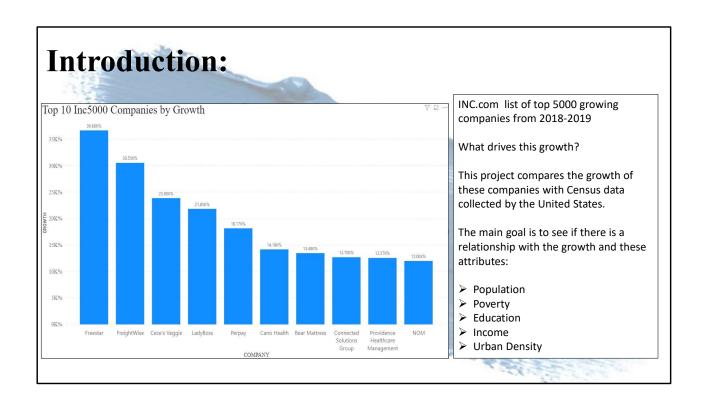
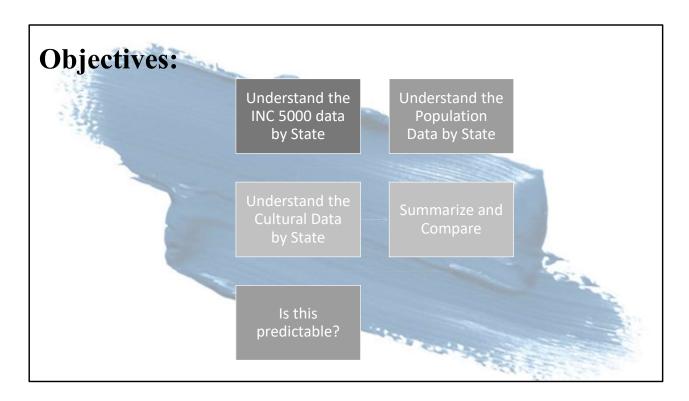


Agenda:

- **≻**Introduction
- **≻**Objectives
- ➤ Research Hypothesis
- ➤ Methodology
- **≻**Findings
- ➤ Conclusion and Recommendations



Main Point: "The goal of this analysis is to investigate what drove the growth of these companies. There are many variables and attributes that go into business growth, but for this analysis it will focus on geography, the education of the people living near these companies, the poverty, and earnings. It is essential to understand how these businesses have grown so fast and this analysis is the initial step in this direction. "Be sure to impress this is an investigation and the initial step in the workflow.



Main points:

- -The objective is to find any unique or unknown relationships or patterns in the dataset.
- -The next part is comparing a summarized dataset to geography in terms of what states the companies are in, the population of the states, and the amount of urbanization in each state.
- -The third part is to then compare this dataset to the summarized poverty, earnings, and education of each state.
- -The final objective is to show how these comparisons can be used to predict which business type will grow the fastest based on this list of 5000 companies across the

country.

By prediction is it possible to use the results from summarize and compare to model the outcome of business growth?

Hypothesis:

Does each variable have normal distribution?

- > Ho: It does not have normal distribution
- > Ha: It does have normal distribution

► Is there a strong correlation between the state population, growth, urban density, and Business Growth?

- ➤ Ho: There is no correlation
- ➤ Ha: There are strong or weak positive or negative correlations.
- > Is there a strong correlation between the population's education, earnings, and poverty with the growth of these companies?
 - ➤ Ho: There is no correlation
 - ➤ Ha: There are strong or weak positive or negative correlations.

> Is it possible to identify a combination of these attributes that would either enhance or discourage company growth?

- ➤ Ho: It is not possible or there is no definitive combination
- ➤ Ha: It is possible to find a specific combination
- This is the key slide for the analysis, this is what we are setting out to answer.
- 1. Is each variable normally distributed? In this case the null hypothesis (Ho) is the data does not have a normal distribution and the alternate (Ha) is the data has a normal distribution.
- 2. There is a strong correlation between the state population, growth, and urban density and business growth. Ho: there is no correlation or a very weak correlation, Ha: The correlation is strong
- 3. There is a clear relationship between the population's education, earnings, and poverty with the growth of these companies. Ho: there is no correlation or a very weak correlation, Ha: The correlation is strong.

