

Understanding the Determinants of Healthcare expenditure Amongst Medicare Beneficiaries

Amanny Abuthuraya

ECON 400: Econometrics

Professor Liliana Lawrence

July 18, 2024

Introduction

CMS.gov is a government website that provides data about Medicare and Medicaid health programs. These programs are type of government aid for US citizens, and recently have been extended to include immigrants in some states like Washington (Kamb, L., 2024). Medicare/Medicaid are health insurance programs based on either income, disability or age (USHHS, 2022). Every year, the Medicare Current Beneficiary Survey is conducted by the CMS and NORC. They select and contact random beneficiaries to interview them to fill out the survey (California Health Advocates, 2023)

Since 1991, the MCBS has been working towards the goal of improving Medicare and understanding the needs, costs and experiences of beneficiaries. By understanding the program's coverage, legislators and policymakers can focus on addressing the issues that arise. In this paper, my goal is to analyze the MCBS sample dataset to understand how a variety of variables like age, sex, race, income, chronic conditions, payments and different types of healthcare events effect the total healthcare expenditure.

Medicaid and Medicare are topics of interest to me, because I would like to broaden my scope of analysis and learn more about medical analysis. My analysis can be used in policy development for those associated with high spending trends and help develop care strategies and preventative measures for them.

To understand healthcare expenditure among beneficiaries of Medicare and Medicaid, it is important to review other papers that are relevant to the topic. I will be discussing two papers, "Quality, Health and Spending in Medicare Advantage and Traditional Medicare" published by AMJC and "Out-of-pocket health spending among Medicare beneficiaries: Which chronic diseases are most costly?" published by PLOS ONE, respectively.

"Quality, Health and Spending in Medicare Advantage and Traditional Medicare" compares Medicare Advantage (MA) and Traditional Medicare based on quality, health and cost outcomes (Xu, W. et al, 2021). It was found that MA plans tend to provide better quality care with lower spending costs when compared to TM. It suggests that the Medicare plan can affect the overall health expenditure.

Meanwhile, "Out-of-pocket health spending among Medicare beneficiaries: Which chronic diseases are most costly?", examines how supplemental costs, that are usually paid by "gap insurances", are associated with higher expenditure on healthcare (Johnston, K et al, 2020). Beneficiaries who use gap insurances tend to use their healthcare services more, even with Medicare.

After reviewing both research papers, it can be said that the type of insurance that is used by Medicare beneficiaries contributed to their spending trend. With the "higher end" of the line insurances having the most medical costs associated with them. Also, the health condition of the beneficiaries is another important factor in spending trends. Together, these findings suggest that the structure of insurance coverage is important when understanding spending patterns.

Data

The data is taken from CMS.gov, Medicare Current Beneficiary Survey – Cost Supplement. A public use file containing information on expenditures and payment sources allowing researchers to conduct analysis on Medicare beneficiaries living only in the community (cms.gov, 2024).

The dataset originally had over 30 columns, many of which were not needed for this research, so they have been removed. Also, the data had a few categorical data columns that were transformed into separate dummy variables. Meanwhile there were a lot of adjusted columns that were not of use and have been excluded from the original dataset. ID columns and other arbitrary columns were also removed before starting analysis.

Models & Methodologies

1) The main columns of the dataset are:

- Independent variable: Totalpayment
- Dependent variables:
 - Categorical variables: sex, race, income and age
 - Numerical variables: numberchroniccond, dentalevent, visionevent, hearingevent, homehealthevent, inpatientevent, medicalproviderevent, outpatientevent, prescribemedicine, medicarepayment, medicaidpayment, medicareadvantagepayment, privateinsurancepayment, outofpocketpayment, uncollectedliability and otherpayments

The word “event” in the variables refer to the places where the patient was treated, and “payment” refers to how their treatment were paid for. The categorical variables needed to be transformed. Each category (for example category 1 & 2) was transformed into a separate dummy variable (0,1) carrying their category name instead for easy analysis.

2) Methodologies

- Min, Mean and Max:

sex	numberchroniccond	dentalevent	
Min. :0.0000	Min. :1.000	Min. : 0.0	
1st Qu.:0.0000	1st Qu.:2.000	1st Qu.: 0.0	
Median :0.0000	Median :2.000	Median : 143.2	
Mean :0.4575	Mean :2.348	Mean : 819.0	
3rd Qu.:1.0000	3rd Qu.:3.000	3rd Qu.: 597.0	
Max. :1.0000	Max. :3.000	Max. :22708.7	
visionevent	hearingevent	homehealthevent	inpatientevent
Min. : 0.0	Min. : 0.0	Min. : 0.0	Min. : 0
1st Qu.: 0.0	1st Qu.: 0.0	1st Qu.: 0.0	1st Qu.: 0
Median : 30.0	Median : 0.0	Median : 0.0	Median : 0
Mean : 280.0	Mean : 129.1	Mean : 288.4	Mean : 2034
3rd Qu.: 189.7	3rd Qu.: 0.0	3rd Qu.: 0.0	3rd Qu.: 0
Max. :14083.4	Max. :9615.3	Max. :27825.2	Max. :110040

