



1

Agenda

- Background
- Introduction and goal
- Data source and description
- Key challenges
- Descriptive analysis(tableau)
- Predictive analysis(SAS EM)
- Key Learnings

2

1

Background of the blight property

Blight: abandoned and/or vacant property, either building or home, that is in an unacceptable condition of disrepair. It is inhabitable and has little to no economic value.

Background:

- US is facing the problems of empty neighborhoods across the country. We can see similar problems in Memphis, TN, and the causes for this are: Poor planning, Poverty and Sub Urbanization are the causes of the Urban Blight. The problems are caused because of – Unemployment, depopulation problems, promote crime, loss of deindustrialization. Blight affected areas are characterized by dirty, cloistered, tumbledown buildings and inhospitable living.

Effects of such properties

- Danger to the public
- Increase in negative incidents
- Negative effect on surrounding property values
- Reduces property tax revenue

3

Business objective and project goal

• Business objective:

Identify properties that are possible to become blighted. This information can help local community and government plan in advance and minimize the bad effects of a blighted property/area.

• Project goal

Identify the data that are more related to to blight property, characterize their relationships and predict the possibility that a property turn into a blight property

• Approaches

- Use demolition fees as an indicator of blight properties.

- Explore relevant data (i.e., crime, sales, blight rating, property tax etc.) with tableau, and try to find the data related to blight property and explore their relationships.

- Explore the data with SAS EM to further characterize the relationships between the relevant data and the target and predict the possibilities.

4

Data source and description

Datasets	Content of each data set
Windshield	The thorough housing information , condition details of all the Memphis and Shelby property, and Memphis and Shelby county regulation and code enforcement, ~150,000 parcel ID. The data set was not cleaned
Tax aggregation	32 statement records of aggregated (sum) of property tax information over the period of 2010-2016, 46,751 parcel ID, data collected from 2010 to 2016.
MLGW aggregation	count how many times MLGW cut off the supply to one property. Interval (discrete) data with data range [0, 2], two types of data recorded (gas cut and electricity cut) per year for 15,678 parcel ID
Sales aggregation	count how many times sale happens each year. Interval (discrete) data with data range [0, 27], one sale record of each year for 127,609 parcel ID
Crime aggregation	count how many crimes happen on each property. Interval (discrete) data with data range [0, 447], three types of crimes(violent, nonviolent, other) document for each year; 98,366 data point

- Data processing approach and flow for further analysis:
 - Randomly sampling parcel IDs in the tax aggregation data and compare the tax value with the information on the Shelby county trustee website
 - Select the tax aggregation data from a year that has the most data points documented
 - Clean up all types of errors in the Windshield data
 - Convert Parcel IDs in the above five datasets into the same format.
 - Inner joined all the dataset and create a dummy variable for demolition fees.
- Target variable: demolition dummy variable
- Input variables for SAS EM:
 - 1st round, all the variables from the above data set except for unary data, ID, zip code and data missing over 50%
 - 2nd and the later round, all the variables that have the P value lower than 0.6 in the regression model statistics

5

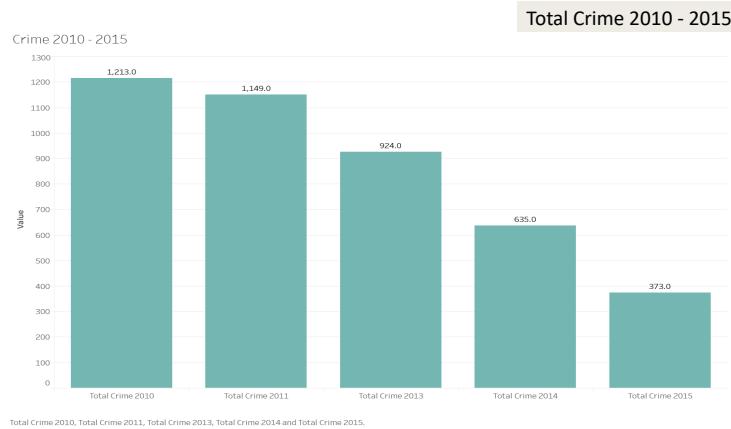
Key Challenges

- Data challenges:
 - Data cleaning, the windshield data provide useful details of housing conditions but it is very messy. It took about 8 hours to clean it to a useful/informative level.
 - The inaccuracy of some given datasets requires to examine the data using higher authority data source, which was not expected when the project started.
 - Need domain knowledge of housing, and Shelby county regulation to understand the meaning and value of related variables and manipulate the input without biasing the modeling process.
 - Too many variables for EM regression models and missing values in a high dimension dataset
- Data techniques:
 - Descriptive analysis with Tableau
 - Joint or merge datasets, Query and computing in SAS EG
 - Decision trees, regression and neural network modeling with SAS EM

