

# Stat 450 Project Report: Employment of the Chinese Government Agencies

*Harry Xu*

*2018-03-06*

## Summary

This project aims to better understand the distribution and trends of employment for two of China's civil government branches, Taxation Administration Services (TAS) and Environmental Protection Agencies (EPA). This report focuses on the visualization and a panel regression analysis using the TAS data. The model suggests that the significant (5% level) variables that could possibly explain the variation in number of staff per province over time are: Public Expenditure and GDP per Capita. The EPA data analysis can be done in the similar manner.

## Introduction

This report presents a graphical visualization and statistical analysis of employment levels for two of China's civil government branches, Taxation Administration Services (TAS) and Environmental Protection Agencies (EPA). The number of staff for each province were collected at each government, age, and education levels from years 1992 to 2013. The main objectives are to first generate visualizations to explore patterns and trends within government staffing based on different social and economic factor levels, and secondly, to build a statistical model to explain the variation in staffing numbers at various levels of the government. The graphs in this report are generated by Tableau Desktop 10.5, and the panel regression methods is used thoroughly in this the analysis to model the data. This analysis only focuses on the TAS dataset, in particular, gives an example visualization and analysis for different age levels. The other levels, such as education and department unit levels, and the analysis for analysis can follow a similar analysis structure, which is not be shown in this report.

In addition, due to the limitation of the data provided and the possibility that the assumptions for the panel regression model may not hold, the result of this analysis can be unsatisfactory. It is advised that the reader should take great attention to the conclusion and recommendation section of the report for clarification due to these limitations and follows the guide to future steps.

The section layout for the report is as follows:

- Data Description
- Methods
- Results
- Conclusions and Recommendations (Future Steps)
- References
- Appendix
- Figures

## Data Description

Two collections of employment data for the two civil service braches (TAS and EPA) were obtained. These panel data include staff numbers for each of the 30 provinces, categorized by different bureau, age, education and department unit levels. The TAS dataset spans 8 years (1996, 2000-2007) and the EPA dataset spans 18 years (from 1992 to 2011, missing year 1994). The third dataset, namely Econ-demographic, has 10 different Chinese provincial economic variables for each province spanning from 2000 to 2013.

These datasets were cleaned and its structures simplified to ensure a balanced panel data, to allow computer software to generate visualization and panel regression analysis. The table below illustrates the structure of the TAS dataset (the EPA dataset are similar, and will not be shown as the analysis is focused mainly on the TAS dataset.)

Dataset TAS:

Variable	Levels
Age	8 levels (<30, 31-35, 36-40, 41-45, 46-50, 51-54, 55-59, >60)
Edu	7 levels (Junior, High, Technical, Other Post, Bachelor, Post-Grad)
Unit	4 levels (Admin, DA Admin, Tax, DA non-Admin)
Bureau	2 levels (STB and LTB)
Years	8 years (1996, 2000-2007)
Response Variable	Staff Number
Number of Provinces	30

Dataset Econ-demographic: contains the 10 different economic measures for each province per year, these measures include: Urban Area, Number of above size industries, public expenditure, GDP, Population, Ratio of Urban Area, GDP per Capita, and Area per Square Km. The structure of the dataset is as follows:

Variable	Levels
Years	14 years (2000-2013)
Response Variable	Staff Number
Number of Province	30

## Methods

The main statistical method used for the analysis of TAS staffing data is the two-way fixed-effect panel regression. A panel regression is a regression method used to analyse two dimensional data. The data is usually collected over a time period for the same entities. In this report, the Provinces would be the entities. Such two-way fixed-effect panel regression model used can be expressed mathematically as:

$$Y_{it} = \beta_1 X_{it} + \alpha_i + \lambda_t + \epsilon_{it}$$

where  $Y_{it}$  is the response variable representing the number of staff for a specific province  $i$  in a specific year  $t$ . The matrix  $X_{it}$  are the possible economic explanatory variables that could explain the response variables, and  $\beta_1$  are the slopes for the explanatory variables of interest for each  $X_{it}$ .  $\alpha_i$  is the entity (Province) fixed effect, and  $\lambda_t$  is the time fixed effect. Both terms have a specific effect on  $Y_{it}$ . The Provincial effect is the effect at a specific year on the mean number of staff for that specific province. The time effect is the effect at a specific year. It is assumed that the Province effect is fixed, meaning that the mean number of staff for a province does not vary over time, but it might be different from other provinces. The time effect is fixed for provinces, and it indicates that the mean number of staff might be different for different years, but this

staffing level difference is the same across all provinces. In this model, time  $t$  takes value 1996, 2000 to 2013, and  $\epsilon_{it}$  is the error term for Province  $i$  in year  $t$ . It is assumed that the errors are iid and  $\epsilon_{it} \sim N(0, \sigma^2)$

