

# PREDICTIVE MODELING

## PROJECT PRESENTATION

ANAND TAMHANKAR | SHAN XU | SHRUTHI TADAMIRI | URVI KOHLI

# AGENDA

- Business Problem
- SEMMA
- Business Value
- Recommendations/Conclusions

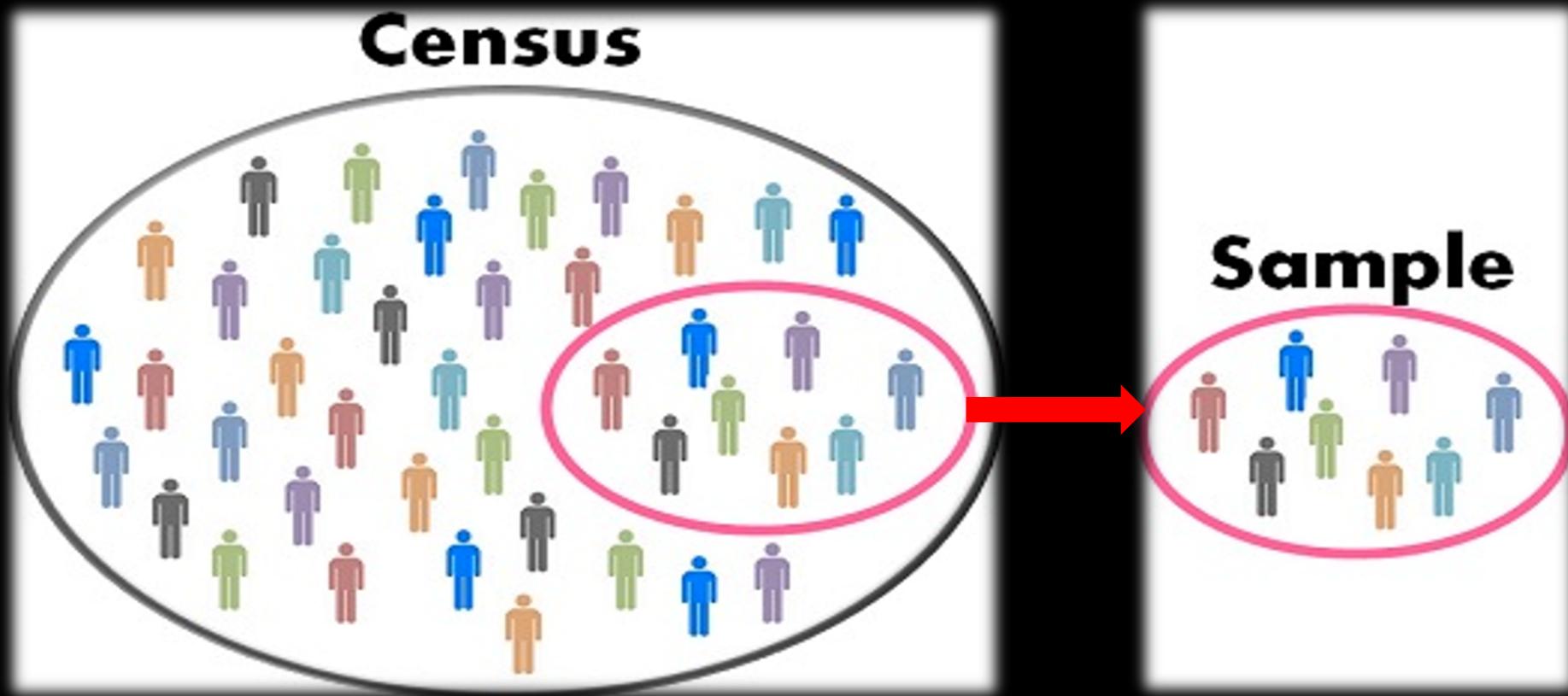
# BUSINESS PROBLEM

**Scenario:** We are a team two Digital Marketing Managers and two Data Analysts at Apple Inc.

