

Obesity in America: A Growing Trend



David Todd

Utilizing Geographic Information Systems (GIS) to explore obesity in America, this study aims to determine statistically significant correlations for the growing trend.

Pennsylvania State University

GEOG 586: Geographic Information Analysis

Winter 2011

Dr. Justine Blanford

3 / 1 5 / 2 0 1 2

I. Introduction

Obesity, according to Merriam-Webster, is a condition characterized by the excessive accumulation and storage of fat in the body. In general terms, obese is the label for a range of weight much greater than what is generally considered healthy for a given height. In adults, obesity is determined by using weight and height to calculate a number called the "body mass index" (BMI). BMI is used because, for most people, it correlates with their amount of body fat. An adult who has a BMI of 30 or higher is considered obese. Obesity has been shown to increase the likelihood of certain diseases and other health problems, e.g. heart disease, diabetes, and hypertension.

While healthy eating and regular exercise is generally considered the best way to combat obesity, there may be other significant correlating factors that can be linked to obesity, e.g. income, education level, or gender. What's interesting about obesity is the rate at which it has risen to epidemic proportions in the United States with no clear cut reason why. Significant correlating factors could help focus on ways to combat obesity across America or which demographics to target with education campaigns.

II. Objectives

This study aims to utilize Geographic Information Systems (GIS) to explore obesity in America by mapping obesity rates in order to visualize any trends and whether or not there is clustering. The objective is to determine statistically significant correlations for the growing trend of obesity in the United States. It will attempt to dissect any patterns that emerge through the 2010 data timeframe and measure spatial autocorrelation. This study looks to answer the question – *Does obesity have statistically significant correlations?*

III. Data

A. Data Sources

The majority of data gathered for this study comes from the Centers for Disease Control and Prevention (CDC). Specifically, the CDC operates the Behavioral Risk Factor Surveillance System (BRFSS), which is the world's largest, on-going telephone health survey system, tracking health conditions and risk behaviors in the United States yearly since 1984. BRFSS data contains over 100 different variables by state, e.g. obesity rates, income levels, education levels, diabetes, alcohol ingestion, race, gender, etc. To aid with visualization, geographic data sources are also provided by ESRI's base map layers.

B. Data Preparation

The CDC has completed the majority of data preparation. They present the BRFSS data ready to be used for analysis in the form of GIS shapefile maps. The files contain the survey data and documentation, and are available in Zip Archive File (ZIP) format. Once acquired, the data was loaded into ESRI's ArcGIS software for additional preparation. The decision was made to focus on the continental United States (CONUS). The shapefiles were edited to delete the both Alaska and Hawaii, leaving 49 records – 48 States plus the District of Columbia. One additional step for data preparation was required before analysis could begin. Some of the files contained attribute cells that were stored as "Text" data type. Because the data in the fields are numeric in nature, the data type needed to be changed in order to support analysis. To turn the data into a numeric data type, the *.dbf file was opened in Excel so it would be possible to change the format of the Text cells to Number. The file was then saved as an *.xls file. Finally, in ArcMap, the original data layer is joined with the *.xls file, making it possible to perform analysis.

C. Data Analysis

The first step in data analysis for this study is performing a Hot Spot Analysis on obesity, as seen in Figure 1. The output of this tool is a map of z-scores which shows statistically hot spot and cold spot clustering in the data.

📴 Hot Spot Analysis with Rendering		
Input Feature Class		^
brfss_state_2010_downloadCopy	- 🖻	
Input Field		
A4409_3	-	
Output Layer File		
C:\Users\dtodd\Documents\School\586\Final_Project\2010\brfss_state_2010_HSpot1.lyr	6	
Output Feature Class	_	
C:\Users\dtodd\Documents\School\586\Final_Project\2010\brfss_state_2010_HSpot1.shp		
Distance Band or Threshold Distance (optional)		
		_
		-
OK Cancel Environments	Show Help >>	

Figure 1: Hot Spot Analysis (Getis-Ord Gi*) tool

In the next steps, data analysis for this study focused primarily on performing regression analysis using the spatial statistics tools in ESRI's ArcGIS software package. Regression analysis is used to attempt to explain phenomena, i.e. obesity, in terms of other variables, e.g. education, income, gender. Utilizing a GIS like ArcGIS, it is possible to perform regression analysis to model and predict these complex phenomena. Regression analysis starts by determining which variables will specify a good regression model.

