

What is the Relationship Between Fast Food and Diabetes?



Group 6
Johanna Pina,
Giulia Tasca,
Rachel Fletcher,
Omar Nasir

Agenda

- Introduction and Initial Thoughts
- Questions We Hoped to Answer
- DataBase Discussion
- Machine Learning Stuff
- Visualizations
- Reflection
- Fast Food Industry: Business Case



Introduction and Initial Thoughts



- Exploring the relationship between number of fast food restaurants and the number of Diabetes cases in each state
- Finding an underlying relationship between the variables
- Topic selected based on common interests among team members
- Various datasets on Kaggle regarding fast food restaurants

- Several other topics were considered such as personality tests, Michelin restaurants, unemployment rates and topics related to the pandemic.
- wanted to do a project relating to predicting and discovering hidden relationships between prevalent topics in our country; fast food restaurants are big in the U.S. and also diabetes is very prevalent
- we settled on Fast Foods and diabetes for logistical reasons, such as access to the necessary data sets; wanted to see if there is a relationship between the amount of restaurants and number of cases; if states with higher number of restaurants saw higher number of cases
- we decided we would probably go down the route of working on a simple linear regression

Questions We Hoped to Answer

- 1- Which states have the most and least fast food restaurants?
- 2- What are the top 10 Fast food restaurants nationally? Which Fast food restaurants have the most locations nationally?
- 3- What is the number of fast food restaurants per capita by state?
- 4- Which states have the most cases of diabetes?
- 5- What is the national number of diabetes cases? Number of diabetes related risk factors (smokers, obesity, non-activity levels)?
- 6- Does the number of fast food restaurants correlate to the number of diabetes cases by state?

DataBase Discussion



Data Sets

The screenshot shows two Google Sheets side-by-side. The left sheet, titled 'CleanFastFoodRestaurants', lists various fast-food restaurant locations with columns for address, city, coordinates, and state. The right sheet, titled 'joined_df', merges this data with CDC diabetes statistics, including state names, population, and diabetes rates per capita.

	A	B	C	D	E	F	G	H	I	A	B	C	D	E	F	G	H	I	
1	address	city	cour	latitude	longitude	name	postalc	prov	State	1	state_name	state_population	number_of_rest	number_of_dia	diabetes_per_c	cases	number_of_over	number_of_in	
2	4402 31st Ave	Astoria	US	40.75949	-73.91374	31st Avenue Gyro	11103	NY	New York	2	Alabama	4874486	241	530018	0.000049441	0.10733105	29356	23640	91653
3	701 Belknap St	Superior	US	46.7209	-92.0878	AAW	54880	WI	Wisconsin	3	Alaska	739700	28	44862	0.000037853	0.060648912	23045	13385	6984
4	1200 N Monroe St	Monroe	US	41.93287	-81.38903	AAW	48162	MI	Michigan	4	Arizona	7044008	494	501839	0.000070131	0.071243389	267683	227196	74093
5	8200 Vineland Ave	Orlando	US	28.38737	-81.49249	AAW	32821	FL	Florida	5	Arkansas	3001345	235	299417	0.000078252	0.065087941	180964	146161	63780
6	324 S Highway 24	Heyburn	US	42.57289	-113.70304	AAW	83336	ID	Idaho	6	California	3600197	176	250151	0.000064043	0.065087958	40203	96803	336933
7	1302 N US Highway 71	Carroll	US	42.07319	-94.87777	AAW	51401	IA	Iowa	7	Colorado	5611885	290	278905	0.000051676	0.048698987	14237	93596	39736
8	240 E Lake Mead Pkwy	Henderson	US	36.04105	-114.78989	AAW	89015	NV	Nevada	9	Connecticut	3573297	147	249166	0.000041138	0.069730000	129025	113182	36613
9	Riverton Crossings	Grandville	US	42.87952	-85.75576	AAW	49418	MI	Michigan	10	Delaware	956823	70	85505	0.00007159	0.083634485	49990	39148	12135
10	1484 W Central Ave	Sutherlin	US	43.38600	-123.33564	AAW	97479	OR	Oregon	11	District of Colum	694906	21	44352	0.000030220	0.063824460	24637	19128	9248
11	6001 E Division St	Nellisville	US	44.5528	-90.5901	AAW	54456	WI	Wisconsin	12	Florida	20963613	1026	1762208	0.000048942	0.084060319	959796	843137	237789
12	2240 W Galena Blvd	Aurora	US	41.80336	-88.32579	AAW	60506	IL	Illinois	13	Georgia	10410330	148	912045	0.000051726	0.087837273	484904	391559	118401
13	1401 N Michigan St	Plymouth	US	41.35081	-86.31043	AAW	46563	IN	Indiana	14	Hawaii	124393	69	121045	0.000048442	0.085017267	52312	38120	185665
14	210 Town Center Dr	Dearborn	US	42.32052	-83.22409	AAW	48126	MI	Michigan	15	Idaho	17177175	148	118417	0.000051726	0.087837273	60747	47102	16672
16	5770 W Irlo Bronson Memorial H	Kissimmee	US	28.32955	-81.51593	AAW	34746	FL	Florida	17	Iowa	125285	72	600234	0.000056421	0.074309245	524633	39152	108377
17	713 Happy Valley Rd	Glasgow	US	37.00903	-85.92107	AAW	42141	KY	Kentucky	18	Kansas	6656078	588	596820	0.000048816	0.089488288	343652	265392	114447
18	1100 N Moore Ave	Oklahoma City	US	35.348	-97.49695	AAW	73160	OK	Oklahoma	19	Kentucky	3141550	266	230654	0.000048472	0.073420445	139382	98557	33225
19	2760 W Chandler Blvd	Chandler	US	33.30649	-111.89025	AAW	85224	AZ	Arizona	20	Louisiana	3474628	588	221817	0.000050770	0.072593689	134125	101362	41363
21	3121 Sherwood Way	San Angelo	US	31.4457	-100.4785	AAW	76901	TX	Texas	21	Maryland	4452268	473	422540	0.0000196238	0.094904440	259237	246622	90408
22	1215 N Commerce St	Andmore	US	34.18947	-97.143211	AAW	73401	OK	Oklahoma	23	Maine	50317	148	110445	0.000048453	0.081969891	56570	45097	18230
23	806 Rostrover Rd	Belle Vernon	US	40.13648	-81.51593	AAW	15012	PA	Pennsylvania	24	Massachusetts	1334612	66	109393	0.000048453	0.081969891	277456	201829	66125
25	2522 E University Ave	Des Moines	US	41.60054	-93.59976	AAW	50317	IA	Iowa	26	Michigan	220386	308	501275	0.000051130	0.083214805	204647	16462	91952
26	W Franklin St	Chapel Hill	US	35.91119	-70.05745	AAW	27514	NC	North Caroli	27	Minnesota	9711161	321	407172	0.000051046	0.082394067	204701	16462	91952
27	250 W Baltimore St	Wilmington	US	41.30548	-78.15374	AJ's Hotdogs & Gyros	60481	IL	Illinois	28	Mississippi	6566230	321	350966	0.000057669	0.063052730	197050	134705	47894
28	1480 Ocean Ave	Rumson	US	40.35016	-73.97301	Ama Ristorante	7760	NJ	New Jersey	29	Missouri	2988510	130	320701	0.000043500	0.107311338	183621	148715	64568
29	5540 Old Cheney Rd	Lincoln	US	40.75546	-96.64543	Amigos/Kings Classic	68516	NE	Nebraska	30	Montana	6106870	473	521416	0.000054158	0.085384670	273089	22036	96101
30	1002 J St	Auburn	US	40.36457	-96.00095	Amigos/Kings Classic	68505	NE	Nebraska	31	Nebraska	1052482	57	71267	0.000054158	0.085384670	35938	29500	13633
31	1411 Q St	Lincoln	US	40.814500	-96.00095	Amigos/Kings Classic	68508	NE	Nebraska	32	Nebraska	191947	173	138305	0.000090230	0.072186235	75470	52128	20201
32	5701 NW 1st St	Lincoln	US	40.86756	-96.23737	Amigos/Kings Classic	68521	NE	Nebraska	33	Nevada	2969905	245	232865	0.000082494	0.078415978	120647	105801	32364
34							7760	NJ	New Jersey	34	New Hampshire	8850525	265	60898	0.000047097	0.075103778	56915	39552	14742
35							7760	NJ	New Jersey	35	New Jersey	18586572	217	1469506	0.000107799	0.070309780	95433	70168	78181
36							7760	NJ	New Jersey	36	New York	18586572	588	1510870	0.000029812	0.077126238	734990	611125	219983
							7760	NJ	New Jersey		North Carolina	10262233	634	884218	0.008111992	0.042922	384294	186606	
							7760	NJ	New Jersey		North Dakota	754947	40	50142	0.000119548	0.069975584	27175	31175	10516