medicalproviderevent	outpatientevent	prescribemedicine
Min. : 0.0	Min. : 0.00	Min. : 0.0
1st Qu.: 685.9	1st Qu.: 21.54	1st Qu.: 253.5
Median : 1904.2	Median : 254.54	Median : 896.6
Mean : 3984.6	Mean : 2278.99	Mean : 5333.7
3rd Qu.: 4598.2	3rd Qu.: 1283.01	3rd Qu.: 4366.3
Max. : 65844.6	Max. : 105355.57	Max. : 201894.3
totalpayment	medicarepayment	medicaidpayment
Min. : 0	Min. : 0.0	Min. : 0.0
1st Qu.: 2707	1st Qu.: 316.5	1st Qu.: 0.0
Median : 6897	Median : 1961.5	Median : 0.0
Mean : 15148	Mean : 8017.0	Mean : 489.6
3rd Qu.: 16116	3rd Qu.: 6337.2	3rd Qu.: 0.0
Max. : 307916	Max. : 215227.4	Max. : 34316.1
medicareadvantagepayment	privateinsurancepayment	outofpocketpayment
Min. : 0	Min. : 0.0	Min. : 0.0
1st Qu.: 0	1st Qu.: 0.0	1st Qu.: 354.2
Median : 0	Median : 0.0	Median : 1139.1
Mean : 2294	Mean : 1095.2	Mean : 2440.3
3rd Qu.: 1156	3rd Qu.: 645.8	3rd Qu.: 2863.8
Max. : 91750	Max. : 62549.5	Max. : 37203.2
uncollectedliability	otherpayment	age_g1
Min. : 0.0	Min. : -168.45	Min. : 0.0000
1st Qu.: 0.0	1st Qu.: 0.00	1st Qu.: 0.0000
Median : 0.0	Median : 0.00	Median : 0.0000
Mean : 303.5	Mean : 507.81	Mean : 0.1678
3rd Qu.: 109.9	3rd Qu.: 13.76	3rd Qu.: 0.0000
Max. : 24358.1	Max. : 41444.82	Max. : 1.0000

Since **otherpayment** has a negative value of -168.45, it has been removed because it isn't logical for a payment to be in negative and should be in positive form. Other than that, the maximum values are going to be kept because one medical bill could easily reach high amounts. It isn't something impossible depending on the patient's health.

- Correlation Plot (Graph 1 in Appendix)

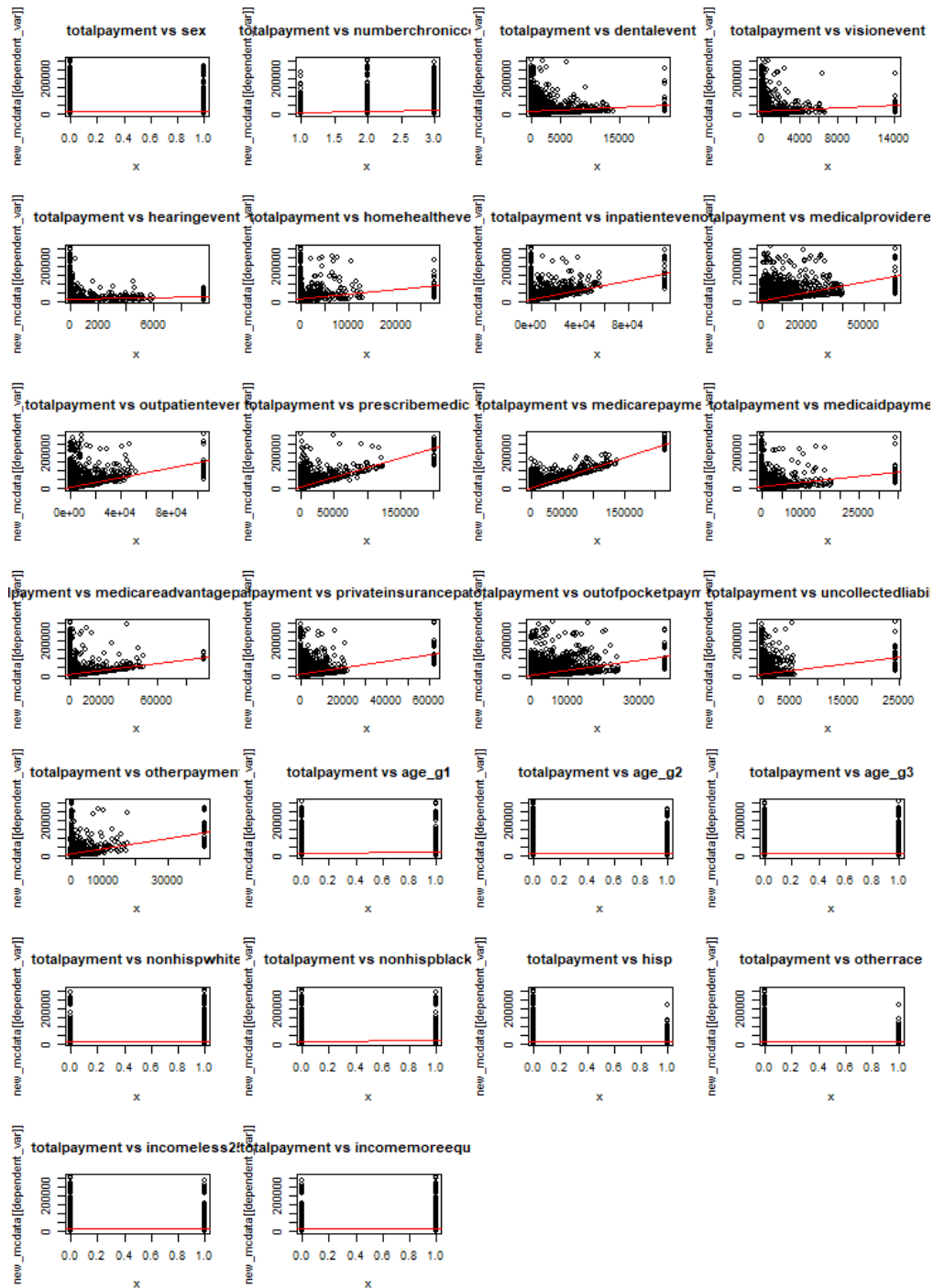
Correlation plots are used to check the multicollinearity issue within the data. Since the data frame was too big to be viewed within R Studio, it was saved as a CSV file to the local device. After going through the columns, there's no issue of severely high correlations (between the values of 0.8 & 0.9). We do have some issues with moderately high correlations, for example:

- 1- Correlation of **homehealthevent** and **prescribemedicine** = 0.717
- 2- Correlation of **inpatientevent** and **medicaidpayment** = 0.798
- 3- Correlation of **medicalproviderevent** and **totalpayment** = 0.717
- 4- Correlation of **prescribemedicine** and **totalpayment** = 0.717

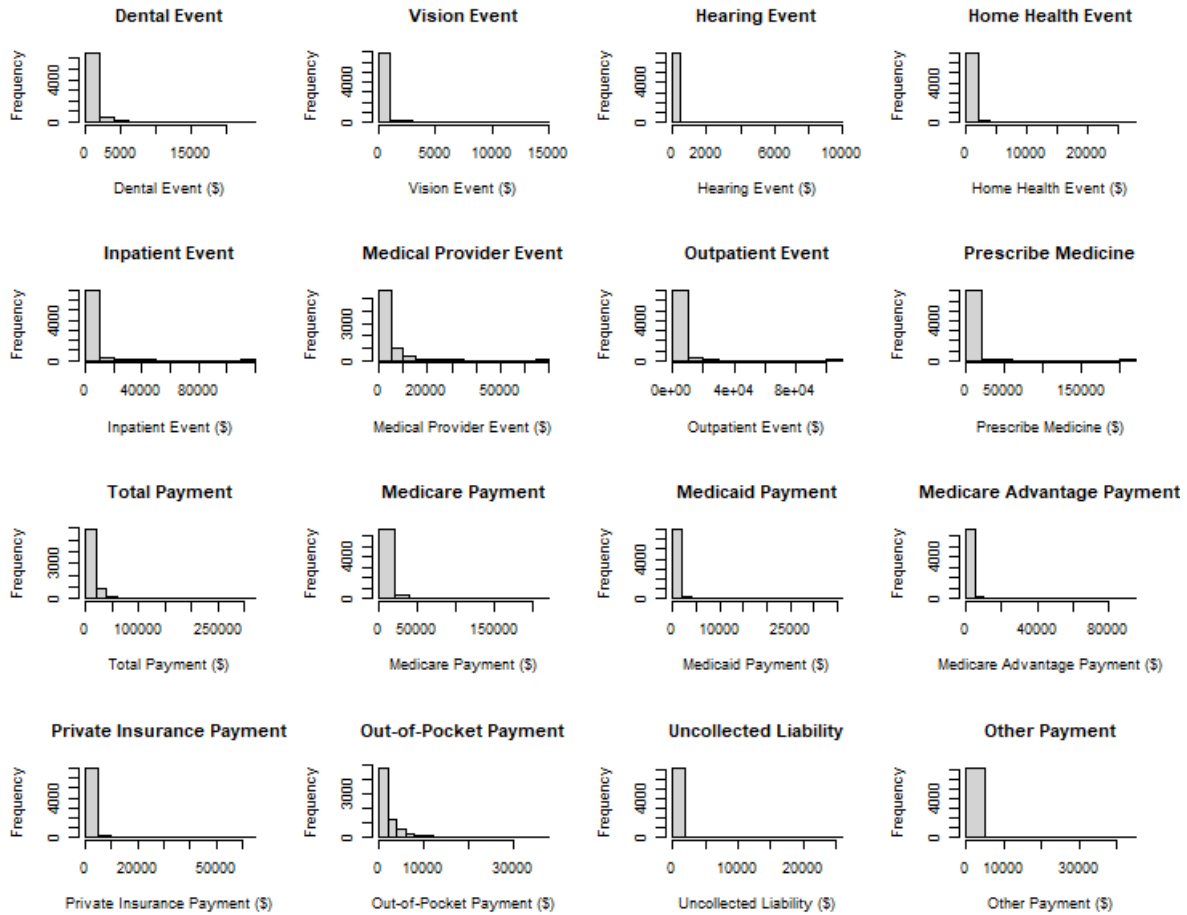
They could cause some problems to a degree, but VIFs of the data need to be checked before omitting any correlated variables.