6

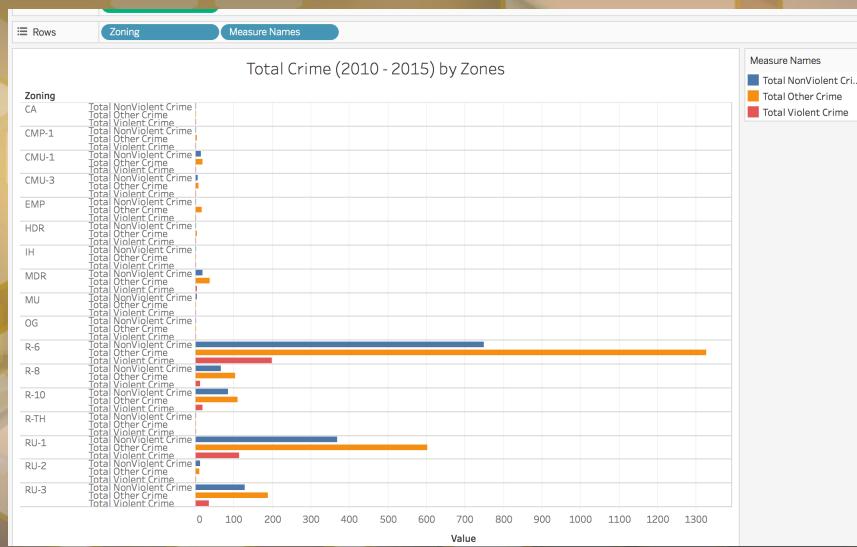
Descriptive Analytics

- Analyze datasets to get basic understanding of simple statistics
- Define the most affected zoning districts
- Analyze and explore relationship between Crime, Sales, Gas/Power cuts, property tax and demolition fees.



7

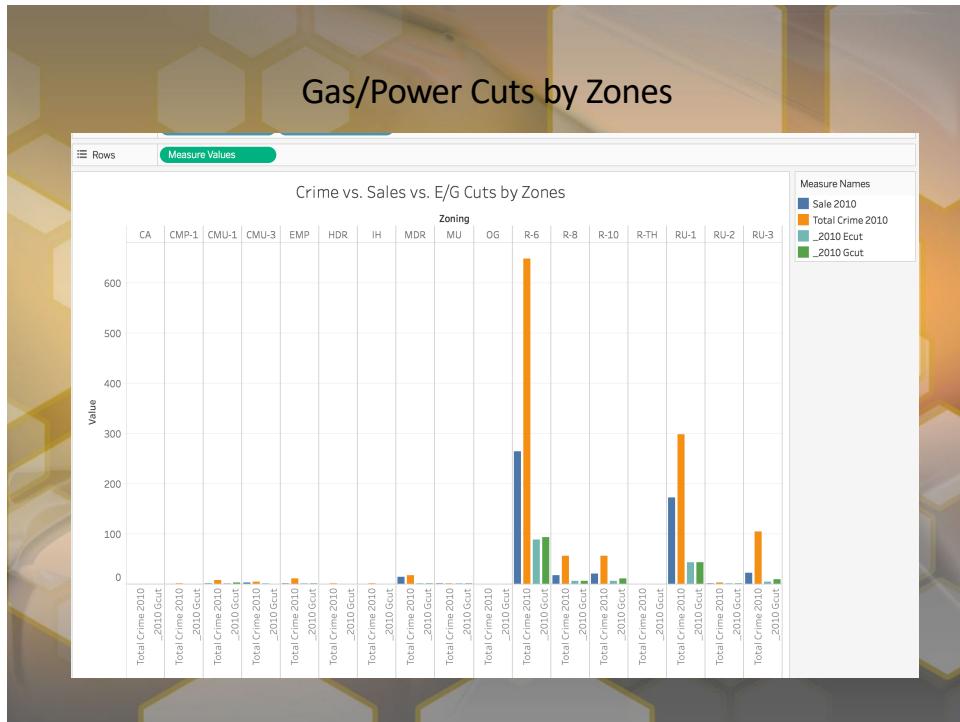
Violent, Non Violent, Other Crime by Zones



