# Understanding the Drivers of Obesity: A Predictive Analysis

Brandan Ly, Gerardo Gutierrez

California State Polytechnic University, Pomona

*Abstract*—Obesity is a complex condition with significant social and psychological implications, affecting individuals across all age groups and socioeconomic backgrounds, and posing a growing challenge to both developed and developing nations [1]. This paper explores and analyzes the factors contributing to obesity using logistic regression and correlation analysis on survey data from Latin America. The results identify key behavioral and demographic predictors of obesity, offering insights that may inform targeted public health interventions.

*Index Terms*—Obesity, logistic regression, predictive modeling, public health, Latin America.

## I. Introduction

Although awareness of obesity has increased, the obesity epidemic remains a major concern in the United States and continues to rise globally. Current estimates suggest that around 42% of U.S. adults and 15–20% of children and adolescents are classified as obese, affecting nearly all segments of the population [2]. Obesity significantly increases the risk of chronic health conditions across age groups and is believed to have developed gradually due to small but sustained energy imbalances over time. While public health initiatives have been implemented to curb obesity, there is limited evidence of widespread progress. The complexity of the issue makes it one of the most challenging public health problems of our time. Although some U.S. subgroups have shown signs of stabilization in recent years, global obesity rates continue to rise, highlighting the need for deeper understanding and more effective interventions.

### A. Hypothesis

By analyzing extensive physical and behavioral data of individuals, we hope to find key indicators of obesity. Our goal is to provide key information on the growing obesity epidemic. Globally, there has been an rapid increase in the consumption of overly processed foods and a heavy reliance on motor transportation. Because of this, we believe the poor eating habits are the key drivers of obesity.

## II. Background and Motivation

Obesity is typically defined using body mass index (BMI), where BMI $= \frac{\text{weight (kg)}}{\text{height (m)}^2}$. A BMI of 30 or higher is considered obese. Although widely used, BMI does not distinguish between fat and muscle mass, and may not reflect true health status in all individuals. Nevertheless, it is a well established population health measure to due its practicality compared to other metrics like body fat percentage which are expensive to precisely calculate [3].
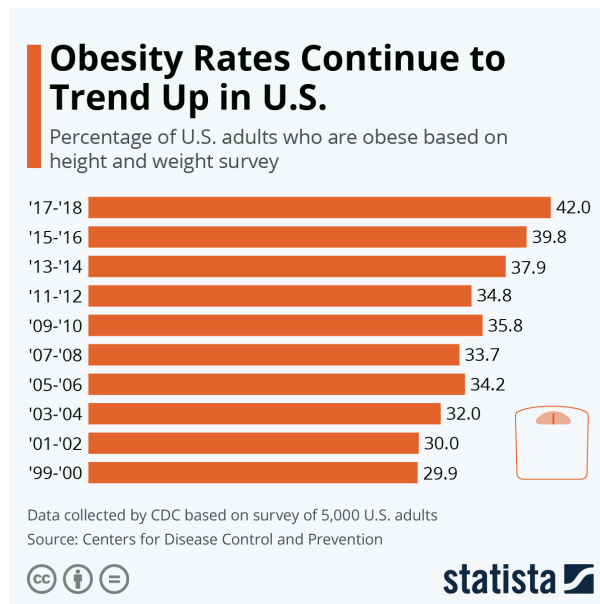


Fig. 1. Obesity prevalence in adults from 1999-2020

Put simply, obesity is a chronic energy imbalance that results from a excessive calorie intake. Beyond calories, obesity is a complex condition with various risk factors and related health risks. So, while we may not be able to change all of our risk factors, it's important to know your risk to work to lower your chances of obesity related risk factors.

