures to target specific lifestyle factors to prevent depression.

Past studies such as "Causal Factors of Anxiety and Depression in College Students: Longitudinal Ecological Momentary Assessment and Causal Analysis Using PCMCI" [Huckins JF \(2020\)](#) have attempted to utilize methods such as the Peter and Clark algorithm to identify significant factors. While this approach is close to ours, we will also use the Fast KCI method for non-linear dependencies and incorporate counterfactual policy evaluation to investigate how changes in certain variables could impact depression outcomes. Other studies like ([Howlett JR \(2013\)](#)) and ([Broomhall AG \(2017\)](#)) have also utilized counterfactual reasoning. In our approach, we use variables that we hypothesize to have a causal relationship with depression. In addition to this, we will incorporate more modifiable lifestyle factors, as opposed to ones they focused on such as self-esteem.

## 2 Methods

### 2.1 LASSO Regression for Feature Selection

For feature selection, we used the LASSO Regression method. LASSO is particularly helpful here because it would help us know which independent variables are useless; it shrinks the slope of the regression line to exactly 0, and by doing so it removes less important features.

---

<sup>1</sup>Liu XQ, Guo YX, Zhang WJ, Gao WJ. Influencing factors, prediction and prevention of depression in college students: A literature review. *World J Psychiatry*. 2022 Jul 19;12(7):860-873. doi: 10.5498/wjp.v12.i7.860. PMID: 36051603; PMCID: PMC9331452.

This would help us know what variables are most relevant to depression and enhance the interpretability of our results.

The idea of LASSO is to try to change the slope of the regression line, in order to avoid overfitting. It introduces a bias term by adding the absolute value of the slope:

$$\text{Min}(\text{sum of squared residuals} + \alpha * |\text{slope}|)$$

Where  $\alpha * |\text{slope}|$  is the penalty term. A larger alpha corresponds to a reduced slope, and thus a more horizontal line. The model then becomes less sensitive to the potential variation of an independent variable.

We also eliminated all professions except for "student" in order to make our results more generalizable to students, and certain cities that only had one or two entries. The remaining cities all had 600+ entries. Using a threshold of  $\alpha = 0.05$ , we were able to reduce the number of features from 111 to 13 (including depression).

## 2.2 The Peter-Clark Algorithm (PC)

The Peter-Clark (PC) Algorithm infers a causal graph with the assumption of no latent confounders or variables. Its general framework is as follows:

1. Starts with a fully connected graph
2. Eliminates as many edges as possible using a specified conditional independence test (e.g., Fisher Z)
3. Given a conditioning set  $S$ , the edges between  $X$  and  $Y$  are removed if  $X$  and  $Y$  are independent. This step is repeated starting with an empty set  $S = \emptyset$  and in each iteration the cardinality is increased by 1.
4. The algorithm then establishes causal directions for each remaining edge, using colliders. Some edges may remain undirected if there isn't sufficient evidence to determine a direction.

Finally, the algorithm produces a structured causal graph. Typically the output is at least a partially directed acyclic graph (PDAG). One thing to note here is that there is some randomness in what the final graph looks like, due to the nature of the algorithm itself. The result won't necessarily be the same every time the algorithm is run, but there will be some sort of consistency in the structure.

### 2.2.1 Fisher's Z Conditional Independence Test

One of the options for the Conditional Independence Test when running PC is the Fisher's Z test.

Fisher's Z assumes that the data follows a Gaussian distribution, and that the relationships between variables can be captured using a partial correlation measure. Given  $X$  and  $Y$  (the

variables being tested) and  $Z$  (the conditioning set), the partial correlation between  $X$  and  $Y$  given  $Z$  (denoted as  $\rho_{XY|Z}$ ) is computed as follows:

$$\rho_{XY|Z} = \frac{\rho_{XY} - \rho_{XZ}\rho_{YZ}}{\sqrt{(1 - \rho_{XZ}^2)(1 - \rho_{YZ}^2)}}$$

Where  $\rho$  is Pearson's correlation coefficient. Next, the null and alternative hypotheses are as follows:

$$H_0 : X \text{ and } Y \text{ are conditionally independent given } Z$$

$$H_a : X \text{ and } Y \text{ are conditionally dependent given } Z$$

Using a threshold previously set for  $\alpha$ , the p-value is computed. If the p-value  $> \alpha$ , then we fail to reject the null and conclude that  $X$  and  $Y$  are conditionally independent given  $Z$ . The edge between  $X$  and  $Y$  is then removed. Otherwise, if the p-value  $\leq \alpha$ , then we reject the null hypothesis and the edge is kept.