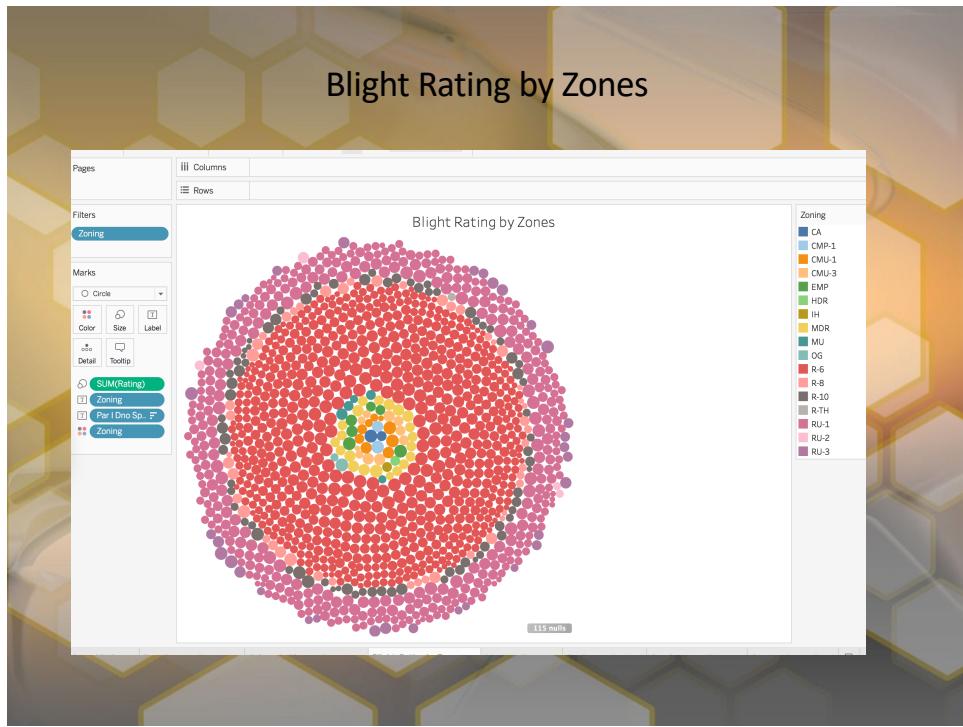
8



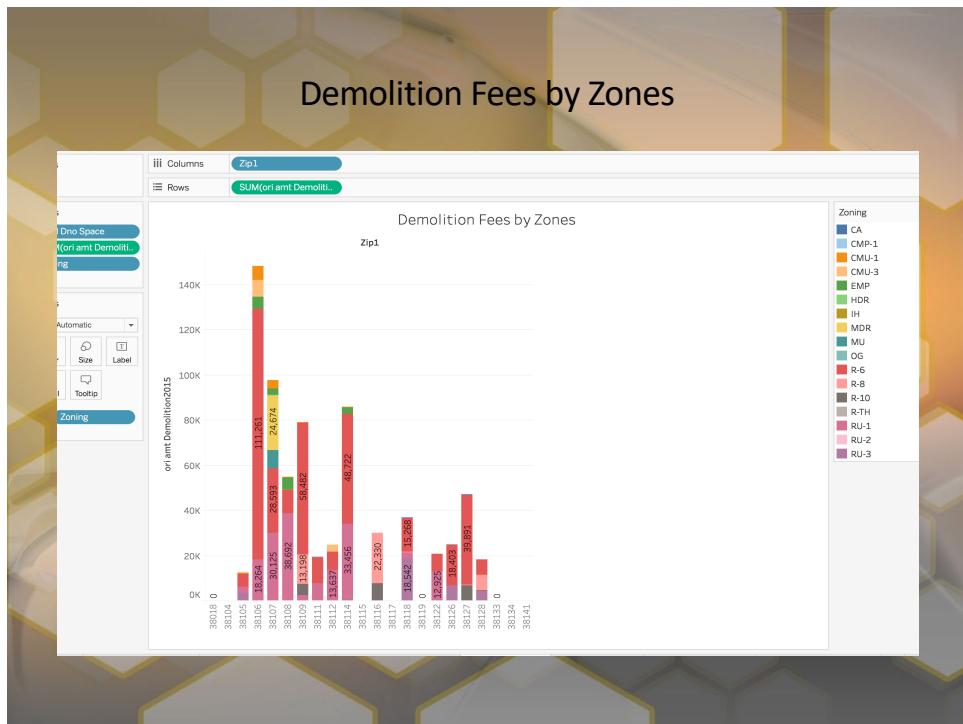
9



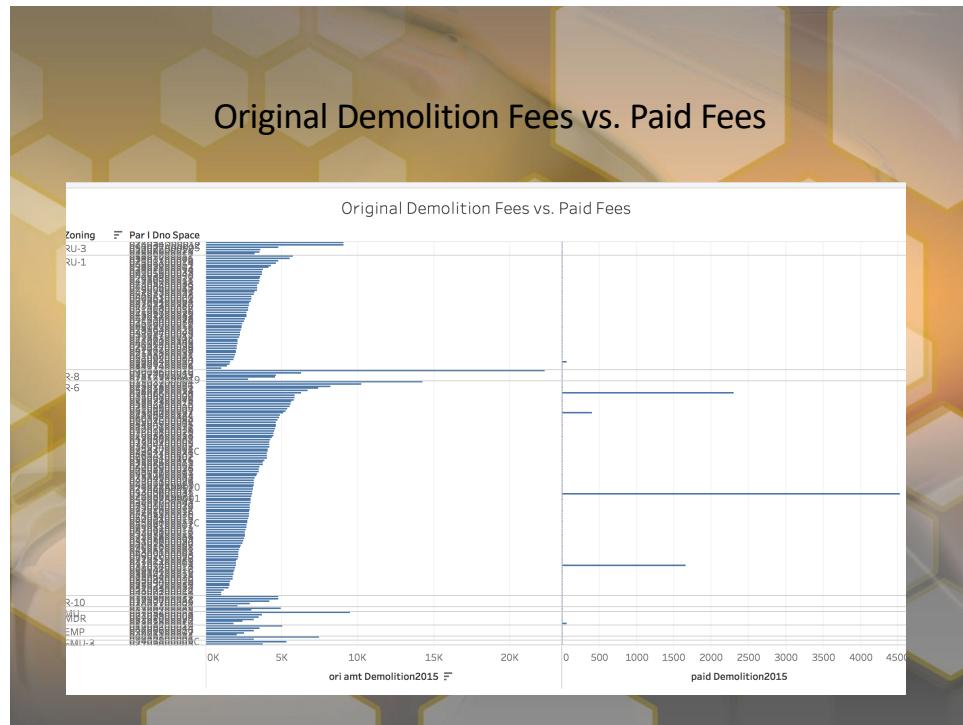
10



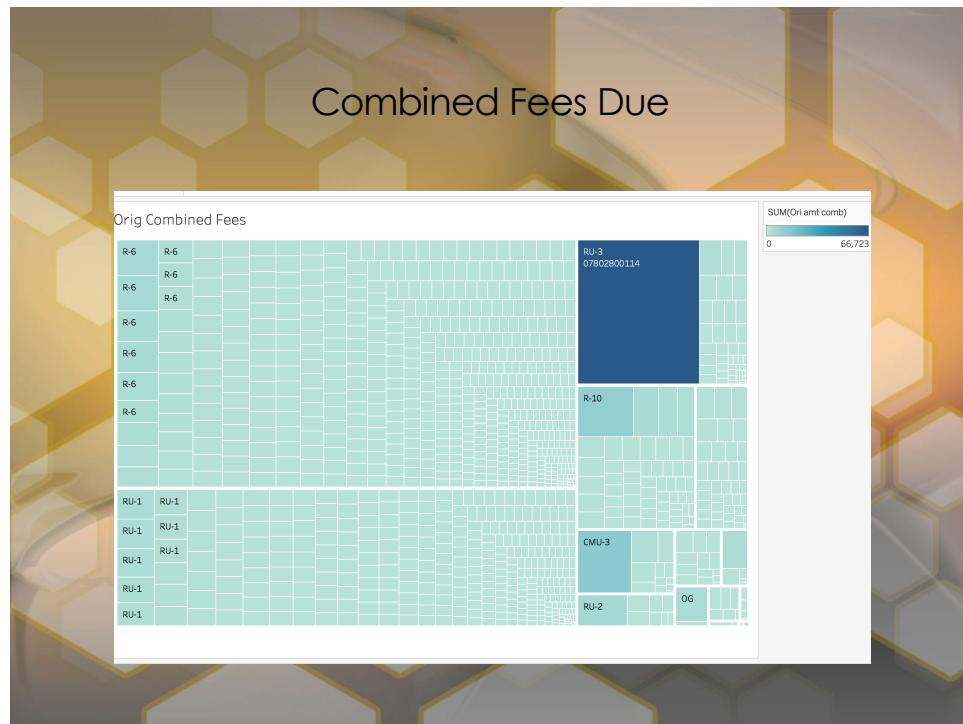
11