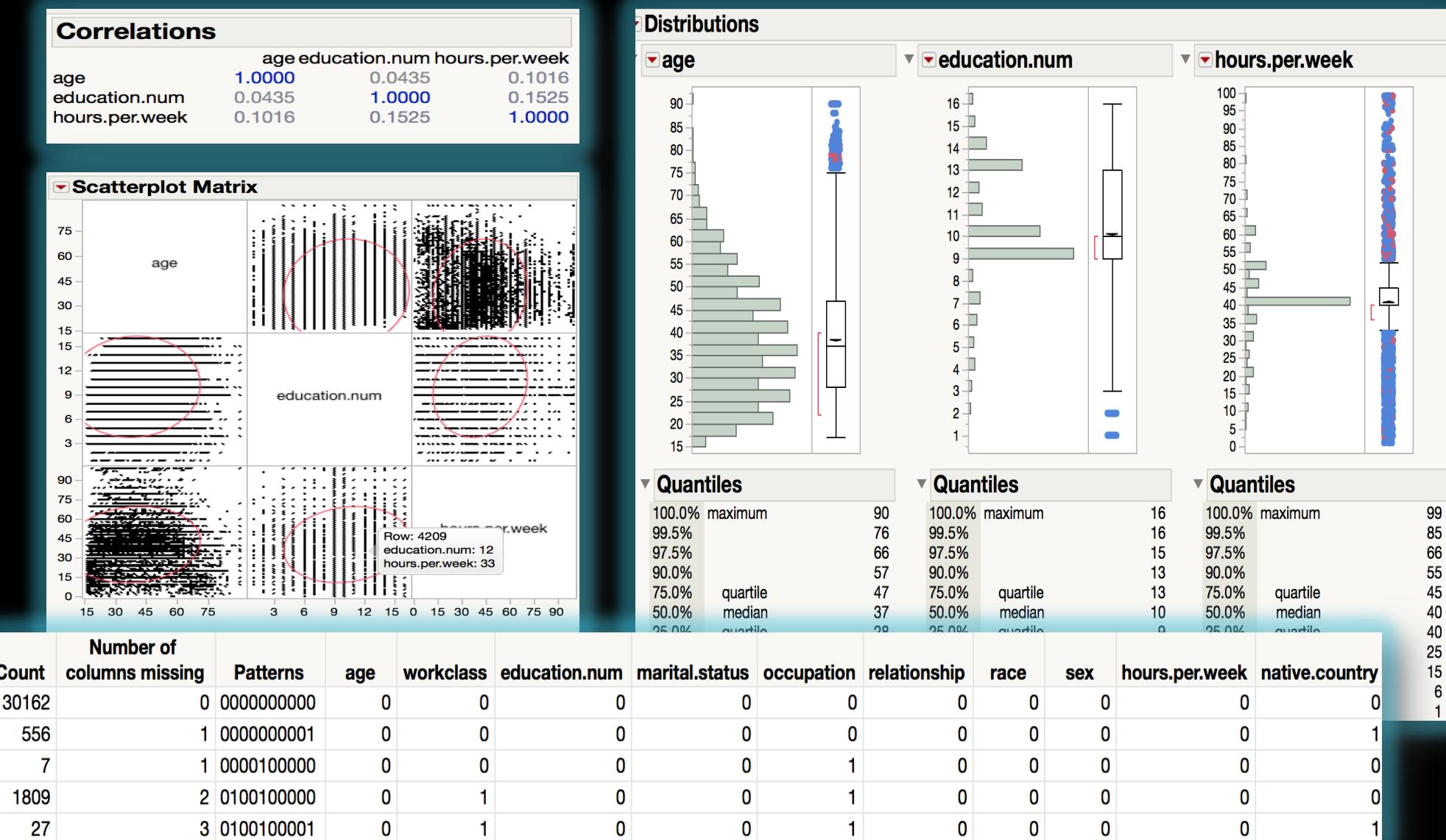
## **Problem to be solved:**

- Limited budget for advertising on social media channels
- Targeted segment: customers with annual income > 50K
- Findings to be reported to VP – Marketing (Jose Cruz), on what variables are the most important in reaching the targeted segment