### 2.2.2 Chi-Squared Conditional Independence Test

The Chi-square test is best used to account for categorical variables in a given dataset, and it is another option when running on PC. Given two categorical variables  $X$  and  $Y$ , again with a conditioning set  $Z$ , a contingency table is constructed and the Chi-Squared test statistic is computed as follows:

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

Where  $O$  is the observed frequency in the contingency table and  $E$  is the expected frequency under the independence assumption. Similar to the Fisher Z framework, the p-value is calculated and if the p-value  $> \alpha$ , then we fail to reject the null hypothesis and the edge between  $X$  and  $Y$  is removed; otherwise, we reject the null hypothesis and keep the edge.

### 2.2.3 Kernel-Based Conditional Independence Test (KCI)

The Kernel-Based Conditional Independence test applies kernel methods to measure dependence. The advantage of the KCI test is that it works for both continuous and discrete data without assuming linearity.

Given variables  $X$  and  $Y$  conditioned on a set  $Z$ , the variables are mapped into a high-dimensional feature space using kernel functions. Then, the Hilbert-Schmidt Independence Criterion (HSIC) is computed and used as a measure of dependence. Then the statistical test is performed as usual, and if the p-value  $> \alpha$ , we fail to reject the null that states that  $X$  and  $Y$  are conditionally independent given  $Z$ , and the edge between them is removed. Otherwise,  $X$  and  $Y$  are assumed to be conditionally dependent given  $Z$  and the edge stays.

This is performed iteratively. An important note regarding this conditional independence test is, because it is a kernel method, its complexity is cubic with respect to the sample size. As such, it is extremely difficult to run on a dataset of 27,900+ entries and requires immense computation power and RAM. The number of samples we ran the algorithm on is further elaborated on in the results section.

## 2.3 Fast Causal Inference (FCI) with KCI (Fast KCI)

The Fast Conditional Independence Algorithm with KCI differs from the PC Algorithm in that it accounts for the presence of latent confounders, while PC assumes no such confounding variables. It is also “faster” than the KCI test used in the PC algorithm, which means it would allow us to implement it with a larger sample size.

Just like the PC, FCI begins with a fully connected graph. Then, it iteratively removes edges using the specified conditional independence test (e.g., KCI). After this, the framework is similar to that in PC; the variables are mapped into a high-dimensional feature space using kernel functions, the HSIC is computed to measure dependence, and the statistical test is performed. If the specified  $\alpha < p$ -value, we fail to reject the null hypothesis and the edge between  $X$  and  $Y$  is removed. Otherwise, the edge between them is kept under the assumption that they are conditionally dependent given  $Z$ . The resulting undirected graph is used to find potential latent confounders, and the output is a Partial Ancestral Graph (PAG) which accounts for such variables.

## 2.4 Counterfactual Computing

Counterfactual Computing is a causal inference method used to infer the result of a potential alternative past. A typical counterfactual question would ask: what would have happened to the outcome of  $Z$ , if I intervened on an observed value  $x$  for the variable  $X$  with a different value? e.g., I observed a score of 1.0 for Depression ( $Z = 1.0$ ), and the Academic Pressure level was ( $X = 5$ ). What would have happened to the Depression score ( $Z$ ) if there was an intervention and the Academic Pressure level was ( $X = 2$ )?

We utilize the method of computing counterfactuals to investigate the potential alternative outcome of Depression when intervening on various variables that are inferred to have a direct causal relationship with Depression, and some that we hypothesized would, but did not make the cut for variable selection.

## 3 Results

### 3.1 Peter-Clark Algorithm Results

Using the 13 features selected by LASSO, we ran the Peter-Clark Algorithm with three independence tests: Fisher’s Z, Chi-Squared, and KCI. When running Fisher’s Z and Chi-Squared, we used all 27,900+ samples. For KCI, 500 samples were used, due to its cubic time complexity.