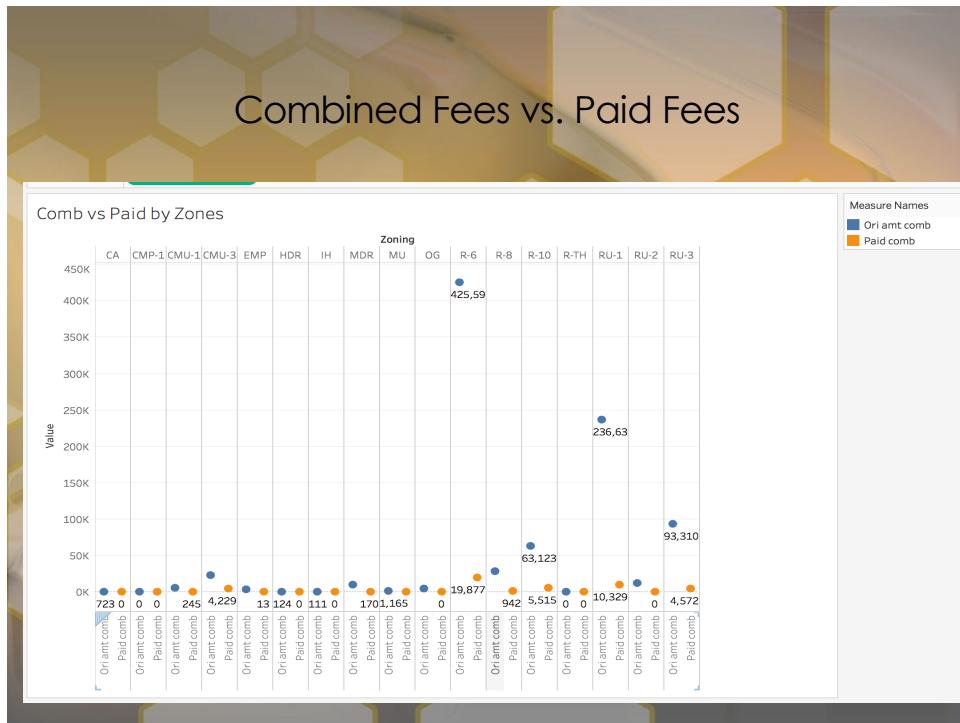
12



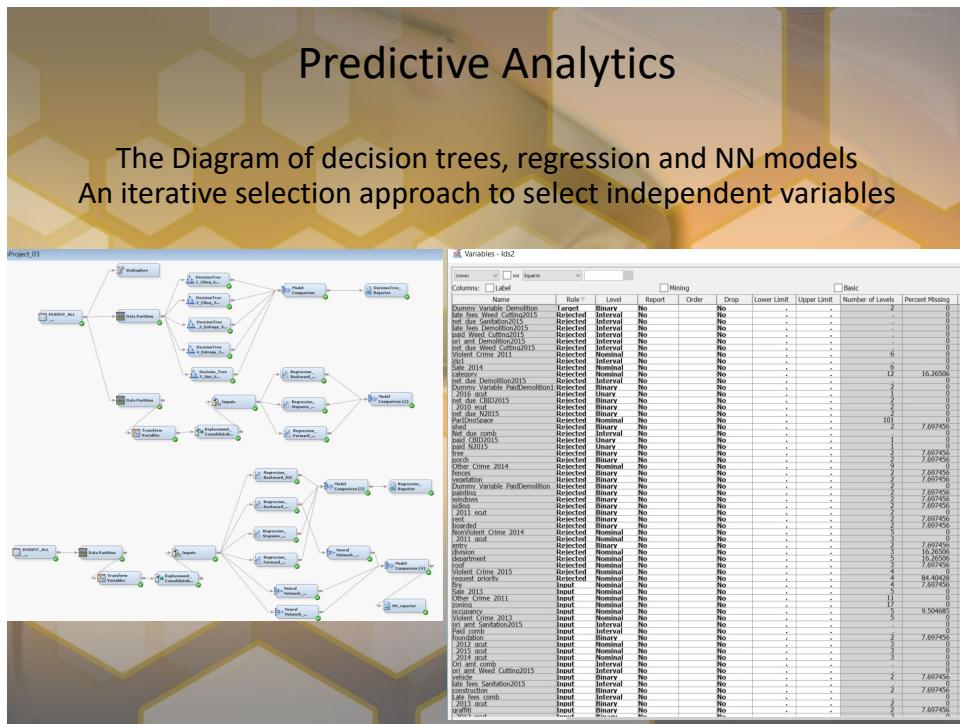
13



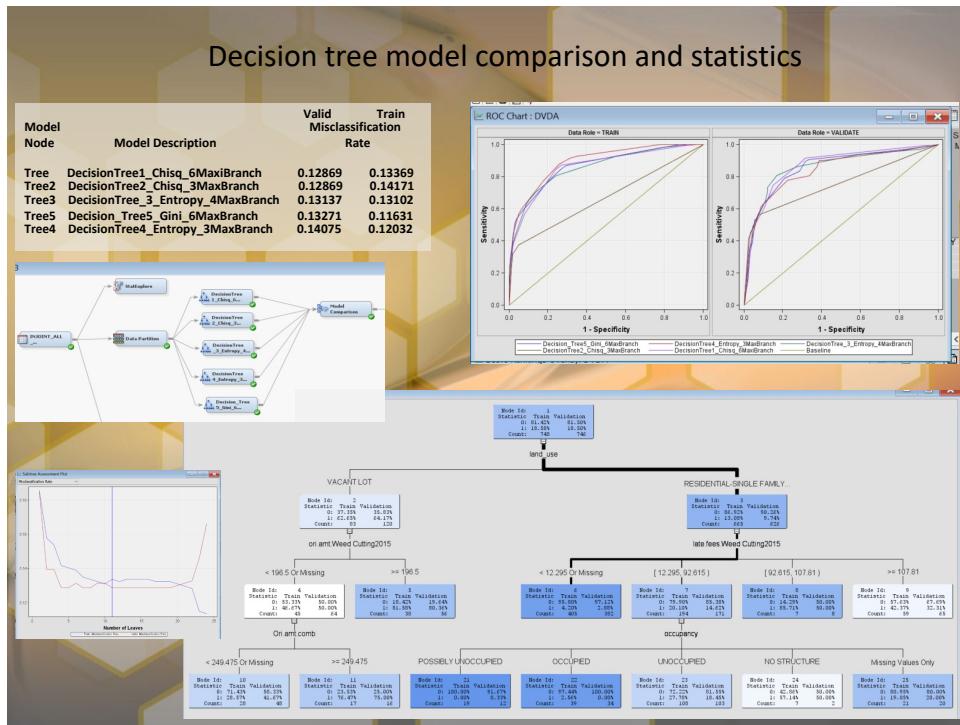
14



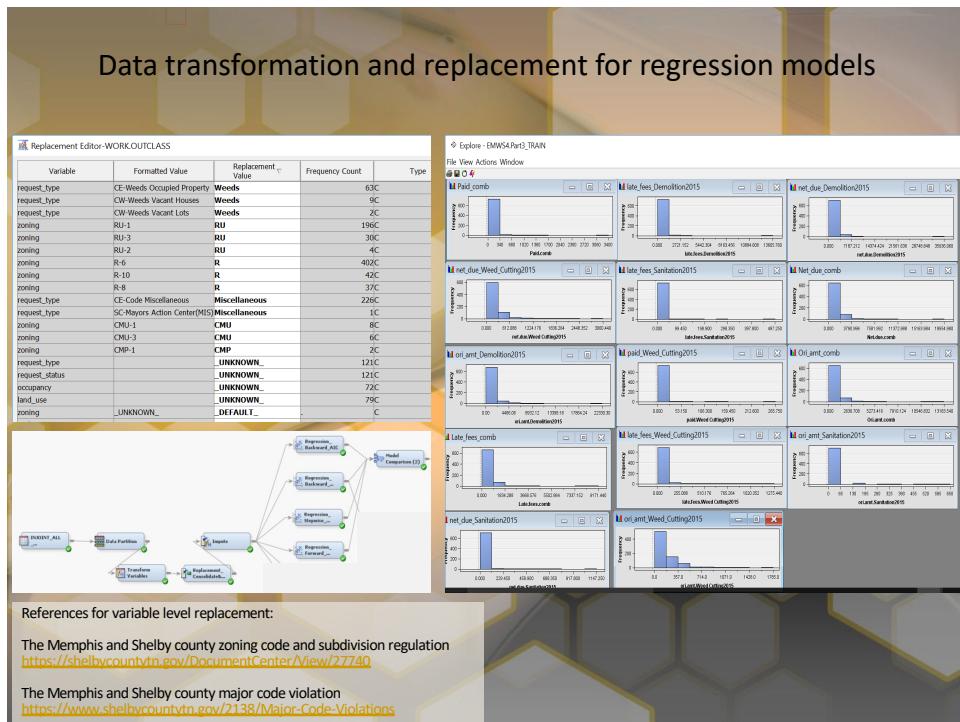
15