Finding a good regression model when over 100 variables are at hand is an extremely iterative process. To aid with this time consuming process, the Exploratory Regression tool (Figure 2), provided by ESRI in their supplemental spatial statistics toolbox was used to determine which explanatory variables were the most statistically significant in explaining the dependent variable obesity. This single process took well over 24 hours to determine models of 1 to 5 statistically significant factors. While 24 hours may seem like a long time to process, it would have taken days or weeks to manually run the thousands of possible model variations.

S Exploratory Regression		X	
Input Feature Class			-
brfss state 2010 downloadCopy	-	2	
Dependent Variable			
A4409 3		•	
Independent Variables			
ST FIPS			
vear			
✓ A187_1			
✓ A187_2			
✓ A187_3			
✓ A187_4			
✓ A187_5			
✓ A187_6			
A559 1		-	
·	- F		
	Add Eigld		
Select All Unselect All	Add Field		
Input Spatial Weights Matrix			
C: \Users\dtodd\Documents\School\586\Final_Project\2010\Regression_Testing\2\brfss2010_S\	VM.swn		
Output Report File			
C:\Users\dtodd\Documents\School\586\Final_Project\2010\Regression_Testing\2\Exp_Reg.txt		6	
Output Table Workspace (optional)			_
C:\Users\dtodd\Documents\School\586\Final_Project\2010\Regression_Testing\2		2	=
Tables to Create			
MIN_MAX_NUMBER_OF_EXPLANATORY_VARIABLES_ONLY		-	
A Fearch Criteria			
Max Number of Explanatory Variables			
	1 I I	1	
3		-	
1		20	
Min Number of Evolapatory Variables			
1		_	
1		20	
		20	
Min Adj. R-Squared			
		0.5	
Max Coefficient p-value		05	
	0	.05	
Max VIF Value		7.5	
		7.5	
Imin Jarque-Bera p-Value		0.1	
Min Contint Autocompletion of unbur		0.1	
Min Spatial Autocorrelation p-value		0.1	
		0.1	Ŧ
OK Cancel Environments	Show He	- - - 	1
	CHOWTIC		
	-	-	

Figure 2: Exploratory Regression Tool in ArcGIS 10.0

Next, as seen in Figure 3, the Scatterplot Matrix tool was used in order to check the possible explanatory variables that were determined by the Exploratory Regression tool. The scatterplots showed whether there was a positive or negative correlation between obesity and the explanatory variables.



Figure 3: Scatterplot Matrix Tool

While the scatterplot matrix assists in visualizing correlations, it's hard to determine just by eye if the variable is statistically significant, which is another reason the Exploratory Regression tool was used first. The tool allowed the study to immediately focus on the significant variables.

The next step in analysis is to perform Ordinary Least Squares (OLS). The OLS tool (Figure 4) performs linear regression to predictions or to model a dependent variable in terms of its relationships to the explanatory variables. It generates different outputs that include a map of regression residuals and a mostly numeric summary report. The map shows the over and under predictions from your model.

I Ordinary Least Squares	- O X
Input Feature Class	*
brfss_state_2010_downloadCopy	- 🖻
Unique ID Field	
ST_FIPS	•
Output Feature Class	
C:\Users\dtodd\Documents\School\586\Final_Project\2010\Regression_Testing\2\OLS_Pos_4shp	
Dependent Variable	
A4409_3	-
vear	
A187_1	
A187_2	
A187_3	
A187_4	
A187_5	
A187_6	-
A559 1	-
Select All Unselect All	Add Field
* Output Options	
Coefficient Output Table (optional)	
(Soo (Final_Project/2010)(Kegression_Lesting)(2)(COEF_Pos_4	
Diagnostic Output Table (optional)	
\586 \Final_Project\2010 \Kegression_I esting \2\DIAG_Pos_4	- E
OK Cancel Environments	Show Help >>

Figure 4: OLS Tool

The other output of OLS is the summary table. It contains a great deal of important information about the explanatory variables in the OLS model. While the summary report has a lot of numeric data, refer to Figure 5, it displays an asterisk next to any variable that is statistically significant. The summary report includes both a probability and robust probability for each variable. The robust probability needs to be used if the spatial relationships vary across the study area. Again, it's easy to determine which probability to use since the summary also includes a Koenker p-value and an asterisk if it's statistically significant. If there is an asterisk, the robust probability should be used. The most notable value in the summary output is the R² value. R² is a value between 0

and 1 that represents the percentage to which the explanatory variables tell the whole story. The closer to 1 the value is, the better.