- screenshot of some of the data
- the table on the left is the restaurant data that we got from Kaggle
- the table on the right is the dataset with state name, population, diabetes cases, smoker cases, obesity cases, and non-active cases
- population and predictor data was found from the CDC website for Diabetes in the U.S.
- The data set was obtained from kaggle.com and was uploaded by Datafiniti's Business Database. The data can be found [here](#). We have two sets of restaurant data that were merged using postgres SQL. A dataframe was created to add the diabetes data. The diabetes data set was obtained from the CDC website and can be found [here](#). Additionally, from the CDC we obtained more data sets of other features (i.e. obesity rates, smoking rates and non-activity rates) that may also contribute to number of diabetes cases. The last data set was obtained from census.gov to obtain the population from each state to calculate per capita rates.

Database Integration

www.quickdatabasediagrams.com

Fast_Food_Restaurants1

City	varchar(40)
Address	varchar(50)
Postal_Code	int
Lat	int
Lon	int
Country	varchar(5)
Category	varchar(50)
Name	varchar(50)

Fast_Food_Restaurants2

City	varchar(40)
Address	varchar(50)
Postal_Code	int
Lat	int
Lon	int
Country	varchar(5)
Category	varchar(50)
Name	varchar(50)

Population

State_Name	varchar(50)
Population	int

Diabetes

State_Name	varchar(50)
Number	int

Obesity

State_Name	varchar(50)
Number	int

Smokers

State_Name	varchar(50)
Number	int

Non_Activity

State_Name	varchar(50)
Number	int

Joins:

- We first did a full join on the fast_food_1 and fast_food_2 datasets found on Kaggle; this way we could have all of the data for the fast food restaurants and try to get as much data for each city that was available
- we then did a full outer join on the new fast_food_complete dataset with the population dataset we found for the years 2017 and 2018; we joined on state_name so that the new table would have all of the information from the giant fast food restaurant table and with the edition of population for each state
- we then first did a full join on the diabetes, obesity, smokers, and non-activity tables on state names to get one table with state_name, each predictor, and their number of cases
- finally, we did a final full join of the restaurant, state, and population table with the predictors table to get a final giant dataset containing all of the restaurants, their information, state name, population, predictor variables and the number of cases for each of them

