# Variables Affecting Obesity in Colombian, Mexican, and Peruvian Populations

### 1. Introduction

Obesity is a growing problem worldwide, and it is important to examine the obesity levels of people to determine how to best remedy this global issue. The World Health Organization states that the fundamental cause of obesity is an energy imbalance between calories consumed and calories expended. In this study, we will examine how people's eating and behavioral habits will affect their obesity levels as determined by mass body index. Our study focuses on data collected via a digital survey from Colombia, Peru, and Mexico. For this study, we chose to focus on the following classification problem: determining if a person will be obese or not given a set of lifestyles and physical health conditions.

### 2. About The Data

Our obesity-related data was collected from populations with ages between 14 and 61 and diverse eating habits and physics conditions in Colombia, Peru, and Mexico. The research team who collected the data used an anonymous web-based survey that was available for 30 days with unbiased questions. At the end of the surveying period, 485 records were received. The survey questions used are shown in Figure 1.

Questions	Possible Answers
¿What is your gender?	Female
	<ul> <li>Male</li> </ul>
¿what is your age?	Numeric value
¿what is your height?	Numeric value in meters
¿what is your weight?	Numeric value in kilogram
¿Has a family member suffered or suffers from overweight?	• Yes
	<ul> <li>No</li> </ul>
¿Do you eat high caloric food frequently?	<ul> <li>Yes</li> </ul>
	<ul> <li>No</li> </ul>
¿Do you usually eat vegetables in your meals?	<ul> <li>Never</li> </ul>
	<ul> <li>Sometimes</li> </ul>
	<ul> <li>Always</li> </ul>
¿How many main meals do you have daily?	<ul> <li>Between 1 y 2</li> </ul>
	<ul> <li>Three</li> </ul>
	<ul> <li>More than three</li> </ul>
¿Do you eat any food between meals?	<ul> <li>No</li> </ul>
	<ul> <li>Sometimes</li> </ul>
	<ul> <li>Frequently</li> </ul>
	<ul> <li>Always</li> </ul>
¿Do you smoke?	Yes
	• No
¿How much water do you drink daily?	<ul> <li>Less than a liter</li> </ul>
	<ul> <li>Between 1 and 2 L</li> </ul>
	<ul> <li>More than 2 L</li> </ul>
Do you monitor the calories you eat daily?	<ul> <li>Yes</li> </ul>
	• No
¿How often do you have physical activity?	<ul> <li>I do not have</li> </ul>
	<ul> <li>1 or 2 days</li> </ul>
	<ul> <li>2 or 4 days</li> </ul>
	<ul> <li>4 or 5 days</li> </ul>
¿How much time do you use technological devices such as	<ul> <li>0–2 hours</li> </ul>
cell phone, videogames, television, computer and others?	<ul> <li>3–5 hours</li> </ul>
	<ul> <li>More than 5 hours</li> </ul>
; how often do you drink alcohol?	<ul> <li>I do not drink</li> </ul>
	<ul> <li>Sometimes</li> </ul>
	<ul> <li>Frequently</li> </ul>
	Always
Which transportation do you usually use?	Automobile
	<ul> <li>Motorbike</li> </ul>
	Bike
	Public Transportation
	Walking

Once data are collected and labeled into different levels of obesity, the data collectors identified that the distribution of data was imbalanced across categories, as shown in Figure 2. This imbalance would pose difficulties for our classification task, as it might result in classifiers 'with a high accuracy but very low sensitivity towards the positive class' (Elhassan et al. 2017). For this reason, synthetic data was generated using Weka and SMOTE (synthetic minority over-sampling technique). The balanced distribution of the synthetic and original data combination is shown in Figure 3.

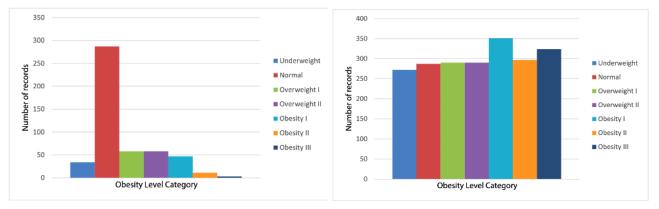


Figure 2: Imbalance in categories of data shown in a bar chart.