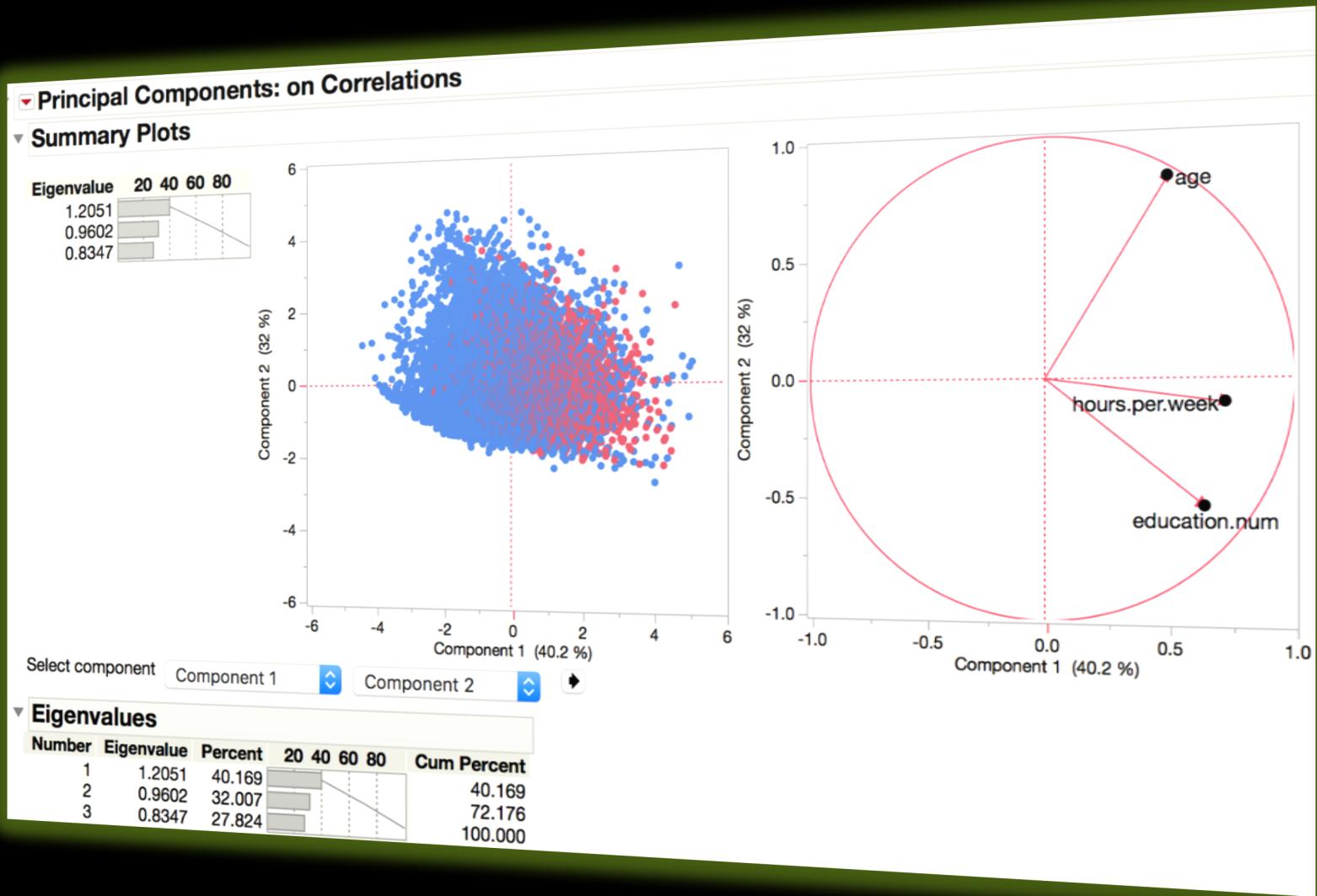
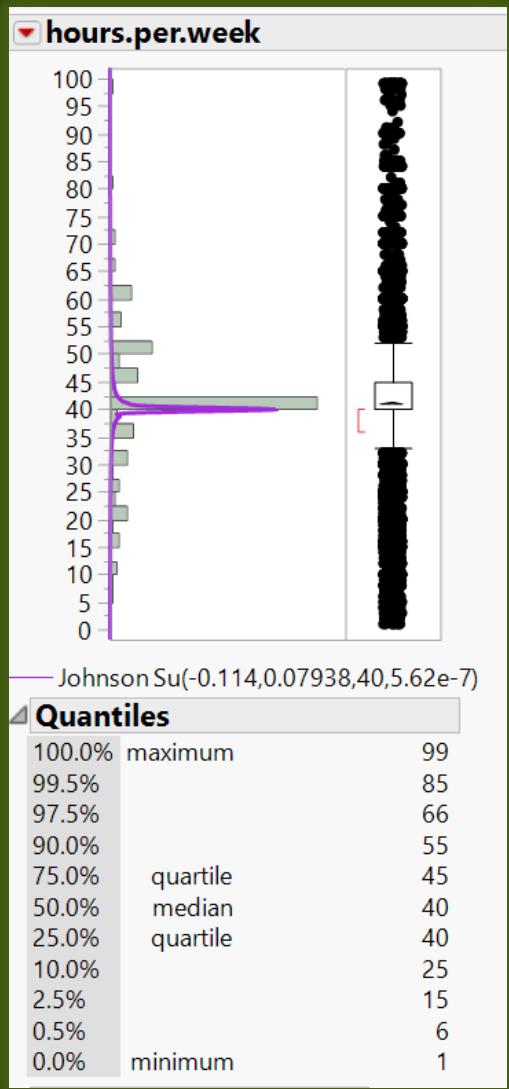
# SAMPLE



# EXPLORE



# MODIFY



# MODEL

- Logistic Regression
- Neural Network
- Partition Tree
- Boosted Tree
- Bootstrap Forest
- K-Nearest Neighbor
- Naïve Bayes