Machine Learning

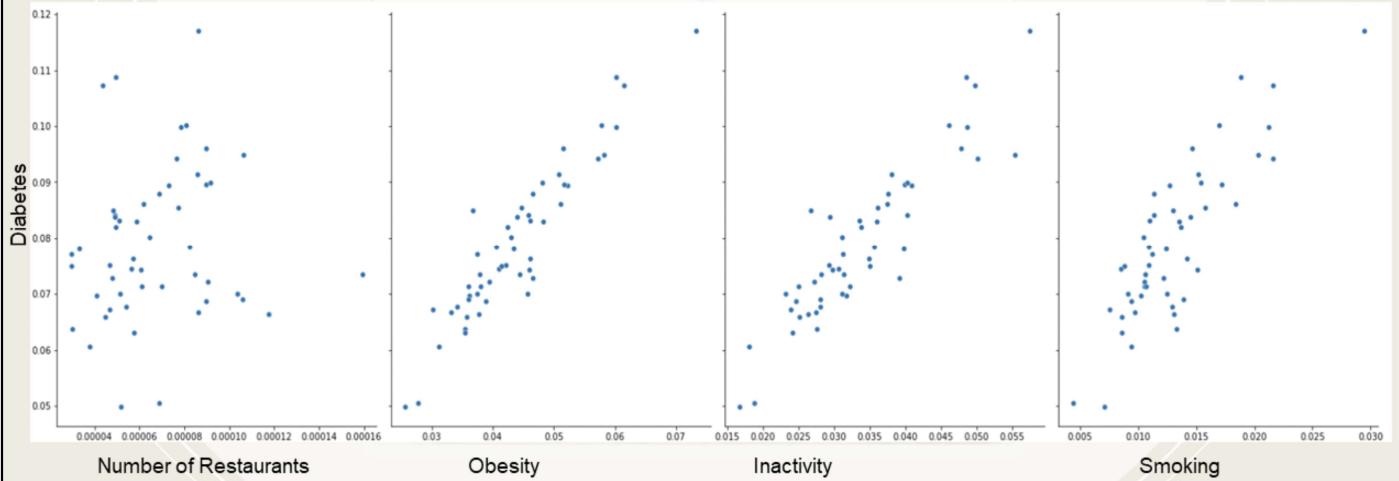


Clean Data Using Heatmap Dataframe

	state_name	restaurants_per_capita	diabetes_per_capita	obesity_per_capita	inactivity_per_capita	smokers_per_capita
0	Alabama	0.000049	0.108733	0.060223	0.048518	0.018803
1	Alaska	0.000038	0.060649	0.031155	0.018095	0.009442
2	Arizona	0.000070	0.071243	0.038002	0.032254	0.010519
3	Arkansas	0.000078	0.099761	0.060204	0.048699	0.021250
4	California	0.000045	0.065988	0.035662	0.025120	0.008615
5	Colorado	0.000062	0.049699	0.025702	0.016678	0.007081
6	Connecticut	0.000041	0.069730	0.036108	0.031674	0.010218
7	Delaware	0.000073	0.089363	0.052214	0.040915	0.012683
8	District of Columbia	0.000030	0.063824	0.035454	0.027526	0.013308
9	Florida	0.000049	0.084060	0.045784	0.040219	0.011343
10	Georgia	0.000069	0.087837	0.046579	0.037613	0.011373
11	Hawaii	0.000048	0.085017	0.036726	0.026762	0.013034
12	Idaho	0.000086	0.066744	0.033167	0.027449	0.009706
13	Illinois	0.000056	0.074359	0.041055	0.030631	0.008481
14	Indiana	0.000090	0.089488	0.051614	0.039860	0.017189
15	Iowa	0.000085	0.073420	0.044367	0.031372	0.010576
16	Kansas	0.000057	0.076259	0.046111	0.034848	0.014220
17	Kentucky	0.000106	0.094904	0.058226	0.055392	0.020306
18	Louisiana	0.000090	0.096020	0.051510	0.047886	0.014613
19	Maine	0.000049	0.081970	0.042387	0.033790	0.013659

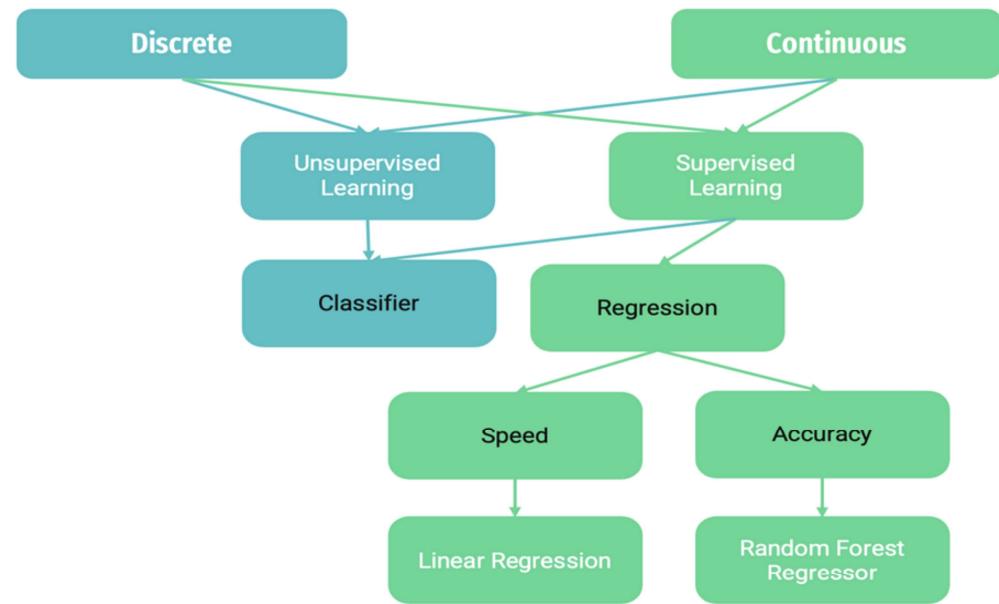
All the values were converted by dividing the number of cases by population per state. The heatmap was created with python using `.style.background_gradient()`. The darker color represents the highest values, which we can already see are similar for diabetes, obesity, inactivity and smoking but not for number of restaurants.