Figure 3: Balanced data categories after synthetic data generation.

The resulting dataset required cleaning, as the generated synthetic data was very noisy and did not follow the specifications as listed in the original survey. For instance, while our multiple-choice answers were recorded using whole numbers, the synthetic data contained values with decimals. As a result, we needed to round the values generated from the synthetic process. Additionally, we added a BMI column, calculated as weight/(height)<sup>2</sup>.

	Gender	Age	Height	Weight	family_history_with_overweight	FAVC	FCVC	NCP	CAEC	SMOKE	CH2O	scc	FAF	TUE	CALC	MTRAN
0	Female	21.0	1.62	64.0	yes	no	2.0	3.0	Sometimes	no	2.0	no	0.0	1.0	no	Public_Transportation
1	Female	21.0	1.52	56.0	yes	no	3.0	3.0	Sometimes	yes	3.0	yes	3.0	0.0	Sometimes	Public_Transportation
2	Male	23.0	1.80	77.0	yes	no	2.0	3.0	Sometimes	no	2.0	no	2.0	1.0	Frequently	Public_Transportation
3	Male	27.0	1.80	87.0	no	no	3.0	3.0	Sometimes	no	2.0	no	2.0	0.0	Frequently	Walkir
4	Male	22.0	1.78	89.8	no	no	2.0	1.0	Sometimes	no	2.0	no	0.0	0.0	Sometimes	Public_Transportation

Figure 4: Snapshot of the final dataset.

There are 10 numerical variables in the dataset.

For variable Age, the minimum value is 14.00 and the maximum value is 61.00. For variable Height, the minimum value is 1.450 and the maximum value is 1.980. For variable Weight, the minimum value is 39.00 and the maximum value is 173.00. For variable FCVC, the minimum value is 1.000 and the maximum value is 3.000. For variable NCP, the minimum value is 1.000 and the maximum value is 4.000. For variable CH20, the minimum value is 1.000 and the maximum value is 3.000. For variable FAF, the minimum value is 0.0000 and the maximum value is 3.0000. For variable TUE, the minimum value is 0.0000 and the maximum value is 2.0000. For variable BMI, the minimum value is 13.00 and the maximum value is 50.81. There are 9 categorical variables in the dataset.

For variable Gender, the categories are 'Female' and 'Male'.

For variable Family History with Overweight, the categories are 'yes' and 'no'.

For variable FAVC, the categories are 'yes' and 'no'.

For variable CAEC, the categories are 'no', 'Sometimes', 'Frequently', and 'Always'.

For variable SMOKE, the categories are 'yes' and 'no'.

For variable SCC, the categories are 'yes' and 'no'.

For variable CALC, the categories are 'no', 'Sometimes', 'Frequently', and 'Always'.

For variable MTRANS, the categories are 'Public\_Transportation', 'Walking', 'Automobile', and 'Motorbike'.

For variable NObeyesdad, the categories are 'Insufficient Weight', 'Normal Weight', 'Overweight Level I', 'Obesity Type II', 'Obesity Type II', 'Obesity Type III'.

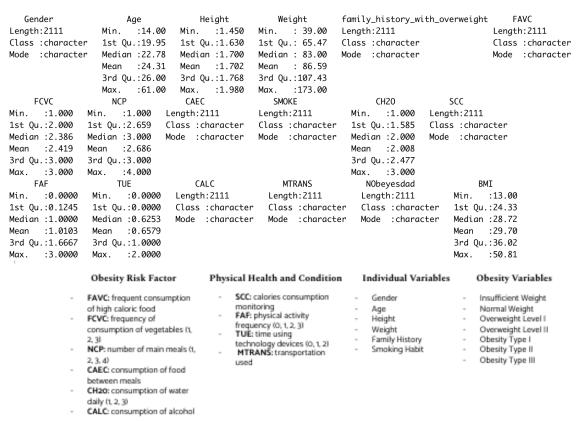


Figure 5: Variables documented in the dataset.