# ASSESS

**Model Comparison Validation=Training**

Predictors

Measures of Fit for income 2

Creator	.2	.4	.6	.8	Entropy RSquare	Generalized RSquare	Mean -Log p	RMSE	Mean Abs Dev	Misclassification Rate	N
Fit Nominal Logistic					0.3602	0.4930	0.359	0.3405	0.2316	0.1706	18097
Neural					0.3785	0.5132	0.3487	0.3355	0.2253	0.1625	18097
Partition					0.3750	0.5093	0.3508	0.3368	0.2272	0.1677	18097
Bootstrap Forest					0.3765	0.5109	0.3499	0.3362	0.2360	0.1644	18097
Boosted Tree					0.3525	0.4844	0.3634	0.3417	0.2455	0.1731	18097
K Nearest Neighbors					.	.	.	.	.	0.1477	18097
Naive Bayes					.	.	.	.	.	0.2069	18097

**Model Comparison Validation=Validation**

Predictors

Measures of Fit for income 2

Creator	.2	.4	.6	.8	Entropy RSquare	Generalized RSquare	Mean -Log p	RMSE	Mean Abs Dev	Misclassification Rate	N
Fit Nominal Logistic					0.3493	0.4808	0.3651	0.3375	0.2290	0.1662	12065
Neural					0.3698	0.5036	0.3536	0.3365	0.2254	0.1653	12065
Partition					0.3630	0.4961	0.3575	0.3382	0.2276	0.1673	12065
Bootstrap Forest					0.3691	0.5028	0.354	0.3370	0.2366	0.1615	12065
Boosted Tree					0.3680	0.5017	0.3546	0.3365	0.2423	0.1687	12065
K Nearest Neighbors					.	.	.	.	.	0.1728	12065
Naive Bayes					.	.	.	.	.	0.2027	12065

*Small variation in values between Training and Validation sets = no over fitting!*

*Low misclassification rate and higher R Square*

# BUSINESS VALUE ANALYSIS

**Base Scenario:** Marketing Campaign (Random Selection)

*Total Budget - \$100,000*

*Cost per person - \$1*

*Total reach – 100,000 people*

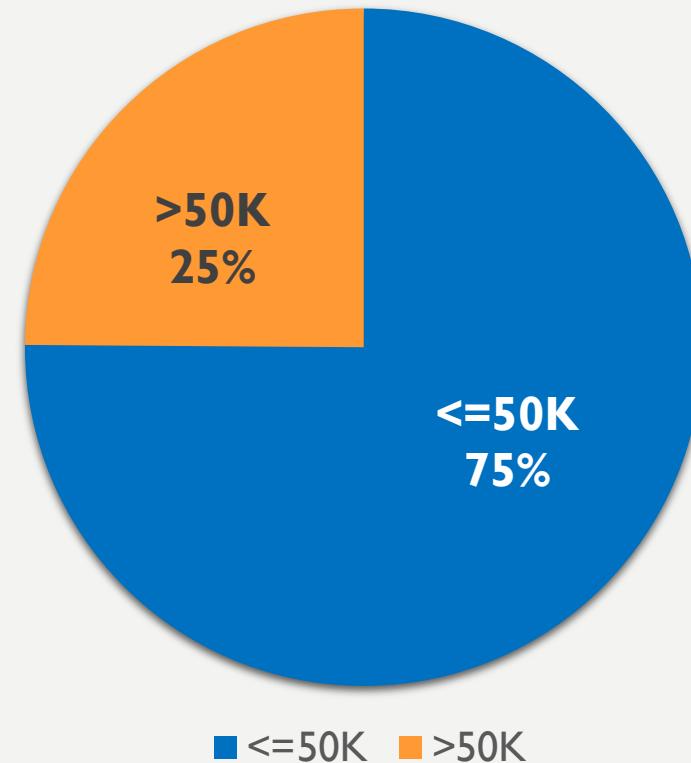
*% of selecting person with income \$ >50K = 7508 / 30162 = 24.9%*

*Revenue per customer - \$1000*

*Conversion rate – 1%*

*Expected Revenue - \$ 249,000*

*Income Distribution*



# BUSINESS VALUE ANALYSIS

Now, lets say we develop a model with a confusion Matrix like this

Actual Count	Predicted Count		
Income	0	1	Total
0	1261	1742	3003
1	452	8610	9062
Total	1713	10352	12065