Relationship Between the Variables



In this visualization we plotted the y and the various x features. Here we can also see that Obesity, inactivity and smoking appear to have a strong positive correlation with diabetes. Whereas, number of restaurants appears more scattered.

Choosing the Model



Continuous: Our predictions would result in continuous data, meaning the data itself would not be classified/discrete. Instead the predictions would be similar values to what we have (which can be any number/value).

Supervised: The machine will not need to independently find which groups or clusters to put the values. The data is already separated by its labels.

Speed: Linear Regression is a simple model to decide if there is a correlation between the x & y variables. This model performs quickly with large data sets.

Accuracy: Random Forest Regressor splits the data into trees, each tree is formed from random rows. Each tree then provides its own prediction and the final result is the average of all of the predictions, making it more accurate by reducing overfitting (models the training data too well).

Define Variables & Split the Data

```
#define the variables
features=['restaurants_per_capita', 'obesity_per_capita',
          'inactivity_per_capita', 'smokers_per_capita']
X = final_df[features].values.reshape(-1, len(features))
y = final_df['diabetes_per_capita'].values

# split data using 80/20 ratio
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1, test_size=0.2)
```

Multiple Linear Regression Model

```
# instantiate and fit
model = LinearRegression()
model.fit(X_train, y_train)

# predict
y_pred = model.predict(X_test)
```

Random Forest Regression Model

```
# instantiate and fit
rf_model = RandomForestRegressor(n_estimators=128)
rf_model = rf_model.fit(X_train, y_train)

# Making predictions using the testing data
predictions = rf_model.predict(X_test)
```

- Here we run the machine learning model. To implement the models, first we define X with the features of number of restaurants, obesity, inactivity and smoking. Y is defined as diabetes cases per capita. The purpose of the model is to be able to predict Y (cases of diabetes) when X is given.
- The next step was to split the data into training and testing sets using a 80/20 ratio. The machine will train with 80% of the data and then will predict 20% of the values. We have the actual values for the 20% which we will then compare to the values the machine learning model predicted allowing us to evaluate the accuracy.
- Next we instantiate the model. This tells the machine which “template” we are using so that it knows what formula to access.
- Fitting refers to how well the model generalizes to data on which it was trained. The model then predicts y, using the X testing data.

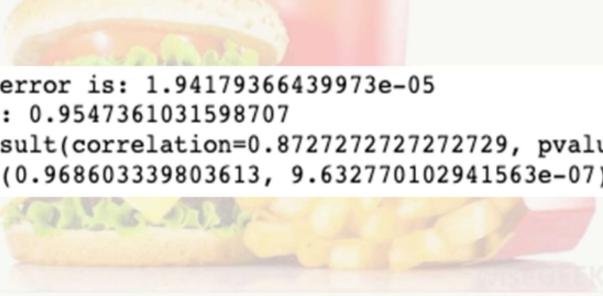
Evaluating the Models

```
Linear Regression mean square error is: 9.501711071642502e-06
```

```
Linear Regression R squared is: 0.8891662637660523
```

```
Linear Regression: SpearmanrResult(correlation=0.8727272727272729, pvalue=0.00045461505140964044)
```

```
Linear Regression pearson coefficient is: (0.9709414518591231, 6.82305350681304e-07)
```



```
Random Forest mean square error is: 1.94179366439973e-05
```

```
Random Forest R squared is: 0.9547361031598707
```

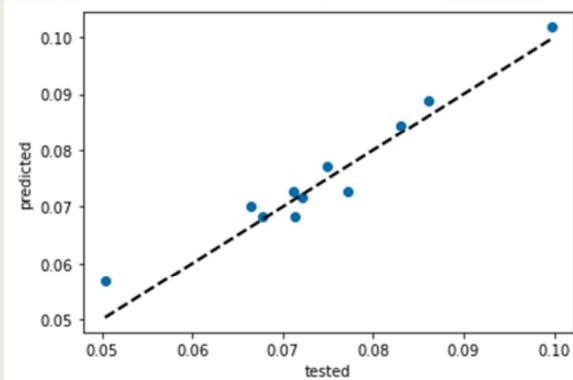
```
Random Forest: SpearmanrResult(correlation=0.8727272727272729, pvalue=0.00045461505140964044)
```

```
Random Forest pearson is: (0.968603339803613, 9.632770102941563e-07)
```

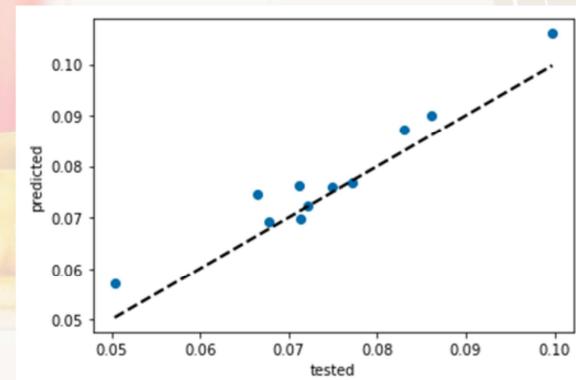
- The mean square error: is the variance around the fitted regression line. The lower the number, the smaller the “errors/ variance”.
- R squared: measures the strength of the relationship between the model and the rates of diabetes. Generally, the higher the better, but if it is too close to 100%, then there could be overfitting.
- Spearman Correlation Coefficient: can take values from +1 to -1. +1 indicates a perfect association, zero indicates no association between ranks and -1 indicates a perfect negative association.
- P-value: The smaller the p-value, the stronger the evidence that you should reject the null hypothesis. The p-value is the probability that the null hypothesis is true.
- Pearson Correlation Coefficient: If the coefficient value lies between ± 0.50 and ± 1 , then it is said to be a strong correlation.