# 3. Exploratory Analysis

## 3.1. Research questions

For our exploratory analysis, we determine to answer the following questions:

- 1. Which factors are most correlated to BMI?
- 2. Can the explanatory variables that are most correlated to BMI be used to classify a person's obesity level?
- 3. How well is our classifier in terms of predictive capability?

## 3.2. Exploratory Results

From the heat map, it appears that there is a correlation between BMI and FAF (Frequency of Physical Activity). For the categorical variables (Figure 7,8), it appears that Frequency of Vegetable consumption and Family History have a greater possible correlation with BMI than the other variables. A separate visualization was created for comparing Family History and BMI as well, as seen in Figure 9. Subsequently, linear regressions were conducted on each of the three variables of interest with their relation to BMI, as well as a regression with all three variables and BMI (Figure 10&11).

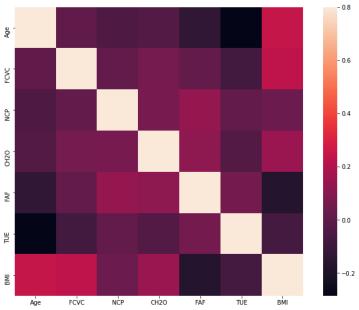
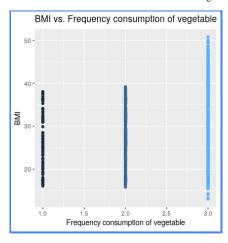
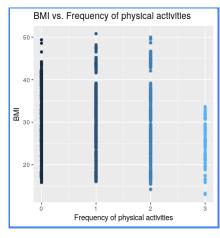
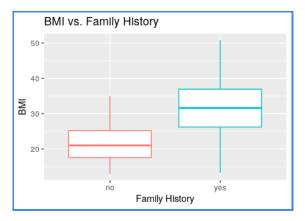


Figure 6: Heat map shows the correlation between numeric variables and BMI

Figure 7&8: Frequency of Vegetable Consumption vs BMI and Frequency of Physical Activities vs BMI Figure 9: Boxplot of Family History vs BMI







BMI vs Family History								BMI vs FCVC								BMI vs FAF							
Dep. Variable:		BMI		R-squa	red:	0.234				Dep. Variable	e:	BA	MI	R-sq	uared:	0.061	Dep. Variat	ble:		BMI	R-s	quared:	0.032
Model:		OLS	Adj.	R-squa	red:	0.233				Mode	el:	OL	S Ac	j. R-sq	uared:	0.060	Mod	delt		OLS	Adj. R-s	quared:	0.031
Method:	Le	east Squares		F-statis	stic:	643.5				Metho	d: Le	ast Square	es	F-sta	itistic:	136.0	Meth	od: I	Least Squ	Jares	F-s	tatistic:	23.45
Date:			Prob (	F-statis	tio):	4.03e-124				Dat	e: Tue, 3	30 Nov 202	21 Prof	(F-sta	tistic):	1.72e-30	Da	ste: Tue	, 30 Nov 2	2021 Pr	rob (F-st	atistic):	6.32e-15
Time:	,	20:47:46		Likelih		-7106.5				Tim		20:47:4		g-Likel		-7321.6	Tir	ne:	20:4		Log-Like	elihood:	-7352.9
			Log-									211		y-Like		1.465e+04	No. Observatio			2111		740.	1.471e+04
No. Observations:		2111				1.422e+04				No. Observation	-				AIC:		Df Residu	als:	2	2107		BIC:	1.474e+04
Df Residuals:		2109		-	BIC:	1.423e+04				Df Residual	S:	210	09		BIC:	1.466e+04	Df Mod	del:		3			
Df Model:		1								Df Mode	el:		1				Covariance Ty	pe:	nonro	bust			
Covariance Type:		nonrobust								Covariance Typ	e:	nonrobu	ist					coef	std err	,	t P>Itl	[0.025	0.975]
				coef	std er	r t	P> t	[0.025	0.9753		coef	-14		P>ltl	[0.025	0.0757	Intercept	31.0088	0.294	105.504	0.000	30.432	31.585
				1.5005				20.799	22.202			std err	t		•		C(FAF)[T.1.0]	-1.0723	0.408	-2.628	0.009	-1.873	-0.272
		Interce	,		0.35		0.000				24.2604		48.892	0.000	23.287	25.233	C(FAF)[T.2.0]	-2.5024	0.460	-5.438	0.000	-3.405	-1.600
C(family_history_v	with_ov	erweight)[T.ye	s] 10	0.0287	0.39	5 25.367	0.000	9.253	10.804	C(FAVC)[T.yes]	6.1540	0.528	11.660	0.000	5.119	7.189	C(FAF)[T.3.0]	-5.7915	0.780	-7.421	0.000	-7.322	-4.261
Omnibus:	27.897	Durbin-Wa	tson:	0.4	177					Omnibus:	124.261	Durbi	in-Watso	in:	0.263		Omnibus	. 115.00	14 Pu	ırbin-Wat		0.224	
Prob(Omnibus):	0.000	Jarque-Bera	(JB):	17.0	120					Prob(Omnibus):	0.000	Jamue	-Bera (JI	Bh 4	7.169		Prob(Omnibus)			ue-Bera		51.159	
Skew:	0.006	Prot	(JB):	0.0002	101					Skew:	0.050	-unque	Prob(JI	-,-	2e-11		Skew	,			(JB): 7.		
Kurtosis:	2.560	Conc	I. No.	4	48												Kurtosis			Cond		5.28	
- Lan Loons.	2.500	Con		-						Kurtosis:	2.275		Cond. N	io.	5.71		Kurtosis	2.00	-	Junu	1400	0.20	