# BUSINESS VALUE ANALYSIS

## With a new model

Total reach – 100,000 people

% of selecting person with income \$ >50K  
=  $1742 / 30162 = 83.2\%$

Revenue per customer - \$1000

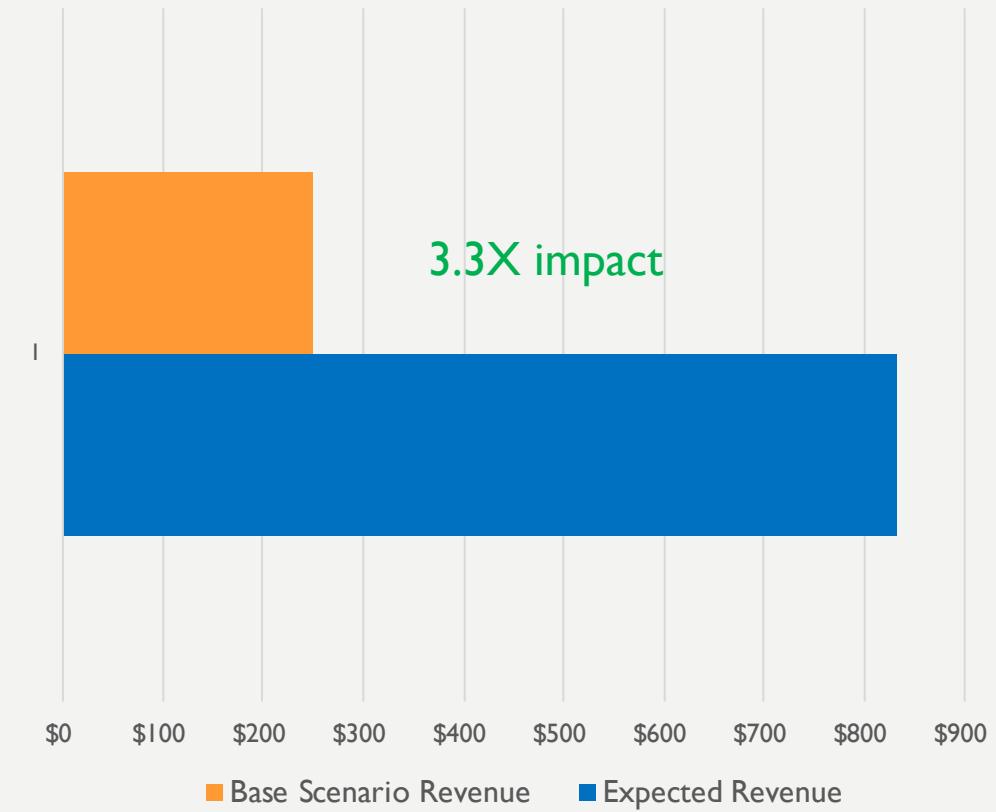
Conversion rate – 1%