In the causal graph for Chi-Squared (Figure 1), we observe the variables directly linked with Depression: Academic Pressure, Financial Stress, Unhealthy Dietary Habits, all with a bi-directed arrow. This suggests that while they are causally linked to Depression, there may be a hidden confounder influencing the relationship. An interesting observation here is the city of Patna which has a directed arrow pointing to Depression; this suggests that this specific city may have a direct influence on Depression. The City of Patna also directly points to Suicidal Thoughts. Depression also has an arrow pointing to the city of Ahmedabad; although there is evidence that mental illnesses are widely spread in the city<sup>2</sup>, it is difficult to say why Depression would lead to the city rather than the other way around.

In the graph result for the Fisher Z independence test (Figure 2), we see Depression having a bi-directed relationship with Academic Pressure, Financial Stress (as seen in Chi-Squared) and Suicidal Thoughts. Depression also has two arrows pointing to Unhealthy and Moderate Dietary Habits, suggesting that Depression may lead to worse Dietary Habits. We also observe that Academic Pressure points to Financial Stress, which points to Suicidal Thoughts. This may suggest that having Financial Issues may lead to having Suicidal Thoughts, which could cause Depression. It is difficult to infer why Depression would also point to the cities Ahmedabad (as seen in Chi-Squared) and the city of Hyderabad, though Depression is prevalent in both<sup>3</sup>.

When running PC with the KCI test, we observe in the resulting graph (Figure 3) the bi-directional relationship between Depression and Academic Pressure, Financial Stress and Suicidal Thoughts (as seen in Chi-Squared and Fisher Z), and that Depression points to Unhealthy Dietary Habits (again, like the Fisher Z graph). We note here that Depression also has a bi-directional arrow with “Other” Dietary Habits. This suggests a potential strong relationship between Depression and Dietary Habits, as it is consistently seen in the resulting graphs.

---

<sup>2</sup><https://timesofindia.indiatimes.com/city/ahmedabad/400-increase-in-anxiety-300-in-depression-cases-in-ahmedabad-hospitals/articleshow/114117258.cms>