Figure 10: Regression results from BMI vs Family History, BMI vs FCVC, and BMI vs FAF.

Dep. Variable:	BMI	R-s	quared:	0.2	75		
Model:	OLS	Adj. R-s	quared:	0.2	73		
Method:	Least Squares	F-s	tatistic:	159	9.4		
Date:	Wed, 01 Dec 2021	Prob (F-st	tatistic):	5.85e-1	44		
Time:	11:17:12	Log-Lik	elihood:	-704	3.6		
No. Observations:	2111		AIC:	1.411e+	04		
Df Residuals:	2105		BIC:	1.414e+	04		
Df Model:	5						
Covariance Type:	nonrobust						
		coef	std err	t	P> t	[0.025	0.975
	Intercept	20.1732	0.551	36.628	0.000	19.093	21.253
C(family_history_w	ith_overweight)[T.1]	9.2656	0.394	23.515	0.000	8.493	10.038
	C(FAVC)[T.1]	3.4552	0.477	7.243	0.000	2.520	4.391
	C(FAF)[T.1.0]	-1.0223	0.354	-2.891	0.004	-1.716	-0.329
	C(FAF)[T.2.0]	-2.0528	0.399	-5.141	0.000	-2.836	-1.270
	C(FAF)[T.3.0]	-4.3431	0.680	-6.384	0.000	-5.677	-3.009

Figure 11: Regression results from BMI vs Family History, FCVC, and FAF

## 4. Modeling

We chose to perform classification tasks to build the model as all of our attributes are categorical and the labels are provided in the data. Overall, we tested five different classifiers: KNN, SVM, Random Forest, Decision Tree, and Naive Bayes. For kNN classification, we need to specify the number of the nearest neighbors k. To find the best-fitted k, we calculate the error rate at each different k given the trained datasets. The results showed that k=3 has the lowest error rate. For each classification, we split the data into two parts: training data and test data. Training data were used to build the model, and test data were used to evaluate the model.