## Results

From the visualizations shown in Figure-1, each age level is plotted against 3 different economic variables: GDP per Capita, number of above size industrial companies, and public expenditure. It is interesting to see that the trend over the years are different at each age level, the test below will show whether there is a significance to these variables at each age level. Note that the variable GDP is excluded as it has strong correlation with GDP per Capita to avoid redundancy. Similar plots can also be generated for the TAS dataset at different education level and office unit level, which are not shown in this report.

Figure-3 shows the plot of staff numbers at each age level in High/Med/Low GDP per Capita groupings. The groupings are grouped manually based on the top plot of Figure-2. The highest 3 Provinces with GDP per Capita are the “High” group, medium 6 provinces are the “Med” group, and the rest are the “Low” group. This plot gives an over all trend within the province groupings which can provide some interesting story to tell.

Figure-4 shows the distribution for the number of staff by age in percentage. This plot provides useful information for comparisons between groups, because it depicts the change of staff numbers per level with a same percentage scale.

For simplification purposes, only the regression output for TAS dataset with Age level less than 30 is shown below. The rest of the results are shown in the Figures section. The two-way fixed-effect panel regression model incorporated the independent variables: urban area, number of above size industrial companies, public expenditure, and GDP per Capita.

The Appendix 1 R output gives the slope estimates for each  $\beta$  with a p-value. The p-values indicated the significance of the slope for each of the explanatory variables. It shows that Urban area, GDP and GDP per Capita are significant covariates at 5% in terms of explaining the variations for the number of staff who 30 years old or less over the span of years. However, the adjusted  $R^2$  value is really low, which suggests that this model is probably not a great fit for all the explanatory variables chosen as the input values for the model and the assumptions in the model are not satisfactory. Nonetheless, the results might provide a useful direction for variable selection in the future.

Figure-5 shows the summary result for all Age levels in the TAS dataset. It gives the significant variables at 5% at each age level separated by bureau that may help explain the variations for the number of staff in TAS.

## Conclusions and Recommendations(Future Steps)

In summary, we are able to conclude from the TAS dataset regarding different employee age levels, Urban area and GDP per Capita are the significant variables for explaining the variation of staff numbers across provinces over the years at a 5% significant level. This method, along with the visualization tools, can also be applied to other datasets with factors such as the education level and department unit levels, to select the important explanatory variables for inference. There are many ways we could possibly explore and visualize and analyze due to the complexity of the dataset.

The following section explain the inherited problems analysis has:

1. One of the fundamental flaw of the this panel regression analysis pointed out by Professor Gabriela Cohen is that there are no way of measuring the variability of number of staff per province, since it is not possible to produce replicated measures for each province.
2. The abnormally low adjusted  $R^2$  value indicates a poor fit of the model.
3. Many parts of the given data contain imputed dataset in which the method of imputation was not explained and proved correct.

4. Due to the nature of the data provided, we are not able to test the dependency of the factor variables between levels (for example, between different ages). Although a correlation test between these factor levels can be done, but the conclusion is hard to interpret.

For future recommendations of panel regression analysis, a random effect for time could be introduced into the model to relax the fixed effect assumptions for time. We can also run multiple diagnostic test for panel regression to ensure the proper type of panel regression model to fit with these assumptions.

Although missing values in the dataset also has potential problems for the analysis, we are advised by the client to look further into the dimensions of the dependent variables itself and explore more about the variations within these 30 provinces. Therefore, we would likely to take the data as is and perform unsupervised learning algorithms to gain more insights. One possible way to achieve the goal is to group the provinces based on their staff number's behaviour over time. For each province's behaviour, a spline can be fitted, and then clustering analysis could be applied to suggest groupings of these splines that possibly indicate similar behaviours.