Expected Revenue - **\$832,000**

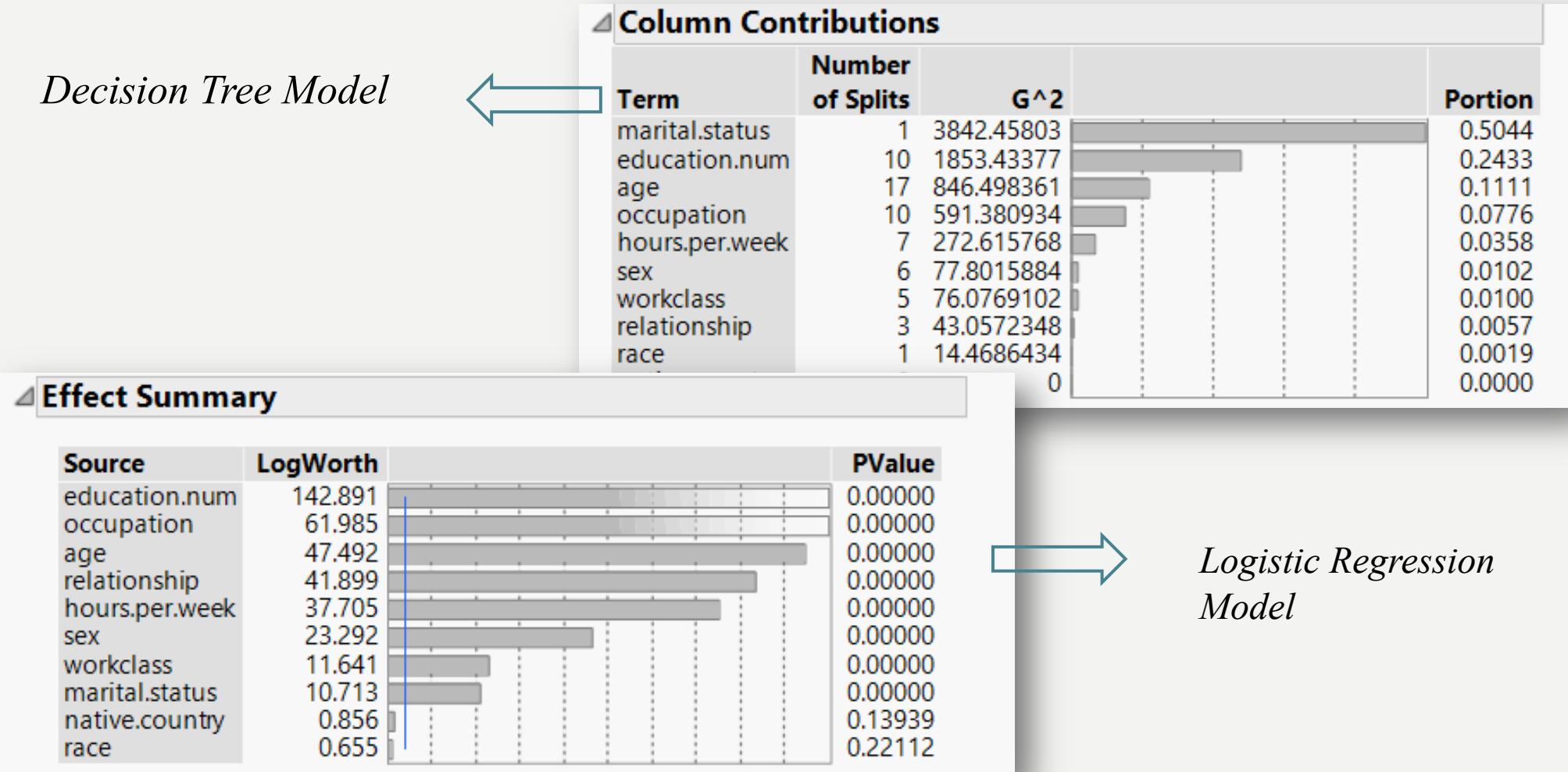
Base Scenario Revenue - **\$249,000**

Impact – **3.3X**

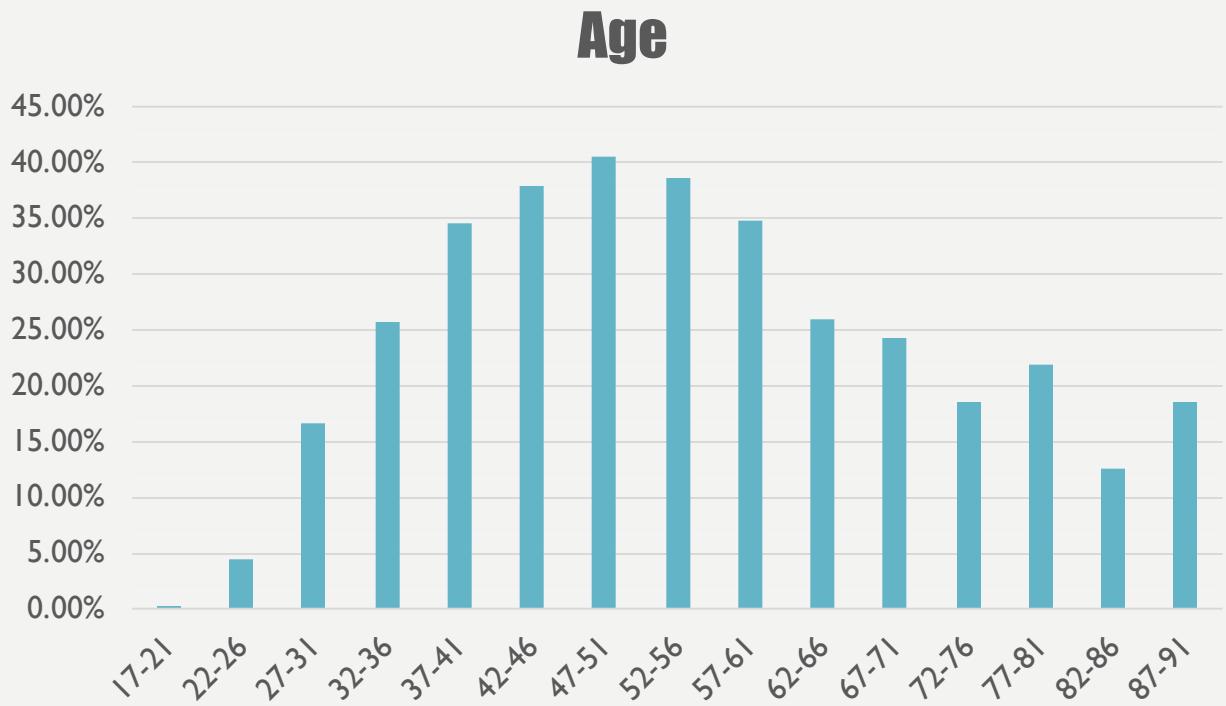
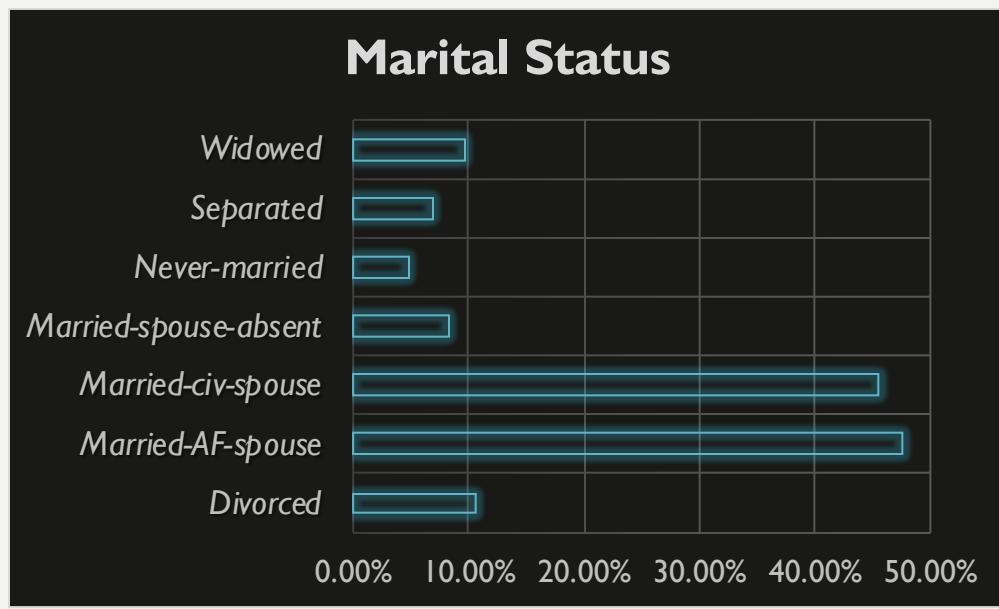
Revenue Comparison