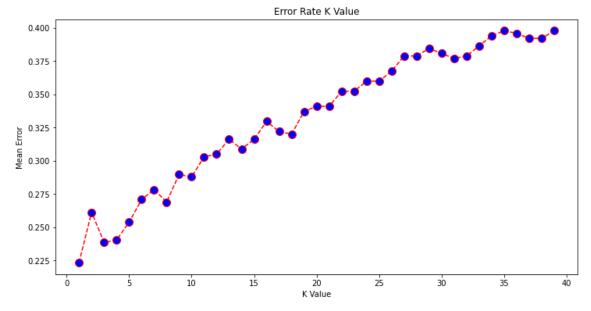


Figure 12: Error rates of various K values in KNN model

### 5. Discussion

For the obesity level, we started our exploratory analysis by investigating the distribution of obesity level by gender: maximum female respondents fell under Obesity Type III and maximum male respondents fell under Obesity Type II. Next, we created a heatmap, in which we identified that there are strong correlations between BMI and Frequency of Physical Activity. We identified relatively strong relations between (1) BMI and Family History, (2) BMI and Frequency of Consumption of Vegetables, (3) BMI and Frequency of Physical Activities, (4) BMI and Alcohol Consumption, (5) BMI and Frequency of Consumption of high-calorie food. In the regression results from BMI vs Family History, BMI vs FCVC, and BMI vs FAF, Family History has the highest coefficient (9.266) compared to other variables, meaning that the change in one unit of Family History will have greater influence on BMI compared to the other variables. At the same time, the adjusted R-squared of the final regression model also points out that only 27.5% of the variation in BMI can be explained by our explanatory variables. We do not, however, have enough evidence to state that any of these factors directly affect the BMI.

The results of our regression further indicate that Family History, FCVC, and FAF alone are not sufficient for the predictive task. Therefore, in choosing the explanatory variables to train our data for classification, we took into account all the possible factors. We calculated the accuracy scores for each classification model and found the KNN model has the highest accuracy score compared to the other 4 models. As a result, we choose KNN as the most fitted classifier for our study. The accuracy score of 76% implies that given a person's set of lifestyle habits (ones that accounted as our explanatory variables), the KNN model will be 76% correct in predicting the person's level of obesity.

The first limitation of our study lies within our selection of the KNN classifier. To be more specific, KNN algorithm is computationally expensive, as we need to calculate the proximity measures during the process. The KNN model is also not good at identifying edge cases, which can adversely be affected by outliers. Additionally, feature selection is critical in KNN, as irrelevant features can dominate a decision, so the model will not perform as well if there is a class imbalance. The second limitation is with our process of data collection. To begin with, 77% of the data is synthetic and consequently might not be representative of the original populations of interest. At the same time, the original survey was only available online and for 30 days, which could cause response and volunteer bias in the data. As a consequence, we lack randomization in the process of collecting data and our findings also lack the external validity as it might not be representative for the population.

## Conclusion

Given the results from our models, KNN appears to be the best fit for our data since it has the highest accuracy. Using our KNN model, we can predict a person's obesity level based on their lifestyle attributes, age, gender, and family history with obesity. To further advance this study, we can try to improve the model's accuracy and sensitivity by using the attributes that are most relevant to the obesity levels instead of using all attributes in its development. Furthermore, to increase the external validity of our study, we can collect more human data and expand our sample size instead of using generated synthetic data.

## Acknowledgment

- Vivian Bui: Actively attended meetings, provided useful insights for the project, implemented R and Python for data visualization and classification modeling, presented project to the class, contributed to the report
- Xinying Qian: Actively attended meetings, provided useful insights for the project, implemented R and Python for data visualization and classification modeling, contributed to the report
- Aishwarya Kotha: Actively attended meetings, provided useful insights for the project, worked hard with the presentation slides, presented project to the class, contributed a lot efforts to the report
- Nikita Mathur: Actively attended meetings, provided useful insights for the project, worked hard with the presentation slides, presented project to the class, contributed to the report
- Yijie Lian: Contributed to the report

### Reference

Elhassan et al. (2017), Classification of Imbalance Data using Tomek Link (T-Link) Combined with Random Under-sampling (RUS) as a Data Reduction Method. Global Journal of G Technology & Optimization.

Fabio Mendoza Palechor, Alexis de la Hoz Manotas, Dataset for estimation of obesity levels based on eating habits and physical condition in individuals from Colombia, Peru and Mexico, Data in Brief, Volume 25, 2019, 104344, ISSN 2352-3409, https://doi.org/10.1016/j.dib.2019.104344.