<sup>3</sup><https://pubmed.ncbi.nlm.nih.gov/35573619/>

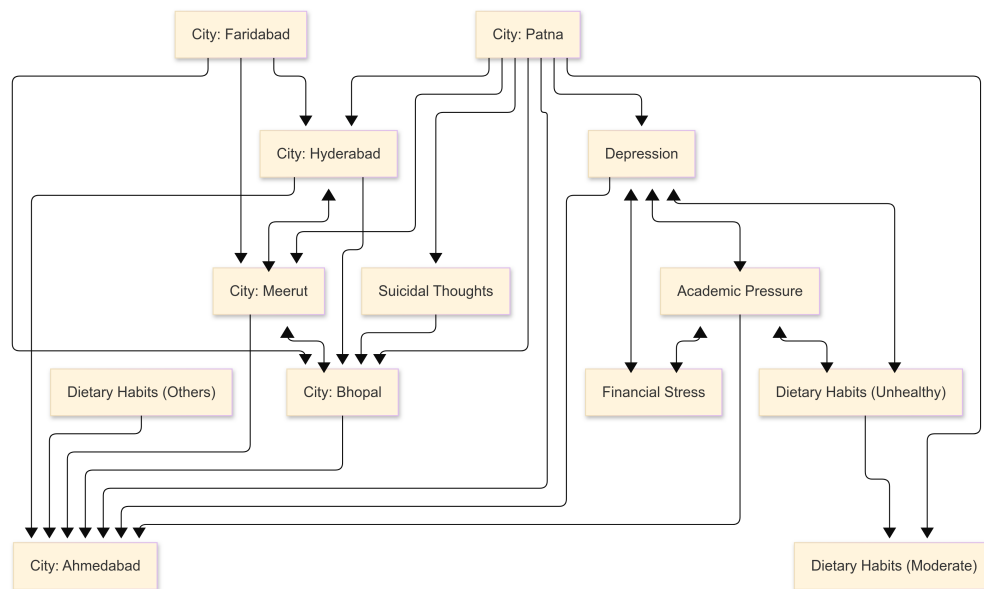


Figure 1: PC with Chi-Squared Graph

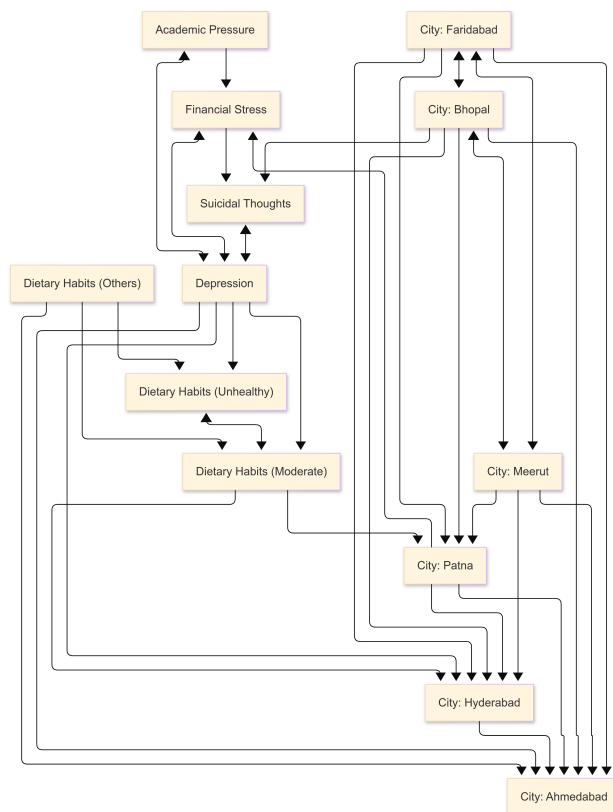


Figure 2: PC with Fisher's Z Graph

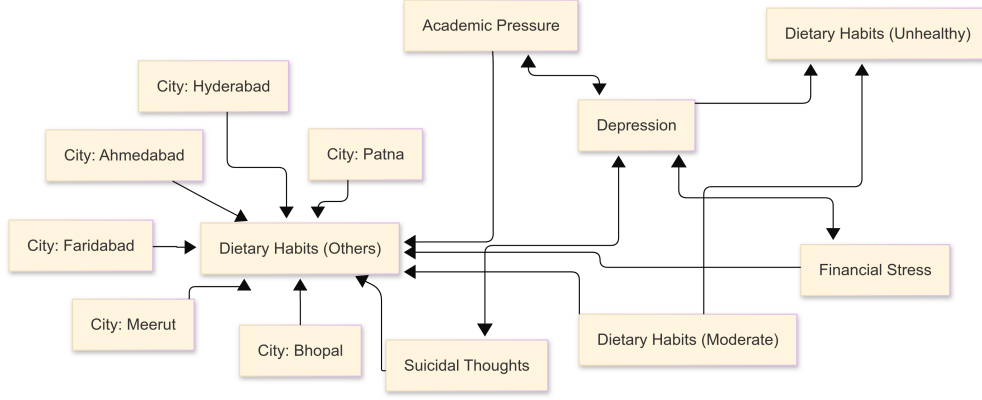


Figure 3: PC with KCI Graph

### 3.2 Fast Causal Inference with KCI Results

The Fast Causal Inference Algorithm applied with the KCI test allows us to use a larger sample from the dataset, as it is not as computationally intensive as the KCI test used with PC. 2000 samples were used, and in the resulting graph from (Figure 4), we see the same variables linked to Depression. There is a bi-directional arrow between Academic Pressure, Financial Stress, Suicidal Thoughts and Unhealthy Dietary Habits. This is mostly consistent with what we've seen in the previous graphs resulting from the three tests we ran with the PC. This suggests the robustness of the features and the algorithms.