4. If these relationships exist, it is then possible to identify the combination of these attributes that would either enhance the growth or discourage the growth of these companies. Ho: It is not possible because there are no correlations, Ha: It is possible to define a positive or negative impact on business growth based on the strong correlations.

The test of normality is mainly to ensure the data should be used for correlation however based on the results we know it is not. This may be a difficult point to sell.

Try and stress that correlation is used on data without a normal distribution daily.

Literature and Limitations: Literature Limitations ➤ Missing Data such as number of employees, > Difficult to find examples of this analysis initial capital, total costs, and other proprietary > However, it is clear small business growth is an essential part of economy There is a problem with the population in Puerto Rico this is not a serious limitation as > Maps have been used to find correlation we are focusing on the lower 48 states. > The following is the criterion to included in this list This led to a very thorough review of all the Be Privately-owned, based in the United States, and independent attributes used to find and mitigate any other Have started earning revenue by March 31, 2016 Had revenue no less than \$100,000 in 2016 Had revenue no less than \$2,000,000 in 2019 Revenue in 2019 exceeds revenue in 2016

Please stress these statistics:

"For this reason, the analysis will only focus on business growth and the following references helped build a framework on constitutes business growth in general.

(Bonsu & Kuofie, 2019) mention how small businesses are the key to the GDP of the United States. For this reason, understanding business growth would help small

businesses be more successful. This analysis can provide the foundation to build on by comparing these relationships with these larger companies identified by INC.com. The second hypothesis deals with population and geography. (Rosenblum et al., 2014) uses geography to correlate the location of Puerto Rican neighborhoods with the introduction of a cheaper type of heroin. A similar analysis can be quickly done by showing a map of where the 5000 companies are in relationship to city density. This idea is what led to this analysis in the first place. To further cement how important small businesses are (Farlie

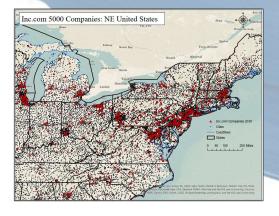
This idea is what led to this analysis in the first place. To further cement how important small businesses are (Farlie et al., 2019) looks at the universe of start-ups in the US and combines it with a panel data set. It also confirms the importance of small business by highlighting how there are nearly 5 million new businesses created annual which leads to over 3 million jobs created annually. "

"There is a criterion used by the census to define poverty and this data set is the count of those groups that meet these criteria that live above or below the 2019 threshold of \$31,275. According to this data set the number of households living below this threshold in Puerto Rico was 1,577,075 but the population of Puerto Rico in 2018 was 1,344,083. This difference led to an error in the results that had to be mitigated. Fortunately, this was the only case."

"

Methodology:

Collecting and Mapping



Discussion:

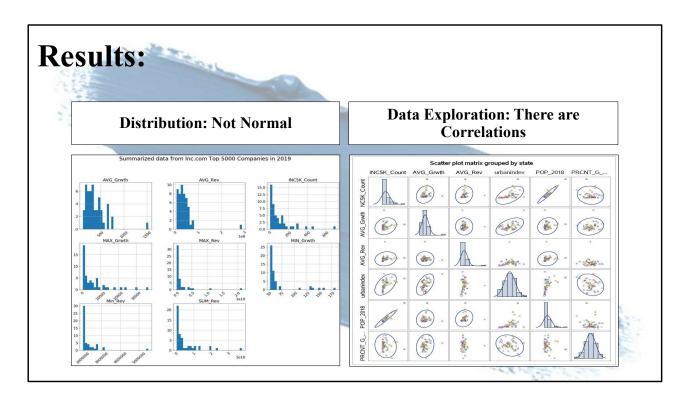
- This process involved scraping the 5000 records from the public Inc.com website into an excel sheet
- > To make the map the names of the cities had to be corrected and when a company used a regional or common name it had to be linked to the closest census designation.
- ➤ This data is at the city level with corresponding state and county.