## References

link to the reference for panel regression documentation: <https://cran.r-project.org/web/packages/plm/plm.pdf>

## Appendix

### Appendix 1: R output for Panel Regression Analysis

```
## Loading required package: Formula

## Twoways effects Within Model
##
## Call:
## plm(formula = age.30.or..below ~ urbanarea + No.of.above.size.industrial.companies +
##      public.expenditure + GDPperCapita, data = agemerged, effect = "twoway",
##      model = "within", index = c("Year", "Province"))
##
## Balanced Panel: n = 8, T = 28, N = 224
##
## Residuals:
##      Min. 1st Qu.  Median 3rd Qu.    Max.
## -2250.0  -434.0   -32.5   291.0  9030.0
##
## Coefficients:
##                                Estimate Std. Error t-value
## urbanarea                    -0.526605    0.592002  -0.8895
## No.of.above.size.industrial.companies -0.060528    0.031908  -1.8970
## public.expenditure            -2.286941    0.690679  -3.3111
## GDPperCapita                   0.112941    0.024492   4.6113
##                                Pr(>|t|)
## urbanarea                     0.374872
## No.of.above.size.industrial.companies 0.059387 .
## public.expenditure             0.001117 **
## GDPperCapita                   7.449e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Total Sum of Squares:      226870000
```

```
## Residual Sum of Squares: 172540000
## R-Squared:      0.23949
## Adj. R-Squared: 0.083275
## F-statistic: 14.5643 on 4 and 185 DF, p-value: 2.3276e-10
```

## Appendix 2:

Git repository: <https://github.com/wenzhengzzz/STAT450-550-Project-GovtAgencies/tree/master/stat450>