Despite increased awareness, obesity rates continue to be a major concern for both developed and developing countries. As shown in Figure 1, the prevalence in obesity in adults has increase by more than 12% since 1999. And in Figure 2 we see that the prevalence of obesity in children has increased by 6% in the same time period. This brings into question the effectiveness of health education and campaigns in the United States. It is important to understand the key risk factors and how to control them due to the numerous chronic conditions associated with obesity, such as coronary artery disease, type 2 diabetes, anxiety, and depression [4].
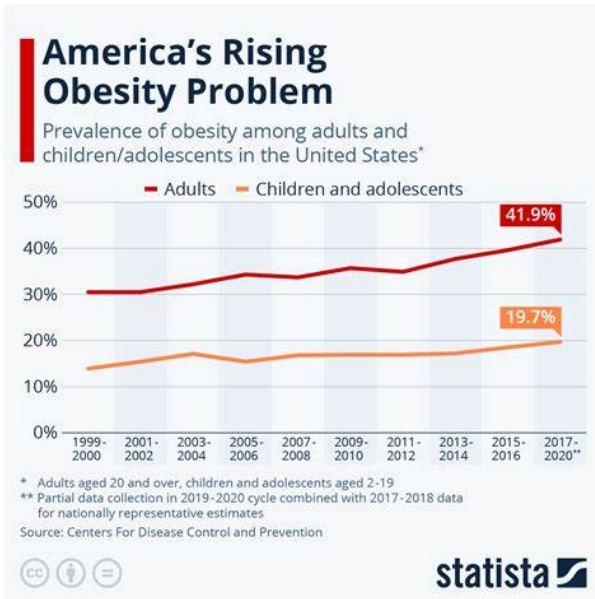
Fig. 2. Obesity prevalence in adults and children from 1999-2020

This paper analyzes data from the UCI Machine Learning Repository [5], which includes behavioral and demographic features such as physical activity, smoking, drinking, nutritional habits, and family history. Our goal is to identify the most influential predictors of obesity to support more targeted health interventions.

### A. Existing Literature

Let's review some existing research that has investigated the risk factors of obesity.

*1) Obesity In Children:* In a cross-sectional study analyzing data from 12 countries. Behavioral variables in 6,025 children aged 9-11 were studied, such as: nocturnal sleep duration, moderate-to-vigorous-physical activity (MVPA), TV time, and an assigned "healthy diet pattern score" [6]. For boys, the odds ratios (95% confidence intervals, CI) for obesity from the multivariable model were 0.79 (0.71–0.90) for nocturnal sleep duration, 0.52 (0.45–0.60) for MVPA, 1.15 (1.05–1.27) for TV time, 1.08 (0.96–1.20) for healthy diet pattern score and 0.93 (0.83–1.04) for unhealthy diet pattern score. For girls, the odds ratios (95% confidence intervals) for obesity from the multivariable model 0.71 (0.63–0.80) for nocturnal sleep duration, 0.43 (0.35–0.53) for MVPA, 1.07 (0.96–1.19) for TV time, 1.05 (0.93–1.19) for healthy diet pattern score and 0.96 (0.82–1.11) for unhealthy diet pattern score. The results of this study indicate that low MVPA, high TV viewing and short sleep duration are significant behavioral correlates of obesity in children.

*2) Obesity In Adults:* In another multinational study across 31 countries analyzing survey data in men and women ages 15-102, behavioral variables were studied such as: smoking status, alcohol consumption, daily fruit and vegetable (FV) intake, physical activity frequency [7]. Light smoking was found to be inversely related to obesity, likely due to appetite suppression,

although not recommended as a preventive measure due to other associated health risks. Regarding alcohol use, drinking 4 or more drinks in a day once or several times a month, was associated with both overweight and obesity. Daily fruit and vegetable consumption was found to be protective against overweight and obesity in men, but no statistically significant association was found in women. Physical activity was positively associated with overweight, but not obesity - likely due to reverse causation.

Overall the most consistently identified risk factors among studies were: low physical activity, poor nutrition, in adequate sleep, and alcohol use.

### B. Data Context and Limitations

The dataset used for analysis contains 17 attributes 2111 records, 23% of which are real and 77% of which were created synthetically using the weka tool and the SMOTE filter. The inital collection

## III. METHODS

This section focused on the methodolgy of our data analysis.