Predicted & Tested number of Diabetes Cases per Capita

Multiple Linear Regression



Random Forest Regression



This plot shows the relationship between the diabetes data that was tested and the diabetes data that was predicted (y variable). The closer the dots are to the line, the higher the accuracy. However, if it is a perfect match, this could signal overfitting.

Result of Analysis

What this means for the population; what we can draw from our findings...

- Strong positive correlation between variables
- Model predicts the number of diabetes cases
- Model is able to calculate reduction in features needed for preferred diabetes outcomes



- There is a strong positive correlation between diabetes and obesity, smoking and inactivity.
- One possible use case for this model includes predicting the number of diabetes cases for a given population when the values of obesity, smoking and inactivity are provided.
- Another use case could be to lower the number of diabetes to a predetermined level, by lowering the cases of obesity, smoking, and inactivity to a matching level.

Visualizations

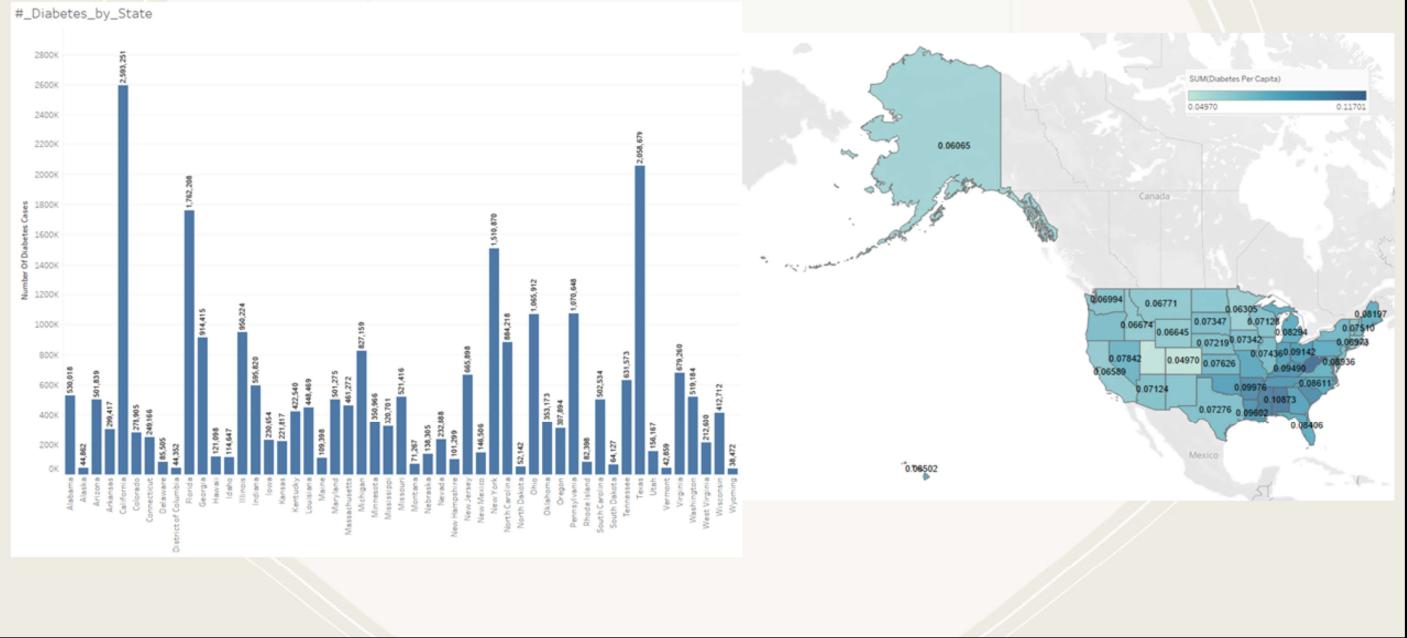


Factors that may increase risk of Type 2 Diabetes



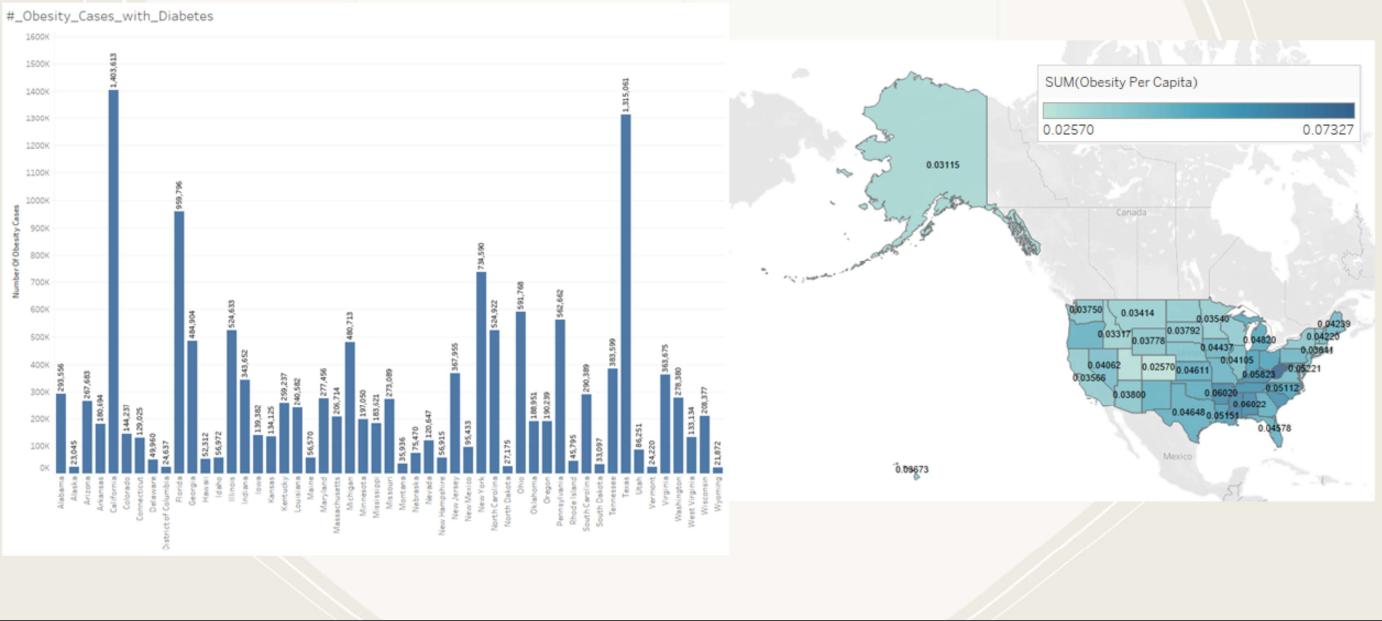
Only Type 2 Diabetes has any relation with fast food. If a person has particular bad eating habits by eating extremely unhealthy food, which causes obesity, doesn't exercise, and smokes will have a increased danger of developing type 2 diabetes. Now on the flip side if a person eats relatively healthy at fast food or on occasion, they are not likely to develop type 2 diabetes.