# RECOMMENDATIONS

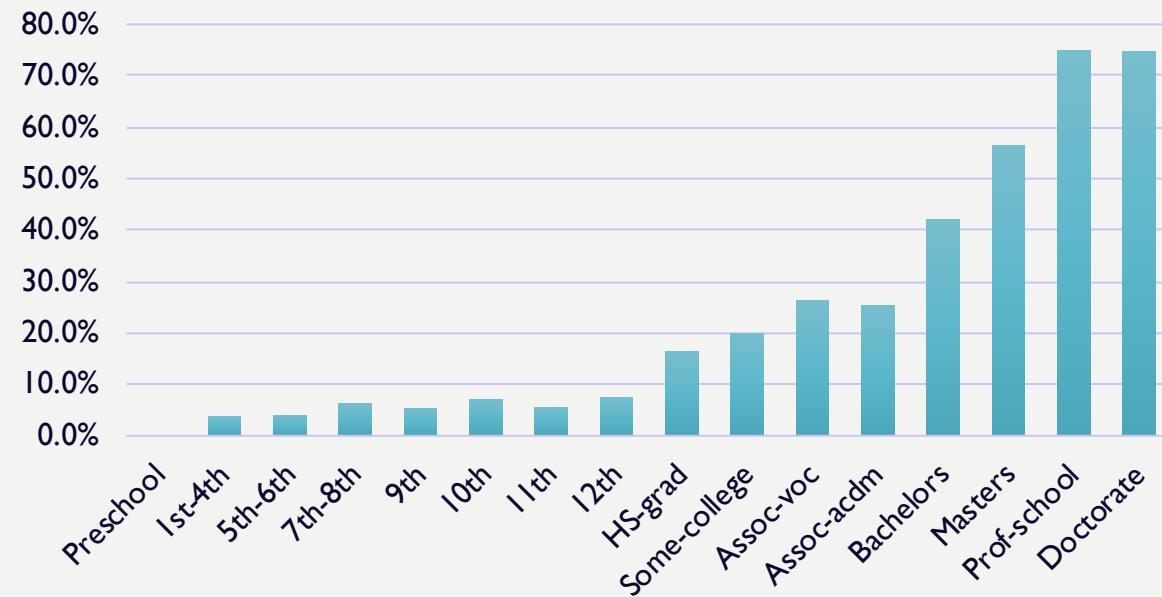


# RECOMMENDATIONS

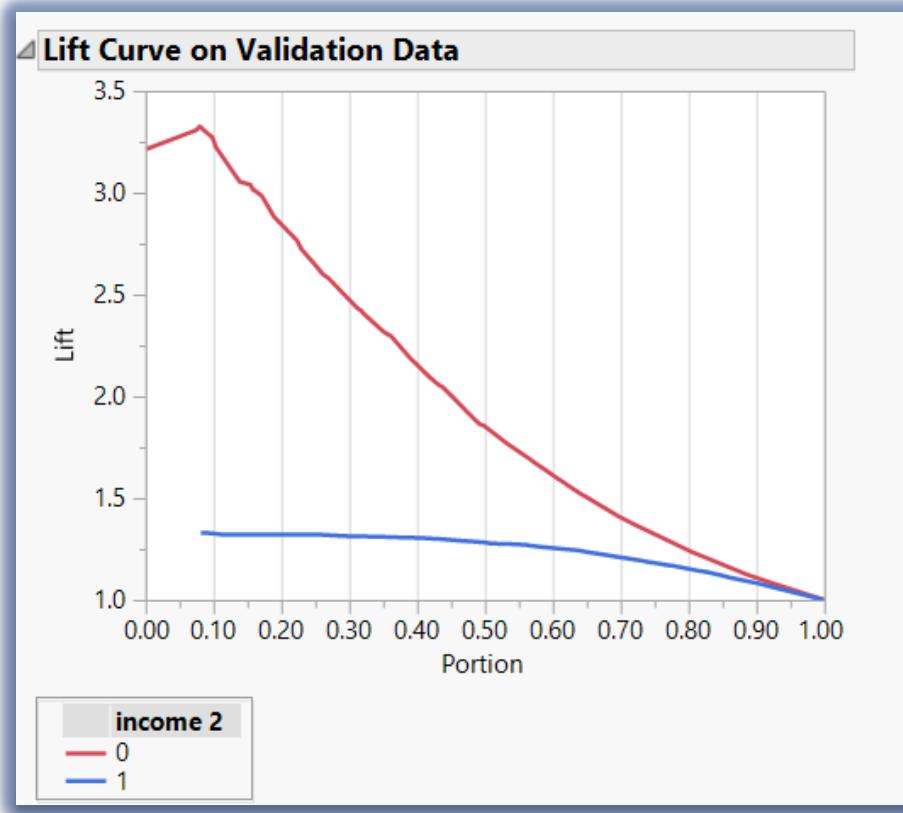


# RECOMMENDATIONS

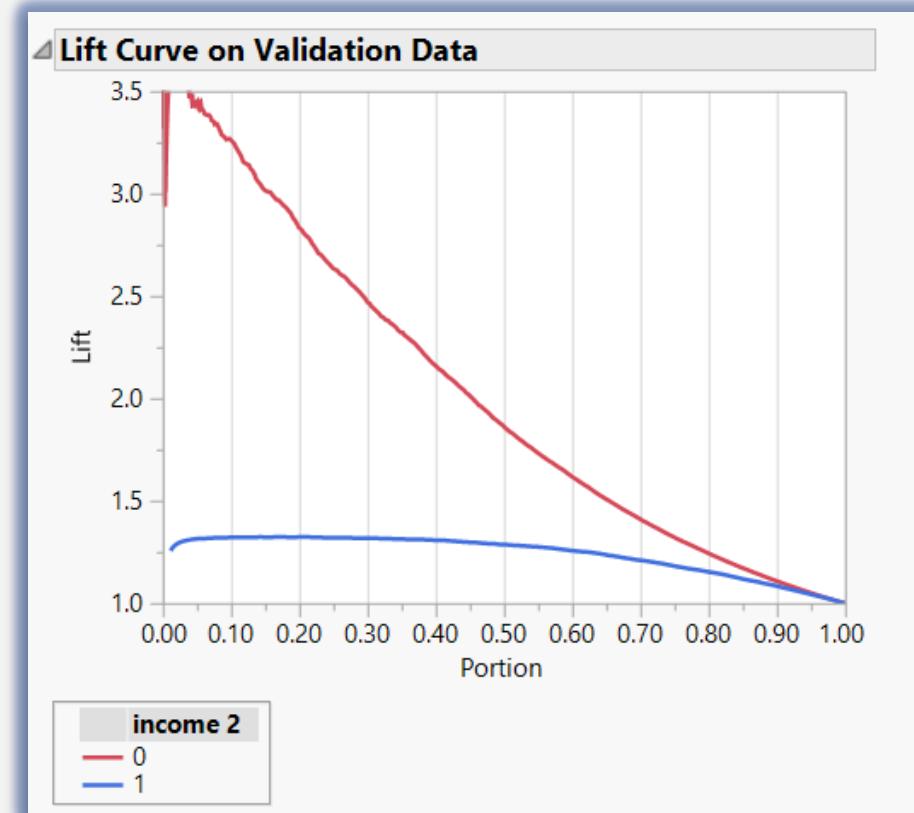
**Education Level/Number**



# RECOMMENDATIONS



*Decision Tree Model*



*Logistic Regression Model*

# RECOMMENDATIONS

Based on the models ran by our analysts, we can say that these four variables are important in targeting the desired customer segment through the marketing campaign

- Marital Status
- Education Number
- Occupation
- Age

*Thank You*