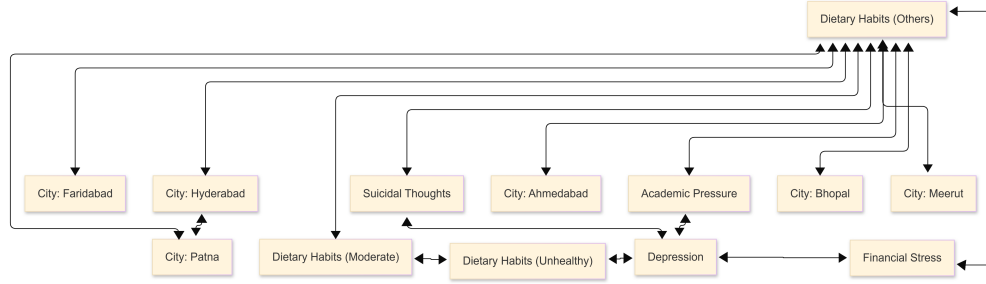


Figure 4: Fast KCI Graph

### 3.3 Counterfactual Computing

To incorporate computing counterfactuals, we intervene on some of the variables that were consistently linked to Depression in our resulting graphs: Academic Pressure, Financial Stress and Suicidal Thoughts. In addition to this, we intervene on Sleep Duration, a variable we hypothesized would play a role in the cause of Depression, but never made the cut for feature selection. This allows us to see how robust the results are from the algorithms, as well as LASSO Regression when eliminating variables which are less useful.

We intervened on Academic Pressure, which originally had an observed value of 5, and fixed the other variables. When intervening with a value of 4 and 1, there is a significant



decrease in the score for Depression, as seen in Figure 5.

Financial Stress seems to have less of a significant impact when compared to Academic Pressure. While controlling for other variables, intervening with the values 2, 4, and 5 makes a smaller difference in the score for Depression (see Figure 6). The original observed value for Financial Stress was 1, but increasing it did not seem to have a big difference in the score for Depression.

Intervening on Suicidal Thoughts (a binary variable) had a similar impact to Academic Pressure. Controlling for other variables again, intervening on Suicidal Thoughts and setting it to 0 significantly decreased the Depression score from the original observed score of 1 (see Figure 7).

We also intervened on Sleep Duration, which was thought to have a role in Depression. This feature did not make the cut when using LASSO Regression for variable selection. The original observed value for Sleep Duration was less than 5 hours, and we intervened with two values: 7-8 Hours, and 8+ Hours. However, the score for Depression did not significantly change with these interventions (as seen in Figure 8). This shows us that LASSO was indeed robust when eliminating this feature.

(intervention) Academic Pressure	Financial Stress	Suicidal Thoughts	Depression
4	1	1	<b>0.88136</b>
1	1	1	<b>0.525441</b>

Figure 5: Depression Score When Intervening on Academic Pressure

Academic Pressure	(intervention) Financial Stress	Suicidal Thoughts	Depression
5	5	1	<b>1.311188</b>
5	4	1	<b>1.233391</b>
5	2	1	<b>1.077797</b>

Figure 6: Depression Score When Intervening on Financial Stress

Academic Pressure	Financial Stress	(intervention) Suicidal Thoughts	Depression
5	1	0	<b>0.579094</b>

Figure 7: Depression Score When Intervening on Suicidal Thoughts

Academic Pressure	Financial Stress	Suicidal Thoughts	(intervention) Sleep Duration	Depression
5	1	1	7-8 Hours	0.966465
5	1	1	8+ Hours	0.925635

Figure 8: Depression Score When Intervening on Sleep Duration

## 4 Discussion

As mentioned, in the causal graphs there are specific variables that we observe to be directly and causally related to Depression; those being Academic Pressure, Financial Stress, Unhealthy Dietary Habits, and Suicidal Thoughts (with the exception of the Chi-Squared graph). Various past studies have shown that academic pressure plays a significant role in various mental illnesses, including Depression. (see [Soh N \(2012\)](#) and [Ye \(2025\)](#)).

Other variables, such as certain cities, were also seen as directly related to Depression. This does not come as a big surprise, as some studies have shown the prevalence of Depression in cities like Ahmedabad ([Solanki SR \(2024\)](#)), and in the state where Hyderabad is located ([Kumar RK \(2022\)](#)). It also does not come as a surprise that Patna points to Depression in a few of the graphs, as studies like [Jha KK \(2017\)](#) have shown that Depression is prevalent among students in the state where Patna is. The study ([Bharati DR \(2022\)](#)) even shows that Depression is more significantly associated with female students, soft drink consumption and fast foods (which indicates a poor diet), and academic dissatisfaction. This aligns with our findings and is consistent with the variables we see as directly related to depression.