OLS_Pos_Results.txt - Notepad					
Eile Edit Format View Help					
Summary of OLS Results Variable Coefficient StdError t-Statistic Probability Robust_SE Robust_t Robust_Pr VIF [1] Intercept -35.031516 37.534376 -0.93318 0.356707 47.761995 -0.733460 0.467899 NIF A187_3 0.221651 0.080976 2.737237 0.009468* 0.074621 2.970342 0.005200* 1.945059 A559_2 0.008475 0.037661 0.225025 0.823199 0.026262 0.322694 0.748744 3.738563 A745_1 0.015039 0.091898 0.163645 0.870903 0.094557 0.159042 0.874502 3.545441 A1363 1 1.030951 0.302372 3.409549 0.001385* 0.268908 3.833848 0.0004735 5.858290 A1737_3 0.117909 0.116249 1.014281 0.317036 0.096090 1.227074 0.227545 3.040253 A2976_4 0.722269 0.207719 3.477145 0.001313* 0.167967 4.300055 0.000119* 4.230214 A3348_2 0.685142 0.4557112 1.498847 0.142397 0.591657 1.158005 0.254283 3.733328 A4347_2 0.530521 0.132817 3.994367 0.000296* 0.120821 4.390951 0.000090* 6.459444 A4396_1 0.018740 0.117556 0.155941 0.874211 0.107846 0.173768 0.862996 3.984866 A4413_2 -0.280501 0.265228 -1.05783 0.297102 0.228236 -1.22893 0.226834 1.828992 A6607_2 0.234164 0.077009 3.040741 0.004320* 0.062452 3.749493 0.00065* 5.714505	*				
OLS DiagnosticsInput Features:brfss_state_2010_downloadcopyDependent Variable:Number of observations:4Akaike's Information Criterion (AICC) [2]:Multiple R-squared [2]:0.852442Adjusted R-squared [2]:Joint F-statistic [3]:19.431657Prob(>F), (11,37) degrees of freedom:Joint Wald Statistic [4]:446.324437Prob(>chi-squared), (11) degrees of freedom:Jarque-Bera statistic [6]:0.294614Prob(>chi-squared), (2) degrees of freedom:	A4409_3 196.520320 0.808573 0.000000* 0.000000* 0.028292* 0.863029				
Notes on Interpretation * Statistically significant at the 0.05 level. [1] Large VIF (> 7.5, for example) indicates explanatory variable redundancy. [2] Measure of model fit/performance. [3] significant p-value indicates overall model significance. [4] Significant p-value indicates robust overall model significance. [5] Significant p-value indicates biased standard errors; use robust estimates. [6] Significant p-value indicates residuals deviate from a normal distribution.					

Figure 5: OLS summary output

Once a model has been specified and OLS has been run, the next step is to perform a spatial autocorrelation test on the OLS residuals to determine if the model is biased. This is accomplished using the Spatial Autocorrelation (Morans I) tool, as seen in Figure 6. The tool outputs an *.html file that shows spatial autocorrelation results (z-score) on a bell curve. If the z-score falls to the far left or far right, then it is significant, but if they fall in the middle, they are determined to be no more significant than random chance. If the z-score is significant, the model is missing important explanatory variables.

Input Feature Class		
2010\OLS		I 🖻
Input Field		
Residual		-
Generate Report (optional)		
Conceptualization of Spatial Relationships		
POLYGON_CONTIGUITY_(FIRST_ORDER)		-
Distance Method		
EUCLIDEAN_DISTANCE		
Standardization		
ROW		-
Distance Band or Threshold Distance (optional)		
Weights Matrix File (optional)		
	OK Cancel Environments.	Show Help >>

Figure 6: Spatial Autocorrelation (Morans I) Tool

Using the above tools to perform data analysis provided enough information to start to answer the question posed by this study.

IV. Results

Does obesity have statistically significant correlations? Yes, 5 explanatory variables account for 83% of the variance of obesity and they are: people aged 35 – 44 that have diabetes making an income of \$35,000 - \$49,999 that have had no physical activity in the last 30 days and have not had any permanent teeth extracted.