### A. Data Collection

Our dataset is sourced from the University of California, Irvine's Machine Learning Repository [5]. The dataset includes behavioral and demographic characteristics of individuals from the countries of Mexico, Peru, and Colombia. The data was collected through an online survey over the span of 30 days. The dataset contains 17 attributes 2111 records - 77% of which was created synthetically using SMOTE.

### B. Data Wrangling

Most of the preprocessing and cleaning was already done by UCI, hence there was no missing values. Our data wrangling was done in python with the Pandas library. First, we renamed all column names to be intuitive and readable for smoother analysis. Then we proceeded by adding features like BMI, a binary obese feature, placing numerical features like 'water consumed daily' into bins to prepare for categorical analysis such as Spearman's Correlation Test. Additionally, we one hot encoded the categorical features to train a logistic regression model in Sci-Kit Learn.

### C. Exploratory Data Analysis

Our data analysis was driven by creating visualizations using the Seaborn and Matplotlib libraries. Additionally, we used Sci-kit learn to train a logistic regression model that helped us identify the greatest predictors of obesity.

*1) Trends and Relationships:* Correlation heatmaps, bar graphs and stacked bar graphs, were used to identify trends and relationships since most of our data was categorical.

*2) Statistical Modeling:* Additionally, we also calculated and plotted the relative risk of being obese for the binary features: *Family History, Monitor Calories, Eat High Calorie Foods Frequently, Smoker Status.*

## IV. RESULTS

This section presents the results of our analysis and logistic regression models. The primary objective of our analysis was to identify the key drivers of obesity.

### A. Correlation Analysis

From Spearman's correlation heatmap in Figure 3, you can observe that the strongest positive correlations were: *Family History (.50), Age(.38), Eating, High Calorie Meals Frequently(.21).*
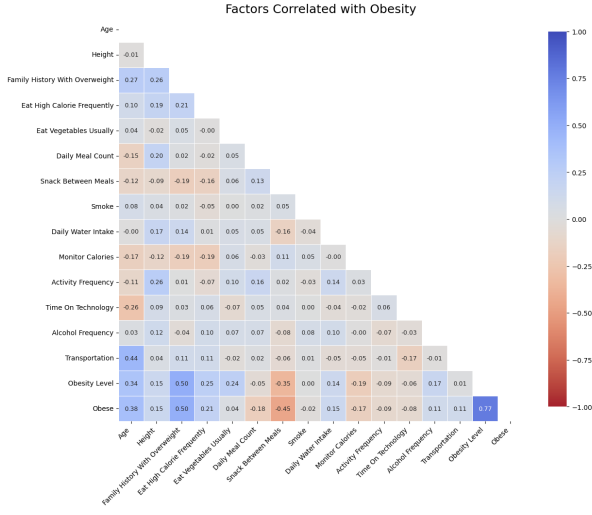


Fig. 3. Spearman correlation heatmap showing relationships among variables.

### B. Stacked Bar Graphs

In Figure 4, Figure 5, and Figure 6 we look at the distribution of the *Family History of Overweight, Eat High Calorie Food Frequently, and Monitor Calories* variables are displayed. We can see that over 90% of individuals who reported no family history of being overweight were not obese, and over 50% of those who reported no family history of overweight obese. 90% of individuals who do not consume high calorie foods were not obese. Over 50% of individuals who consume high calorie foods were obese. Over 90% of individuals who monitor their calorie intake were not obese. Over 40% of individuals who did not monitor calories were obese.