## Figures

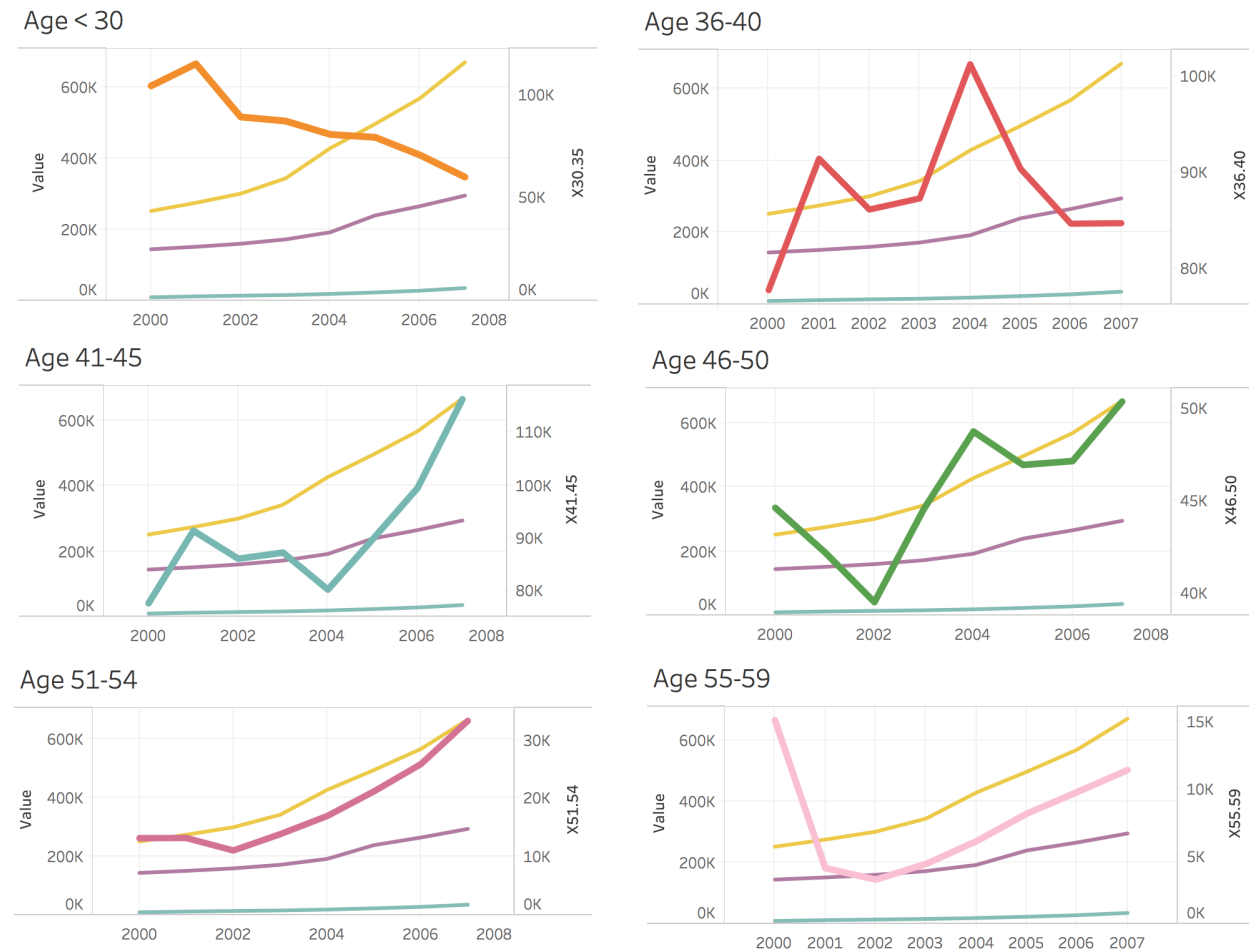


Figure 1: Age with Econ Variables

## Econ Grouping Ideas

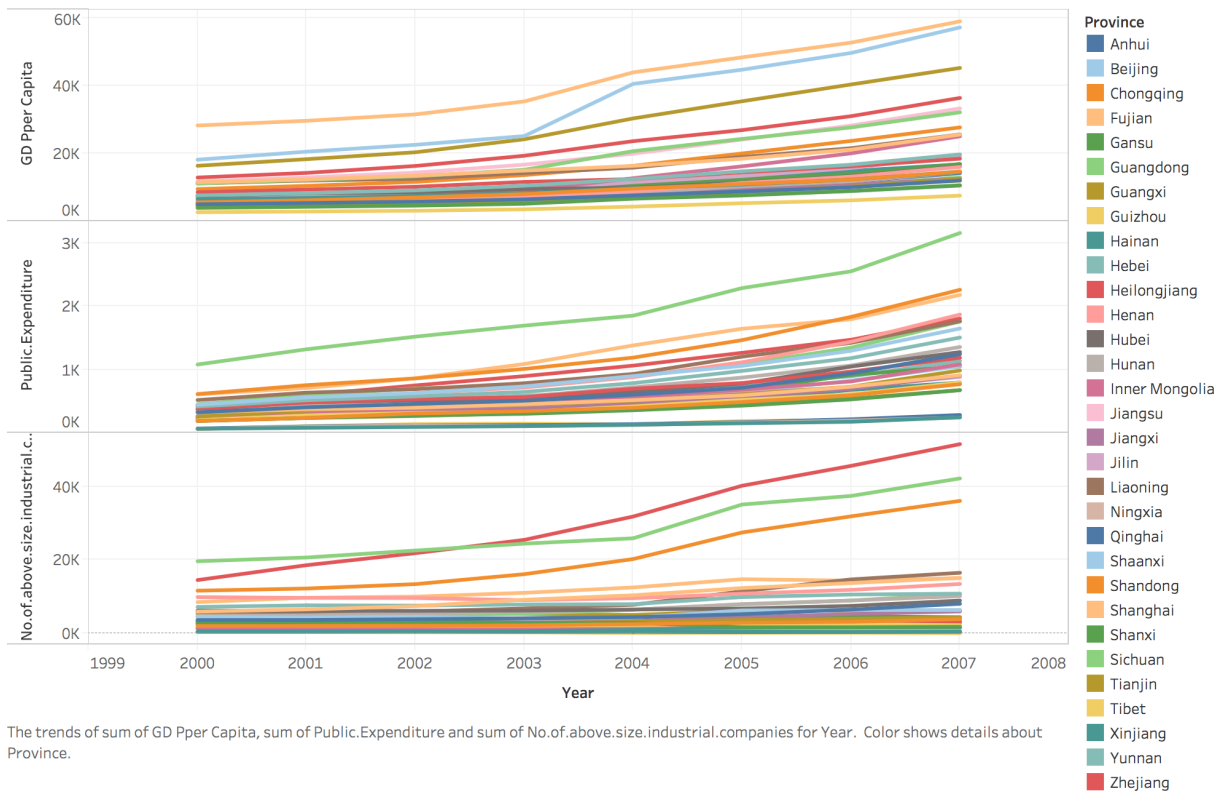


Figure 2: Econ Grouping Ideas

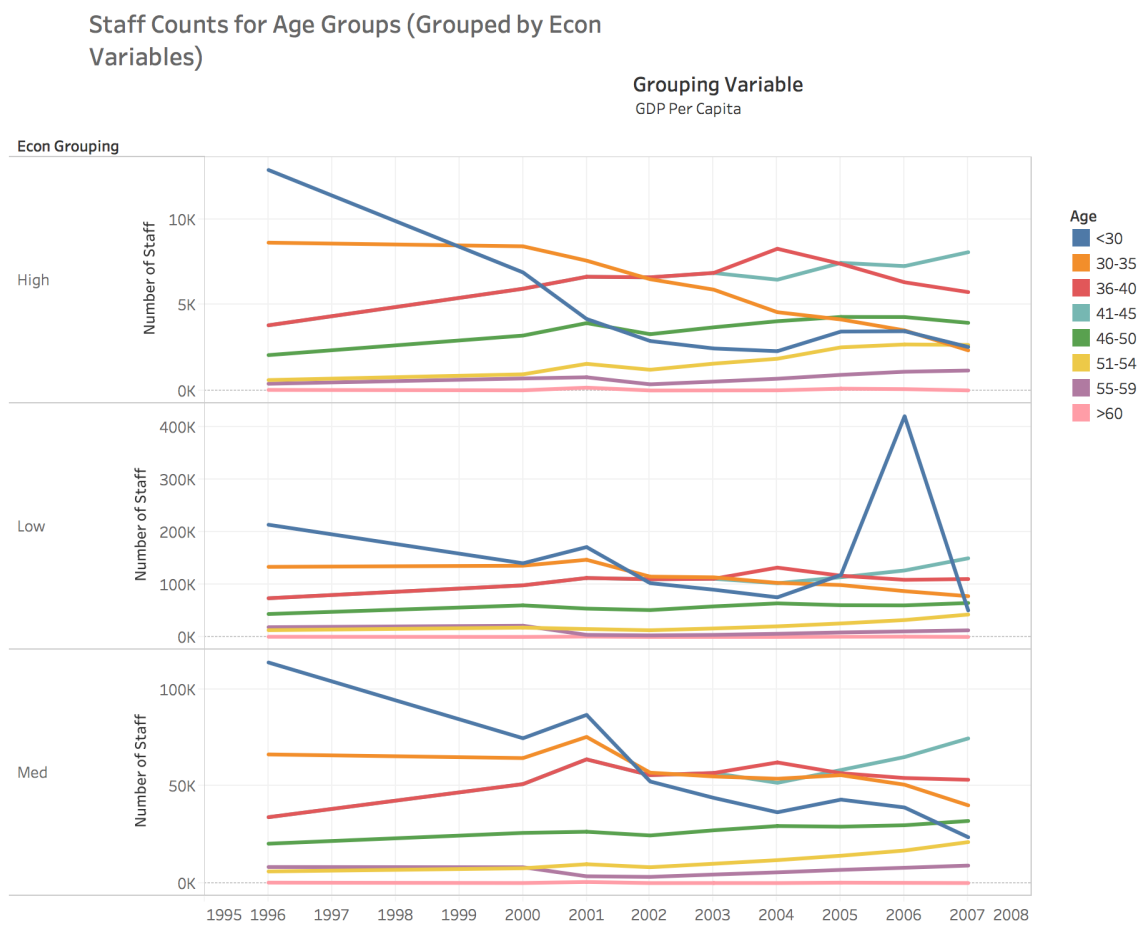


Figure 3: Age Dashboards

## Distribution of Staff by Age

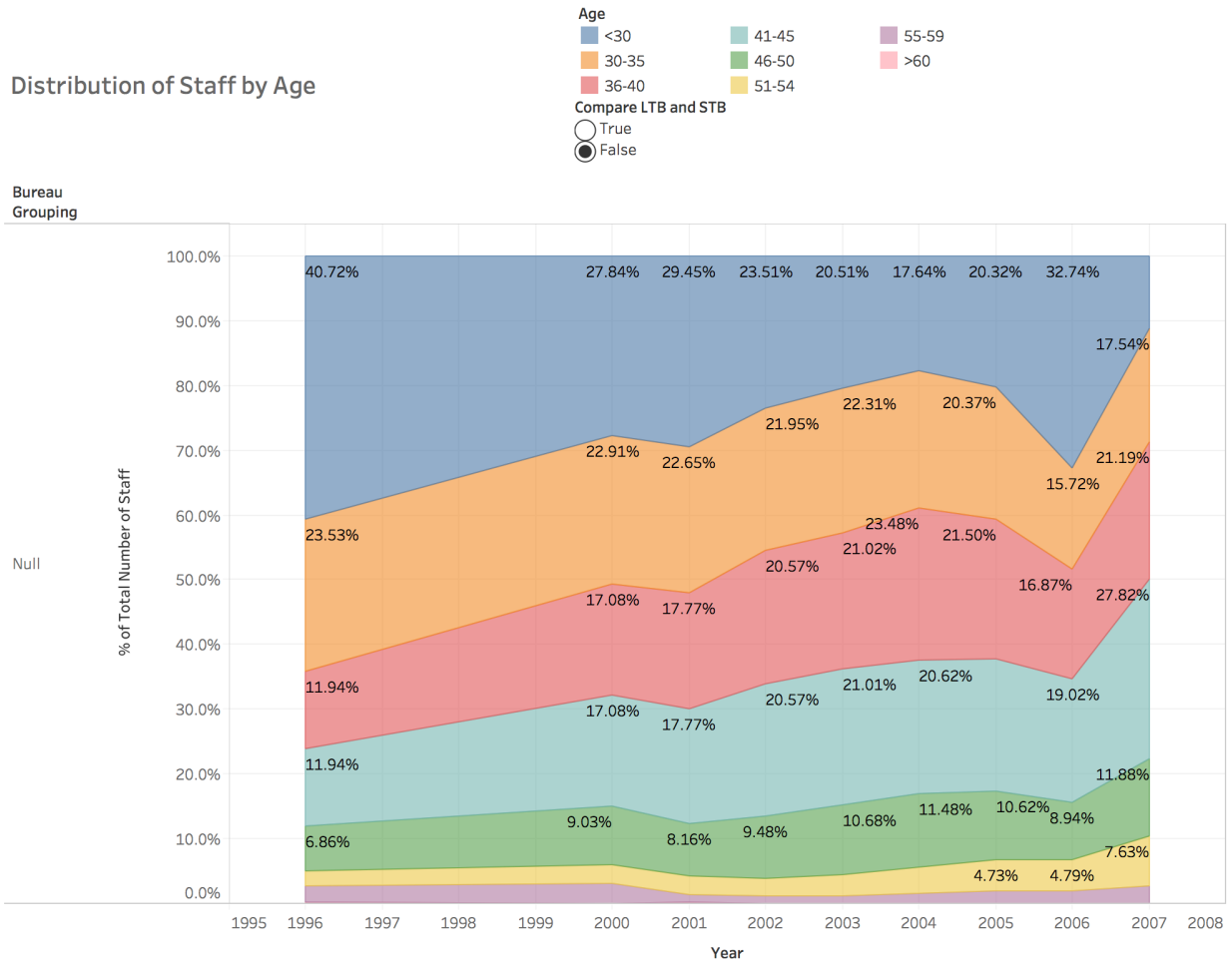


Figure 4: Age Distribution



#### Age - STB

- Under 30
  - o GDP Per capita ( $6 \times 10^{-9}$ )
  - o Public expenditure (0.036)
- 30-35
  - o Urban Area ( $1.92 \times 10^{-6}$ )