This was the initial process that led to the entire project. After spending several weeks scraping the records and then performing an initial analysis it was clear this unique data set and warranted

Further investigation.

The decision was made to compare it to census data because this data is publicly available and also has multiple scales from city to state.

A quick list of the software and analysis



H1:

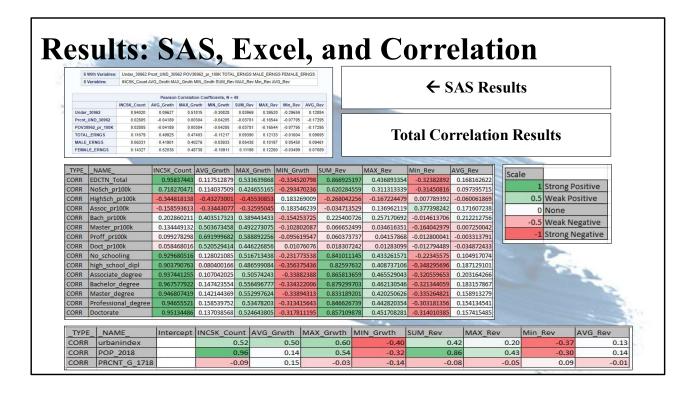
"The results of this analysis rejected the null on some of the hypotheses but not all of them. From the distribution of the attributes (Figure 3) and the following Shapiro-Wilks Tests (Figure 4) this data does not have a gaussian distribution. Ideally, this data should have a normal distribution before attempting correlation but because of the nature of this data finding attributes that have this distribution is very difficult. The null of the first hypothesis was not rejected."

Not normal distribution

H2: and H3:

"To test the second and third hypothesis the dataset was combined and imported into SAS Studio. Figures 5 and 6 show the results of the data exploration process. This process compares the correlation of all the selected variables. In Figure 5 there is a strong correlation with count of companies and total state population in 2018."

"Figure 6 (image on right) brings out some shocking results. It seems strange that the count of companies per state would have such strong positive correlations with all levels of education obtained by those over 25 and households living under an annual salary of \$30,962."

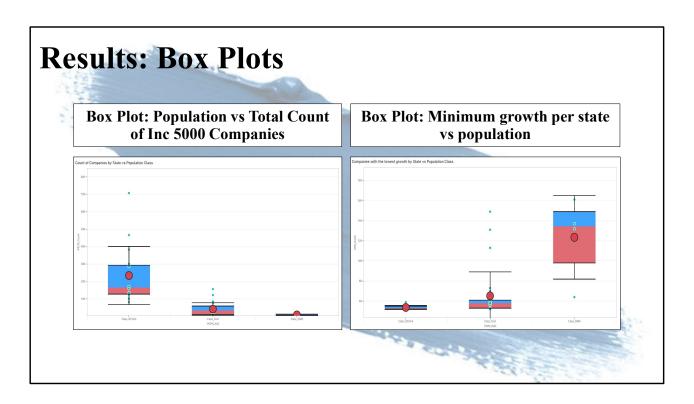


"The correlation tool in SAS Studio was used to dig deeper. Figure 8 shows a more detailed analysis of the correlations in excel using conditional formatting. The Dark Green is the greatest positive correlation, and the dark red is the most negative correlation."

"Based in this portion of the analysis both hypothesis two and three rejected the null and pushed the alternate as there are indeed some strong correlations between the summarized growth, revenue, and counts of the companies. However, more investigation is needed into why the count of education obtained seems to have such strong correlation with the count of companies per state and not with the growth. "

"However, the population has no correlation with average growth, but does have a weak negative correlation with minimum growth. Figure 9 (Bottom Table) shows this correlation table. "

Make a note of how the education correlations seem strange and need further investigation.



Please note the statistics and the outliers. It's important to stress this is one case and there are many more to be investigated.

"Based on this Figures 10 and 11 show the relationship with the new population class, the count of companies by state and the minimum growth by state, respectively. This analysis has shown that states with a

population more than 5 million contain most of the fastest growing companies, but the four states with population lower than 500,000 have the best performing companies out the companies with the slowest growth.