of Diabetes Cases in USA



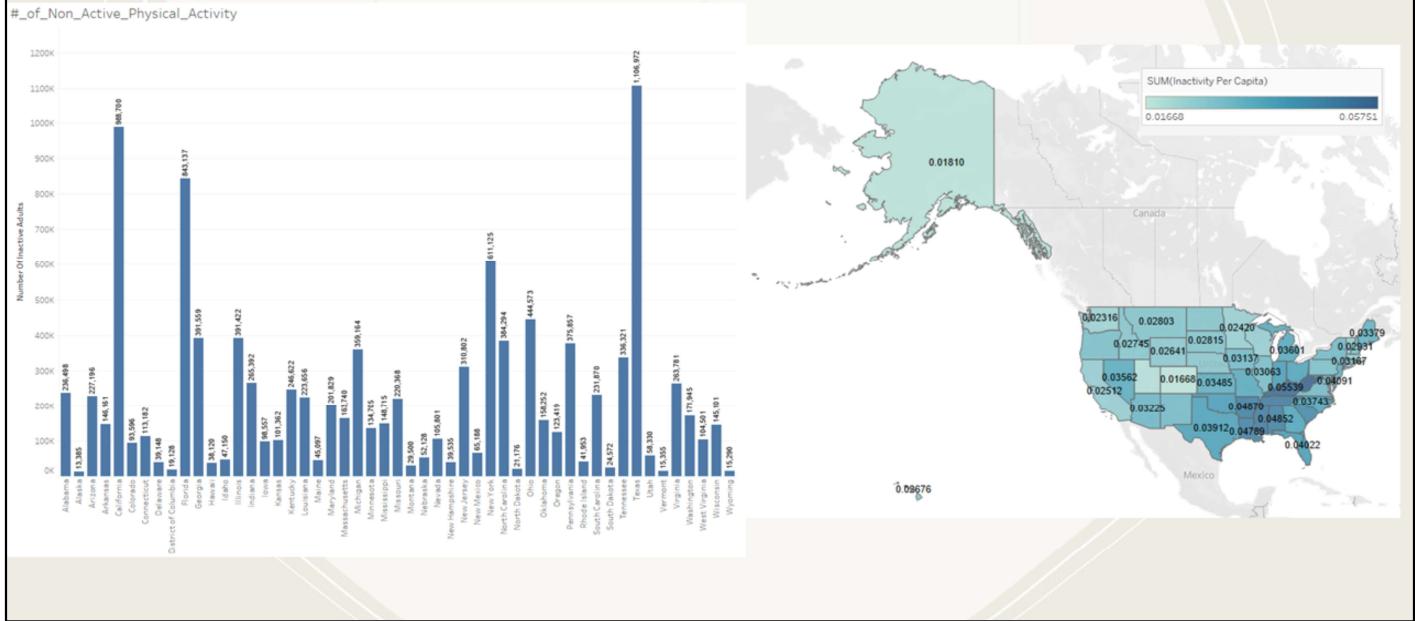
- State with Minimum cases: Wyoming 38,472
- State with minimum cases per capita: Colorado 0.04970
- State with maximum cases: California 2,593,251
- State with maximum cases per capita: West Virginia 0.11701
- Data set was pulled from CDC website:
<https://gis.cdc.gov/grasp/diabetes/DiabetesAtlas.html>
- Data was merged creating a csv file that included data that was pulled regarding diabetes cases in USA, how many cases are obese, how many cases are smokers, and how many cases do not do physical activity.
- The merged dataset that was used to build this map was from the “joined_df” data set.
- Longitude was dropped in the Columns field, Latitude was dropped in the Rows field, SUM(Number of Diabetes Cases) was dropped in the Color field under Marks, and State Name was dropped in the detail field under Marks