  - o Public expenditure (0.0055)
  - o GDP per capita (0.0038)
- 36-40
  - o Above size industrial (0.014)
  - o GDP Per capita (0.01)
- 41-45
  - o Public expenditure 0.003
  - o GDP per capita ( $8.75 \times 10^{-6}$ )
- 46-50 (none)
- 51-54
  - o Urban Area (0.00013)
  - o Above size industrial (0.0006)
  - o GDP per capita ( $6.41 \times 10^{-6}$ )
- 55-59 (none)
- 60 and above (none)

#### Age - LTB

- Under 30
  - o Urban area: 0.0456
  - o GDP per capita:  $2.34 \times 10^{-6}$
- 30-35
  - o Public expenditure:  $7.16 \times 10^{-6}$
- 36-40
  - o Urban area: 0.0417
  - o Public expenditure: 0.0094
  - o GDP per Capita: 0.0041
- 41-45
  - o Urban area: 0.02
  - o Public expenditure: 0.0055
  - o GDP per Capita: 0.0099
- 46-50
  - o Urban area: 0.00076
- 51-54
  - o Urban Area (0.0003911)
  - o Above size industrial (0.0048)
  - o Public Expenditure (0.038)
  - o GDP per capita (0.00057)
- 55-59 (none)
- 60 and above (none)

Figure 5: Significant Vars