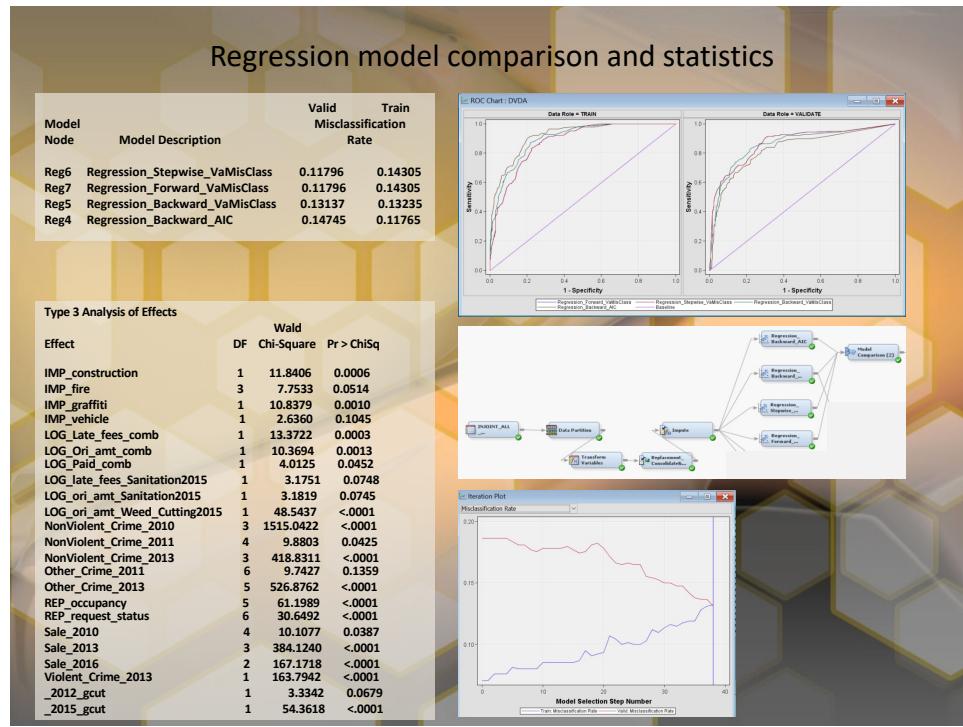
16



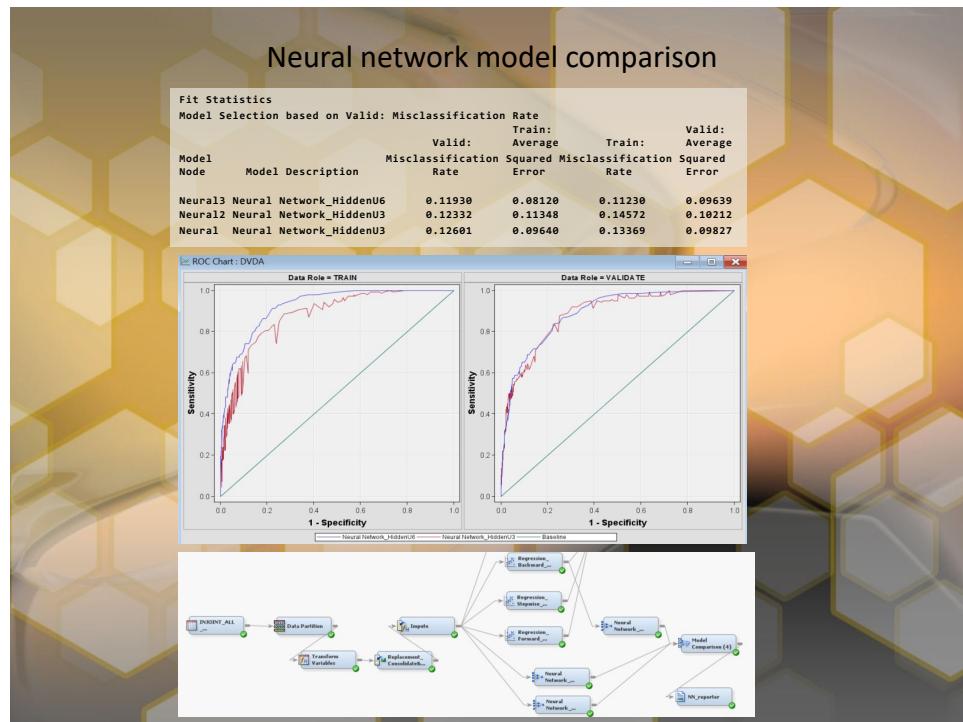
17



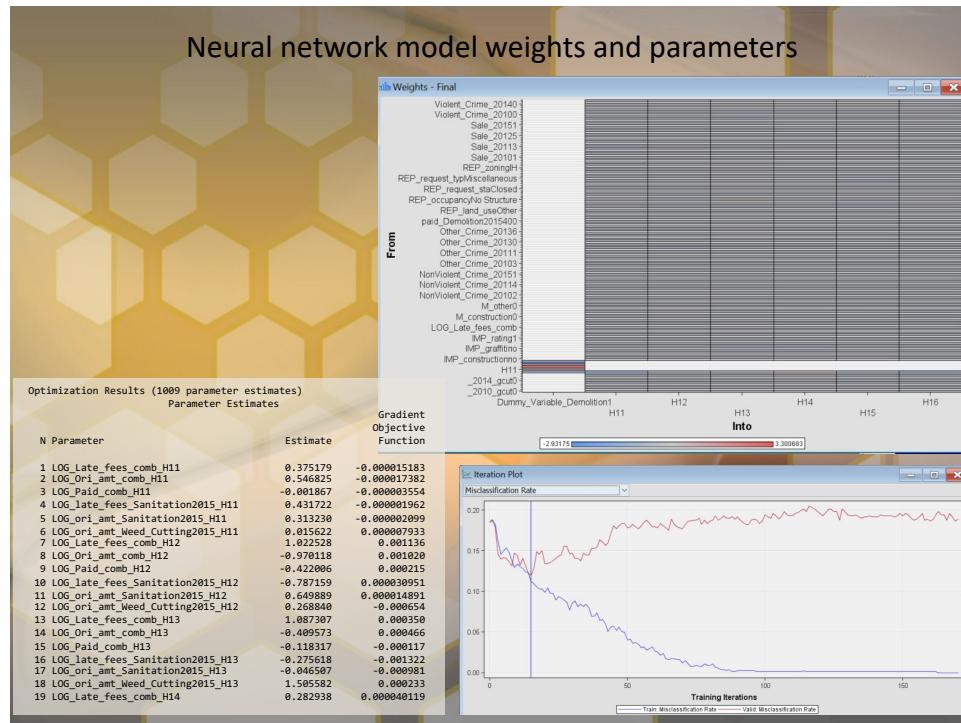
19



19



20



21

Key learning

- Never assume the datasets are accurate or clean from any source without examining them.
- The Identification of the related variables and reduction to the unnecessary data complexity will help the SAS EM perform better---- it is not a good idea to shovel the data into EM and rely on it to perform well.
- The output log file provides a thorough documentation with statistic calculation, model parameters and iteration steps, which are helpful to trouble-shooting model optimization. Therefore, it is worth to learn how to read the file and understand certain amount of those statistics .

22



23