Mapping obesity rates across the United States clearly allows for visualization of clustering, as can be seen in Figure 7. According to the map, states displayed as red have the highest rates of obesity in the country and appear to be clustered in the South.

Since the simple mapping of the data appears to be clustered, running the Hot Spot Analysis tool will identify significant hot and cold spot clustering. Figure 8 shows that there is a statistically significant hot spot cluster in the South consisting of the states OK, AR, LA, MS, AL, GA, TN, KY, IN, and OH. Likewise, it also shows two cold spot clusters – one out West consisting of WY, UT, and NV and the second in New England consisting of NY, NJ, CT, RI, MA, NH, VT, and ME.

With the Hot Spot Analysis (Getis-Ord Gi) proving that there are statistically significant hot and cold clusters, the next step is settling on a model that describes and accounts for the clustering of obesity. This is accomplished with a linear regression OLS test of the dependent and explanatory variables. For this OLS, the dependent variable is Obesity (A4409_3) and as a result of the Exploratory Regression tool, the explanatory variables are Age = 35 - 44 (A187_3), Diabetic (A1363_1), Income = \$35,000 - \$49,999 (A2976_4), No Physical Activity Last 30 Days (A4347_2), and No Permanent Teeth Extracted (A6607_2). Figure 9 displays the resulting map of OLS residuals. The red states indicate that the actual observed values are higher than values predicted by the model. On the other hand, blue states show where actual values are lower than predicted by the model.

The other output of OLS is the summary file. While Figure 9 is a way to visualize the model, the summary table (Figure 10) contains the meat for determining how well suited is the model to the dependent variable. First, all the explanatory variables have a positive coefficient, which means each variable has a positive correlation with obesity. Second, all the explanatory variables appear to be statistically significant at the 0.05 level. Statistical significance is shown with an asterisk, as can be seen in Figure 10. Next is the Koenker Statistic which checks for non-stationarity. Since it's not significant (no asterisk), then probability should be used over the robust probability. Fourth, determine if any of the explanatory variables display redundancy by reviewing the VIF column. The general guideline is if the VIF is above 7.5, then the variable displays multicollinearity, but smaller is better. All values appear to be under the guidelines and therefore multicollinearity is not an issue.

Now that all the explanatory variables have been individually determined to be significant, the model itself needs to be evaluated. This comes in the form of the R^2 value, in this case $R^2 = 0.832578$. This says that the model explains ~83% of the variance of the dependent variable, which is very good. Finally, the summary report also displays the Jarque-Bera Statistic which determines whether or not the model has bias. The statistic is not significant (no asterisk), which means the model residuals are normally distributed.



Figure 7: 2010 Obesity Rates by State



Figure 8: Hot/Cold Spot Analysis of 2010 BRFSS Obesity data by State

With the Hot Spot Analysis (Getis-Ord Gi) proving that there are statistically significant hot and cold clusters, the next step is settling on a model that describes and accounts for the clustering of obesity. This is accomplished with a linear regression OLS test of the dependent and explanatory variables. For this OLS, the dependent variable is Obesity (A4409_3) and as a result of the Exploratory Regression tool, the explanatory variables are Age = 35 – 44 (A187_3), Diabetic (A1363_1), Income = \$35,000 - \$49,999 (A2976_4), No Physical Activity Last 30 Days (A4347_2), and No Permanent Teeth Extracted (A6607_2). Figure 9 displays the resulting map of OLS residuals. The red states indicate that the actual observed values are higher than values predicted by the model. On the other hand, blue states show where actual values are lower than predicted by the model.

The other output of OLS is the summary file. While Figure 9 is a way to visualize the model, the summary table (Figure 10) contains the meat for determining how well suited is the model to the dependent variable. First, all the explanatory variables have a positive coefficient, which means each variable has a positive correlation with obesity. Second, all the explanatory variables appear to be statistically significant at the 0.05 level. Statistical significance is shown with an asterisk, as can be seen in Figure 10. Next is the Koenker Statistic which checks for non-stationarity. Since it's not significant (no asterisk), then probability should be used over the robust probability. Fourth, determine if any of the explanatory variables display redundancy by reviewing the VIF column. The general guideline is if the VIF is above 7.5, then the variable displays multicollinearity, but smaller is better. All values appear to be under the guidelines and therefore multicollinearity is not an issue.