of Diabetes Cases who are Obese

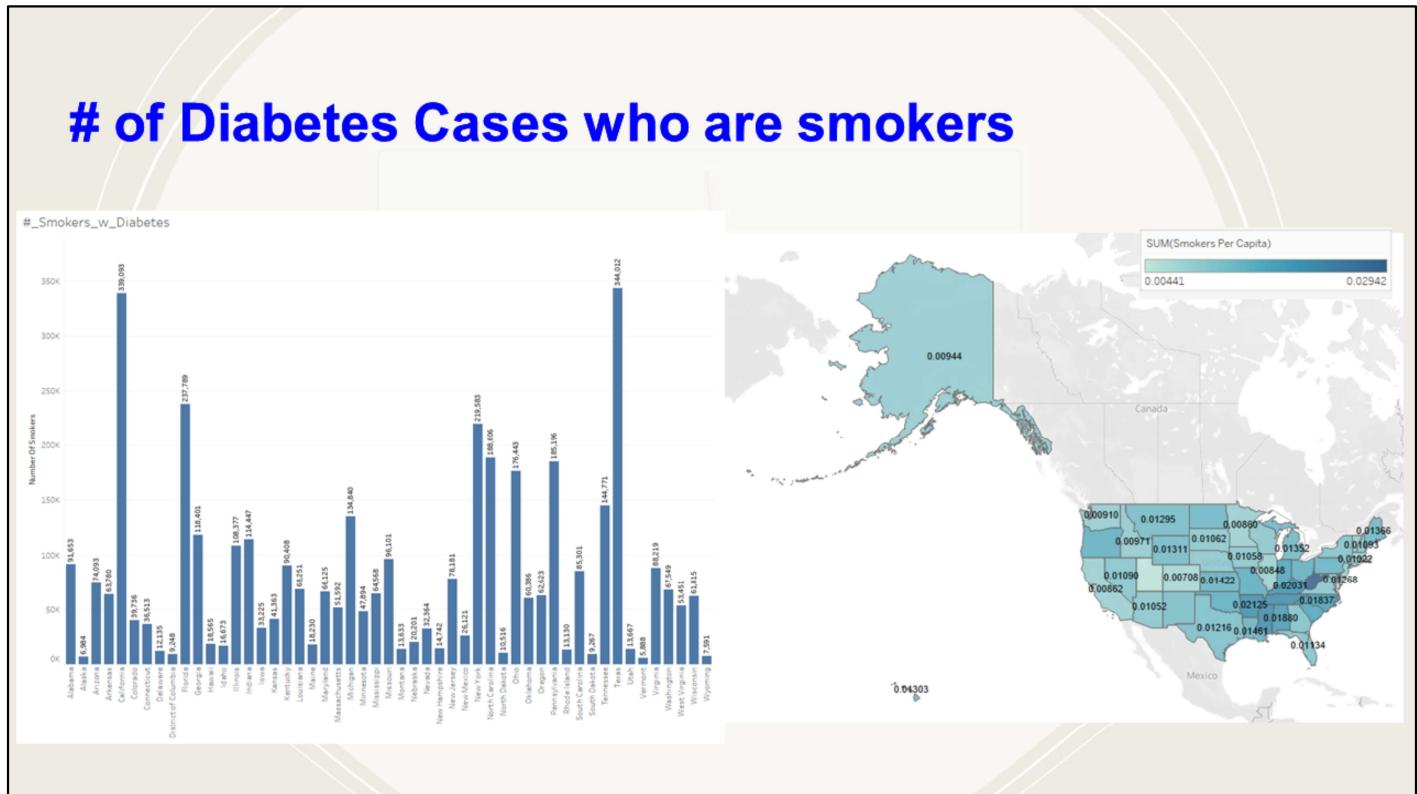


- State with Minimum cases: Wyoming 21,872
- State with minimum cases per capita: Colorado 0.02570
- State with maximum cases: California 1,403,613
- State with maximum cases per capita: West Virginia 0.07327
- Data set was pulled from CDC website:
<https://gis.cdc.gov/grasp/diabetes/DiabetesAtlas.html>
- Data was merged creating a csv file that included data that was pulled regarding diabetes cases in USA, how many cases are obese, how many cases are smokers, and how many cases do not do physical activity.
- The merged dataset that was used to build this map was from the “joined_df” data set.
- Longitude was dropped in the Columns field, Latitude was dropped in the Rows field, SUM(Number of Obesity Cases) was dropped in the Color field under Marks, and State Name was dropped in the detail field under Marks

of Diabetes Cases who are not physically active



- State with Minimum cases: Alaska 13,385
- State with minimum cases per capita: Colorado 0.01668
- State with maximum cases: Texas 1,106,972
- State with maximum cases per capita: West Virginia 0.05751
- Data set was pulled from CDC website:
<https://gis.cdc.gov/grasp/diabetes/DiabetesAtlas.html>
- Data was merged creating a csv file that included data that was pulled regarding diabetes cases in USA, how many cases are obese, how many cases are smokers, and how many cases do not do physical activity.
- The merged dataset that was used to build this map was from the “joined_df” data set.
- Longitude was dropped in the Columns field, Latitude was dropped in the Rows field, SUM(Number of Inactive Adults) was dropped in the Color field under Marks, and State Name was dropped in the detail field under Marks

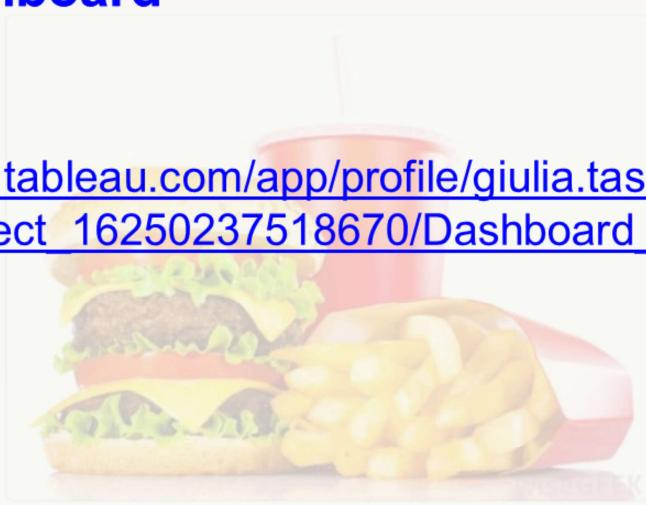


- State with Minimum cases: Vermont 5,888
- State with minimum cases per capita: Utah 0.00441
- State with maximum cases: Texas 344,012
- State with maximum cases per capita: West Virginia 0.02942
- Data set was pulled from CDC website:
<https://gis.cdc.gov/grasp/diabetes/DiabetesAtlas.html>
- Data was merged creating a csv file that included data that was pulled regarding diabetes cases in USA, how many cases are obese, how many cases are smokers, and how many cases do not do physical activity.
- The merged dataset that was used to build this map was from the “joined_df” data set.
- Longitude was dropped in the Columns field, Latitude was dropped in the Rows field, SUM(Number of Smokers) was dropped in the Color field under Marks, and State Name was dropped in the detail field under Marks