There are more cases to build and this analysis has mode this possible.

The Green points represent outliers and the 2 state with highest counts that have a population greater than 5 million are California and Texas. California has more than 700 of the companies on the list and they are all clustered around Los Angeles. Colorado and Utah have the most companies on the list for states with population between five hundred thousand and 5 million. In terms of minimum growth, the outliers for the state in *Class_5mil* were West

Virginia, Rhode Island, and New Mexico. The companies in the list in West Virginia averaged a growth of 149%. The outliers reveal how further investigation is necessary as well as a possible change in scale to get more granular in terms of counties or cities. "

Conclusion:

- > Of the original 4 hypotheses this analysis rejected the null on the last three but not the first. The data does not have normal distribution but there are strong correlations, and it is possible to combine them to build a case
- > There is a negative correlation with low growth and population. The strongest correlation showed the more population the more companies on the list you have
- This analysis has shown the even though lower population decreases the count, the states with a population less than 500,000 improve the companies with a lower growth rate
- This is the initial pass on the data and because of this workflow many more cases will be discovered
- > With the established workflow the data can be refined, and more relationships investigated

The main goal of this entire workflow is to sell it as something that can be applied to other hypothesis. There is more to be discovered based on these 5000 companies and we need to secure resources so we can continue this analysis.

"This analysis set out to test four specific hypotheses. The first dealt with the distribution and the null was not rejected. The second and third pushed the alternate hypothesis in terms of strong or weak correlations and the fourth pushed the alternate hypothesis because these correlations can be used to build cases about the company

growth. However, there are many more hypotheses that can now be answered because of the process has been defined and the analysis can be applied to any other set of variables. "

"Going forward, a different test about the distribution should be utilized and tested as the results of the education correlations may be a product of using a Pearson correlation on data that does not have normal distribution. The relationship between education per 100,000 and the growth should be investigated more as there is an expectation of higher education should warrant faster growth, but the findings of this analysis do not support that assumption. The goal of this analysis was not only to test the hypotheses but to also build a workflow than can now be applied to other variables and tests."

References:

- Bonsu, S., & Kuofie, M. (2019). Small business survival. Journal of Marketing & Management, 10(1), 51–63.
- Brownlee, J. (2018, August 14). 17 Statistical hypothesis tests in Python (cheat sheet). Machine learning mastery. https://machinelearningmastery.com/statistical-hypothesis-tests-in-python-cheat-sheet/
- US Census Bureau. (2020, January 22). Poverty thresholds: the united states census bureau. https://www.census.gov/data/tables/time-series/demo/income-poverty/historical-poverty-thresholds.html
- Fairlie, R. W., Miranda, J., & Zolas, N. (2019). Measuring job creation, growth, and survival among the universe of start-ups in the United States using a combined start-up panel data set. *ILR Review*, 72(5), 1262–1277. https://doi.org/10.1177/0019793919862764
- Fivethirtyeight/data. (n.d.). GitHub. https://github.com/fivethirtyeight/data/tree/master/urbanization-index
- Inc. 5000 2020: An exclusive guide to America's most inspiring entrepreneurs. (n.d.). Inc.com. https://www.inc.com/inc5000
- Rosenblum, D., Castrillo, F. M., Bourgois, P., Mars, S., Karandinos, G., Unick, G. J., & Ciccarone, D. (2014). Urban segregation and the US heroin market: A quantitative model of anthropological hypotheses from an inner-city drug market. International Journal of Drug Policy, 25(3), 543-555. https://doi.org/10.1016/j.drugpo.2013.12.008
- Rupasingha, A., & Wang, K. (2017). Access to capital and small business growth: Evidence from CRA loans data. Annals of Regional Science, 59(1), 15–41.
 https://doi.org/10.1007/s00168-017-0814-9
- Wach, K. (2020). A typology of small business growth modelling: A critical literature review. Entrepreneurial Business & Economics Review, 8(1), 159–184.
 https://doi.org/10.15678/EBER.2020.080109

Although these were not used in the slides, there are references in the notes.