Now that all the explanatory variables have been individually determined to be significant, the model itself needs to be evaluated. This comes in the form of the R^2 value, in this case $R^2 = 0.832578$. This says that the model explains ~83% of the variance of the dependent variable, which is very good. Finally, the summary report also displays the Jarque-Bera Statistic which determines whether or not the model has bias. The statistic is not significant (no asterisk), which means the model residuals are normally distributed.



Figure 9: OLS Model of 2010 BRFSS data

🔲 Untitled - Notepad	
<u>File E</u> dit F <u>o</u> rmat <u>V</u> iew <u>H</u> elp	
Summary of OLS Results Variable Coefficient StdError t-Statistic Probability Robust_SE Robust_t Robust_Pr VIF [1] Intercept -24.119805 6.690074 -3.605312 0.000806* 5.886651 -4.097373 0.000182* A187_3 0.202563 0.064736 3.129057 0.003147* 0.045652 4.437114 0.000063* 1.273297 A1363_1 1.151611 0.262046 4.394683 0.000072* 0.243186 4.735509 0.000024* 4.506732 A2976_4 0.669754 0.126221 5.306200 0.000004* 0.124742 5.369132 0.00003* 1.599890 A4347_2 0.561788 0.092851 6.050399 0.000000* 0.093117 6.033133 0.000000* 3.233553 A6607_2 0.258958 0.064894 3.990498 0.000253* 0.058617 4.417789 0.000066* 4.156409	~
OLS DiagnosticsA4409_3Input Features:2010 BRFSsNumber of Observations:49Multiple R-Squared [2]:0.832578Adjusted R-Squared [2]:0.813110Joint F-Statistic [3]:42.767070Yolk Wald Statistic [4]:328.272534Statistic [5]:2.416182Prob(>chi-squared), (5) degrees of freedom:0.789062Jarque-Bera Statistic [6]:2.159437Prob(>chi-squared), (2) degrees of freedom:0.339691	
Notes on Interpretation * Statistically significant at the 0.05 level. [1] Large VIF (> 7.5, for example) indicates explanatory variable redundancy. [2] Measure of model fit/performance. [3] Significant p-value indicates overall model significance. [4] Significant p-value indicates robust overall model significance. [5] Significant p-value indicates biased standard errors; use robust estimates. [6] Significant p-value indicates residuals deviate from a normal distribution.	~

Figure 10: OLS Summary Results

Since the OLS model has been determined to be a good fit, there just needs to be a final process that checks the regression model's over and under predictions are not clustered. As seen in Figure 11, the output of the Spatial Autocorrelation (Morans I) tool shows that the residuals display a random spatial pattern and are not clustered. Taking all these factors into account, it appears that the OLS model is well specified.

So *Does obesity have statistically significant correlations*? Yes and the variables hint at who is most likely to be obese and that conventional wisdom still stands – physical activity is important to control obesity. The two most surprising results are the income level and permanent teeth extraction results. One hypothesis about income level is that people making less than \$35,000 may receive some Federal assistance, e.g. food stamps, that cannot be used at fast food places. Likewise, people that make \$50,000 or more might be more inclined to spend more on fresh fruits and vegetables. The no permanent teeth extraction variable is still a mystery.



Figure 11: Spatial Autocorrelation (Morans I) of the OLS Residuals

V. Summary

Upon reflection, I feel this study was a practical use of the knowledge and tools learned throughout the course. Finding the Exploratory Regression tool was hugely helpful in checking variables contained within the data set that I would not have had time to do otherwise. I had hoped to continue back in time to provide trending analysis, but between inconsistencies in the data and the amount of time needed to determine variables, it was not to be in this study. Probably the biggest disappointment with this study was not being able to obtain finer level data to do more detailed distance analysis, i.e. distance to grocery store, proximity to fast food restaurants, etc. Time permitting, this study could certainly be extended and supplemented with more detailed data over a longer period of historical time. It could be effective in creating obesity education programs, as well as knowing exactly where to concentrate those education efforts.

VI. References

O'Sullivan, David and Unwin, David. Geographic Information Analysis – 2nd ed. Hoboken, New Jersey: John Wiley & Sons, Inc.

GEOG 586: Geographic Information Analysis. https://www.e-education.psu.edu/geog586/

CDC BRFSS GIS Data Files. http://www.cdc.gov/brfss/maps/gis_data.htm

Merriam-Webster. http://www.merriam-webster.com/dictionary/obesity