These graphs further highlight the strong link to obesity of *family history of overweight, not monitoring calories, eating high calorie foods frequently*.
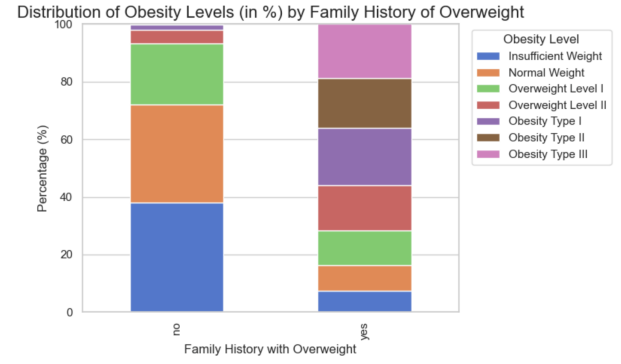


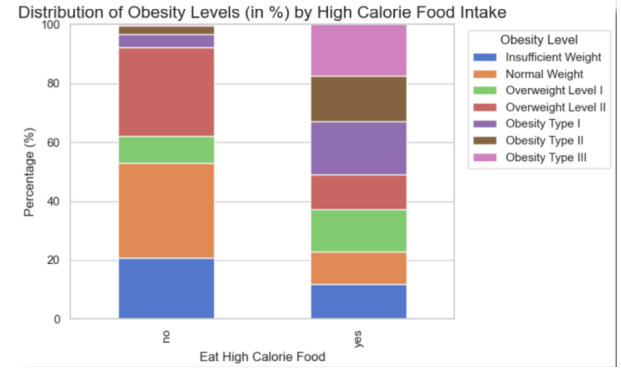Fig. 4. Obesity Level Distribution of Family History Of Overweight



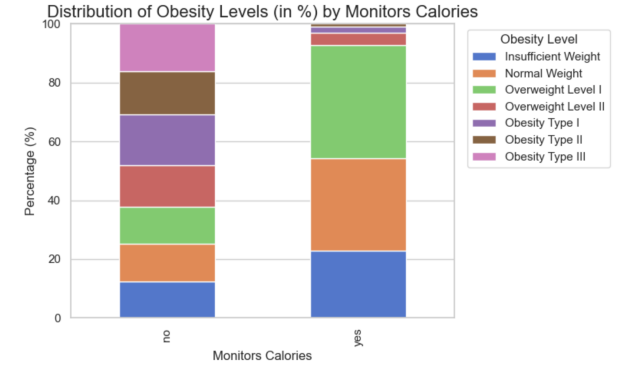Fig. 5. Obesity Level Distribution of Eat High Calories Frequently



Fig. 6. Obesity Level Distribution of Monitor Calories

### C. Relative Risk

In Figure 7 we observe the relative risk of obesity for the binary features of our dataset. For each feature (relative risk of obesity) we have: *Family History of Overweight (3.18), Not Monitoring Calories(1.94), Eating High Calorie Food Frequently (1.63), Being a Male (1.12), Active Smoker(0.93).*

Once again, this emphasizes the strong link having family history or poor nutrition to obesity.

### D. Logistic Regression

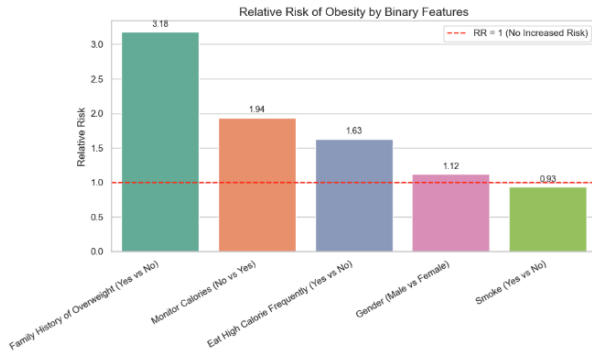And again in our Figure 8 you can see that family history is the strongest risk factor for obesity.