Link to Dashboard

https://public.tableau.com/app/profile/giulia.tasca/viz/FastFoodProject_16250237518670/Dashboard_Graphics

1



Dashboard consists of:

- # of Diabetes cases in the US
 - The merged dataset that was used to build this map was from the “joined_df” data set.
 - Field used SUMNumberofDiabetesCases was dropped in the Text field under “Marks”
- Risk Factors for Type 2 diabetes
 - The merged dataset that was used to build this map was from the “joined_df” data set.
 - Field used SUMNumberofInactiveAdultsCases was dropped in the Text field under “Marks”
 - Field used SUMNumberofObesityCases was dropped in the Text field under “Marks”
 - Field used SUMNumberofSmokers was dropped in the Text field under “Marks”
- Heat Map of Diabetes by State for 2017 and 2018
 - The dataset that was used to build this heat map was the DiabetesAtlasData dataset.
 - Columns included the Year field, Rows was the State Field, SUM(Number) was dropped in Color and Detail under MARKS
- Top 10 restaurants for sample of dataset
 - The merged dataset that was used to build this map was from the “CleanFastFoodRestaurant_df” data set.

- Name count field was put in the Color field of Marks, CNT(Name) was put in Angle field of Marks, Max(Name) was put in Text field of Marks, CNT(Name) was put in Text field of Marks.
 - A filter was done to limit the # of restaurants to the Top 10
- Restaurants per capita using sample set
 - The merged dataset that was used to build this map was from the "CleanFastFoodRestaurant_df" data set.
 - Longitude was dropped in the Columns field, Latitude was dropped in the Rows field, SUM(Restaurants Per Capita) was dropped in the Color field under Marks, and State Name was dropped in the Detail field under Marks
- Diabetes cases per capita
 - The merged dataset that was used to build this map was from the "joined_df" data set.
 - Longitude was dropped in the Columns field, Latitude was dropped in the Rows field, SUM(Diabetes Per Capita) was dropped in the Color field under Marks, and State Name was dropped in the Detail field under Marks
- Count of fast food restaurants by state using a sample from a data set
 - The merged dataset that was used to build this map was from the "CleanFastFoodRestaurant_df" data set.
 - CNT(Name) was dropped in Columns field, State was dropped in Rows field

Changes During the Project

- Questions
- Data and Predictors
- Models



Analytics Project Life Cycle

The 5 Phases



- questions: working through our project, we realized that we weren't able to work with county data, due to the difficulty of breaking each state into counties because of overlapping issues, so we decided to focus solely on each state; this led us to take out the questions we had at the beginning pertaining to counties and any questions that seemed to be too similar; we could only get population data for the years 2017 and 2018; so we would have tried to be more specific and find a lot more data including more years
- data and predictors: working through the machine learning section of our analysis, we realized that having only fast food restaurants as the x variable wasn't going to help us create an accurate regression model for predicting and analyzing the relationship with diabetes; there are numerous factors that come into play regarding diabetes, various predictors, and we couldn't dive into the relationship we were interested in without exploring other predictors; this led us to going back to the CDC website and finding more data to work with; we then got data on 3 different risk factors for diabetes, obesity, smoking, and inactivity, and the number of cases in each state, and we added that data into our tables and completed our analysis with the new addition of data
- models: we initially thought to solely use simple linear regression for our model; we quickly realize that we needed to have different models to complete a sufficient analysis; we decided to then move forward with multiple linear regression, with our newly added predictors; we also created a random forest regression model to add to our analysis, which was different from our initial thoughts

Reflection

- Recommendation for Future Analysis
- Final thoughts



- Given more time for an analysis, we would recommend further research into predictors of Diabetes and include those variables in our data frames and machine learning models. Obviously there are various factors that contribute to Diabetes in the U.S., like location, genes, health and physical habits, and more. We recommend including as many X variables as possible to create a sufficient model that can accurately predict Diabetes and portray the strength of each correlation between the dependent variables and independent variable. Another recommendation for future analysis would be to include bar charts for the different predictors of Diabetes and their prevalence in each state. Creating these visuals would show how the severity of each predictor varies across the country. Finally, finding and including relevant data sets from other countries could help broaden and strengthen the analysis.
- We decided to work with data from a sample found on Kaggle and the CDC website for Diabetes in the United States. After working with the Diabetes and restaurant data and comparing it to the population size in the United States, we decided that if we had more time to complete this analysis, we would try to get bigger data sets that are more representative of the population sizes in each state. Another aspect of the project that we would have done differently is incorporating various X

variables from the beginning. As we worked on the machine learning section of the analysis, we realized that we should have incorporated more predictors of Diabetes. We can't formally comment on and find a direct relationship between Diabetes and fast food restaurants, as there are a ton of other factors that play into that relationship. We were able to incorporate various predictors in our analysis, but we wish we brought them in sooner. In addition, with more time we could have added more questions with relation to the newly added variables and visuals to explore these relationships. Finally, if we had more time for the project, we would have spent more time on the introductory analysis in Python, to make the database integration smoother.

- overall, we had a successful project and learned how to work as a team and work on an analytics project that constantly changes and learn how to adapt to one another; as you go through analyses, questions, data, processes, and more change and you have to be okay with updating and changing as you go, and sometimes starting fresh. this entire bootcamp has given us the tools to create sufficient analysis and answer any questions that we are interested in in personal or professional worlds