## 5 Conclusion

Mental health remains a relatively neglected yet important aspect in our society, especially in the lives of students. Educational institutes should aim to provide valuable resources for their students, as we've seen that students seem to constantly struggle with maintaining good mental health. Between being pressured in school, worrying about finances (especially for college students), and struggling to maintain a healthy diet, students are at an increased risk of having mental health issues. The city in which a population lives in should be considered as well - providing targeted and localized mental health resources would benefit those who live in a city where mental health is more widely spread.

Moving forward, there is more that can be done with our results and findings. As many of the variables had bi-directional relationships, there is potential for confounding variables. As such, some further investigation and analysis of such potential confounders would widen the scope of our results.

## References

- Bharati DR, Prasad N Choudhary SK Kumar S Pal R., Kumari S.** 2022. “Correlates of depression among school going adolescents in the urban area of Patna in eastern India..” *J Family Med Prim Care.* [\[Link\]](#)
- Broomhall AG, Hine DW Loi NM., Phillips WJ.** 2017. “Upward counterfactual thinking and depression: A meta-analysis..” *Clin Psychol Rev.* [\[Link\]](#)
- Howlett JR, Paulus MP.** 2013. “Decision-Making Dysfunctions of Counterfactuals in Depression: Who Might I have Been?” *Front Psychiatry.* [\[Link\]](#)
- Huckins JF, Hedlund EL Murphy EI Rogers C Wang W Obuchi M Holtzheimer PE Wagner DD Campbell AT, DaSilva AW.** 2020. “Causal Factors of Anxiety and Depression in College Students: Longitudinal Ecological Momentary Assessment and Causal Analysis Using PCMCI.” *JMIR Ment Health.* [\[Link\]](#)
- Jha KK, Nirala SK Kumar C Kumar P Aggrawal N., Singh SK.** 2017. “Prevalence of Depression among School-going Adolescents in an Urban Area of Bihar, India..” *Indian J Psychol Med.* [\[Link\]](#)
- Kumar RK, Biradar N Reddy KS Soubhagya M Sushma SA., Aruna G.** 2022. “The prevalence of depression, anxiety, and stress among high school adolescent’s children in public and private schools in Rangareddy district Telangana state: A cross-sectional study..” *J Educ Health Promot.* [\[Link\]](#)
- Soh N, Lampe L Hunt G Malhi G Walter G., Ma C.** 2012. “Depression, financial problems and other reasons for suspending medical studies, and requested support services: findings from a qualitative study..” *Australasian Psychiatry.* [\[Link\]](#)
- Solanki SR, Shukla RP Rathod VG Solanki RN, Dave VR.** 2024. “Assessment of anxiety and depression among tuberculosis patients of Ahmedabad, India..” *Indian J Tuberc.* [\[Link\]](#)
- Ye, Zhang Z. Tao Z. Liping C. Wang Y. Chen H. ... Zhou J., Y.** 2025. “Academic Pressure and Psychological Imbalance in High School Students: Predictors of Depression via Polynomial Regression and Response Surface Analysis..” *Psychology Research and Behavior Management*, 18, 15–23.. [\[Link\]](#)

## 6 Contributions

**Maryam Almahasnah:** Wrote and edited Project Proposal (including dataset search and specification of methods and objectives of the project), implemented FCI with KCI and Counterfactual Computing, literature review and exploration of past approaches. Wrote report, created poster, worked on website.

**Hunter Brownell:** Worked on PC-algorithm (mainly focused on KCI). Used DSMLP resources to run KCI as a background task. Wrote abstract, discussion and conclusion as well as parts of intro, methods and results sections. Wrote the content for the website.

**Kevin Chan:** LASSO regression, PC-algorithm(fisherz, chisq), helped set up instructions and requirements for Github Repository, replication functions, website writing.

## Appendices

A.1 Project Proposal . . . . .	A1
A.2 Dataset Description . . . . .	A1

### A.1 Project Proposal

[Link to our project proposal.](#)

### A.2 Dataset Description

The [dataset](#) includes 27,902 entries, with each row representing an individual person. It contains 18 columns, each one representing a certain feature. The target variable is Depression Status (Binary, Yes/No). Other features include: ID, Age, Gender, City, CGPA (Grade Point Average, scaled out of 10.0), Profession, Academic Pressure, Dietary Habits, Sleep Duration, and more.