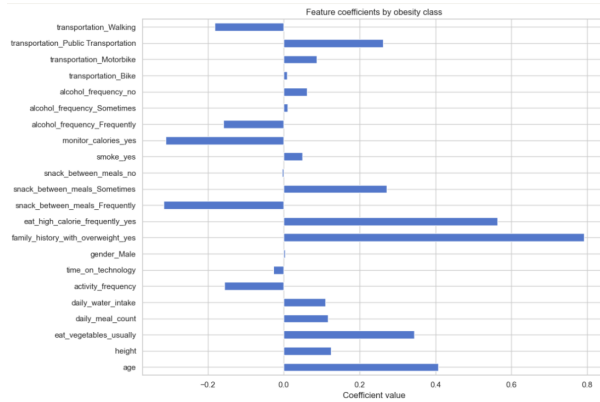Fig. 7. Relative Risk of Obesity for Binary Features



Fig. 8. Logistic Regression Model showing strongest predictors of obesity.

## V. DISCUSSION

Our results show that the two biggest predictors of obesity are *family* history of overweight and *not monitoring your calories*. This highlights that, overall, poor nutritional habits are the overall greatest risk factor of obesity. Since, apart from the genetic implications, those with family history of overweight are likely to adopt the same habits of their parents and family.

### A. Limitations

We must consider that the data was self reported so there could be some reporting bias. Additionally, since 77% of the data was synthetic, the trends and patterns in the graphical summaries may be inflated.

### B. Policy Implications

Using this data, we can look to improve health education and campaigns in Latin America. Many people may not have healthy habits because they have never been made aware of how high calorie food affects their body.

## VI. CONCLUSION

Our hypothesis was proven to be true. Poor nutritional habits are one of the main drivers of obesity. However, it's important to investigate how we can improve public policy and education

to control the non-genetic implications of having a family history of overweight.

### A. Future Work

In the future, we would like to take a more granular approach to investigate why people have poor nutritional habits. We can analyze why people make the food choices and ask questions like "Have you ever been educated on the macronutrients of food and calorie density?"

REFERENCES

[1] World Health Organization, "Controlling the global obesity epidemic," https://www.who.int/activities/controlling-the-global-obesity-epidemic, 2023, [Online].

[2] K. M. Flegal, M. D. Carroll, C. L. Ogden, and L. R. Curtin, "Prevalence and trends in obesity among us adults, 1999–2008," *JAMA*, vol. 303, no. 3, pp. 235–241, 2010.

[3] Centers for Disease Control and Prevention, "About adult bmi," https://www.cdc.gov/bmi/about/index.html, n.d., accessed: 2025-05-07.

[4] M. P. Manoharan, R. Raja, A. Jamil, D. Csendes, S. D. Gutlapalli, K. Prakash, K. M. Swarnakari, M. Bai, D. M. Desai, A. Desai, and S. S. Penumetcha, "Obesity and coronary artery disease: An updated systematic review 2022," *Cureus*, vol. 14, no. 9, p. e29480, Sep. 2022.

[5] C. A. Fernandes, M. H. C. Fernandes, and J. M. C. Rodríguez, "Estimation of obesity levels based on eating habits and physical condition," https://archive.ics.uci.edu/dataset/544/estimation+of+obesity+levels+based+on+eating+habits+and+physical+condition, 2019, uCI Machine Learning Repository. [Online].

[6] P. T. Katzmarzyk, T. V. Barreira, S. T. Broyles, C. M. Champagne, J.-P. Chaput, M. Fogelholm, G. Hu, W. D. Johnson, R. Kuriyan, A. Kurpad, E. V. Lambert, C. Maher, J. Maia, V. Matsudo, T. Olds, V. Onywera, O. L. Sarmiento, M. Standage, M. S. Tremblay, C. Tudor-Locke, P. Zhao, T. S. Church, and for the ISCOLE Research Group, "Relationship between lifestyle behaviors and obesity in children ages 9–11: Results from a 12-country study," *Obesity*, vol. 23, no. 8, pp. 1696–1702, 2015. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/oby.21152

[7] S. Pengpid and K. Peltzer, "Associations between behavioural risk factors and overweight and obesity among adults in population-based samples from 31 countries," *Obesity Research Clinical Practice*, vol. 11, no. 2, pp. 158–166, 2017. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1871403X16300722