- Linearity & Skewness

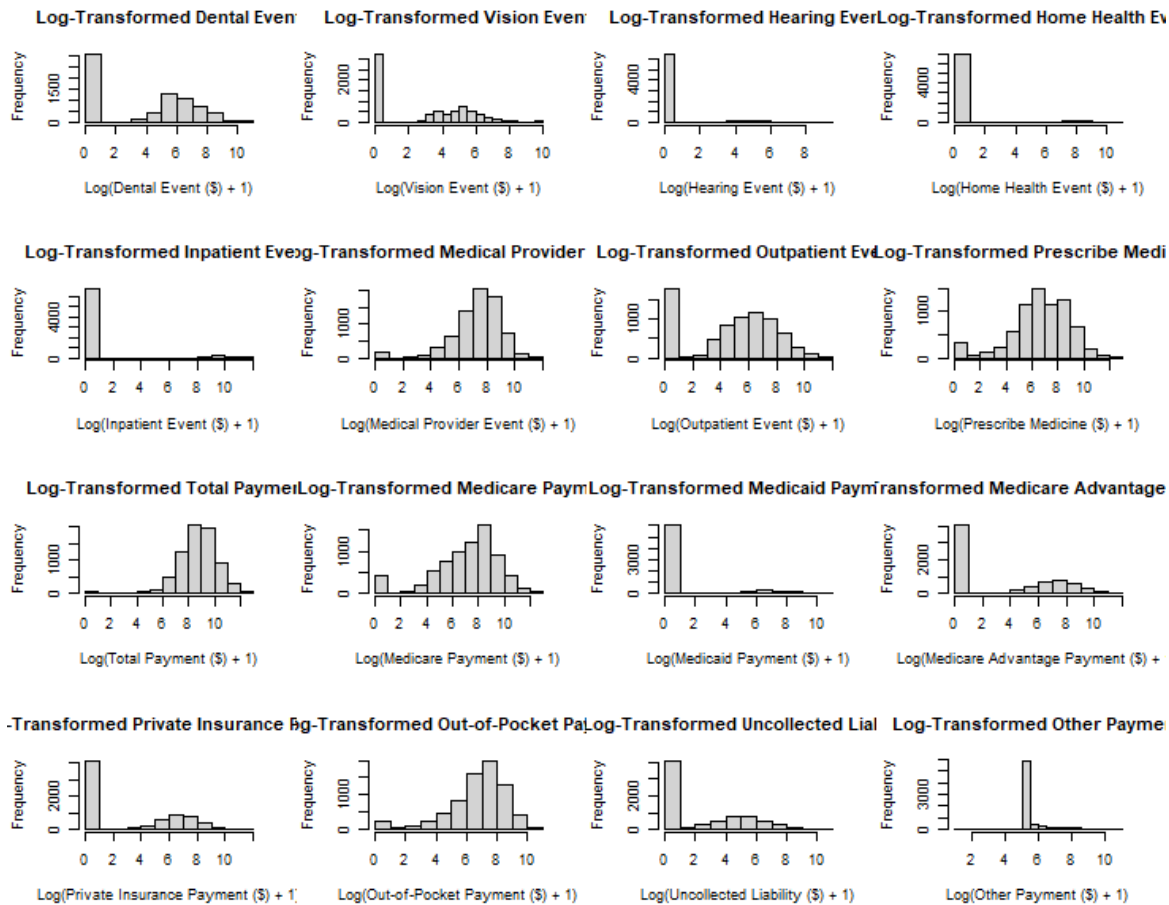
This is how the data looks like when plotting



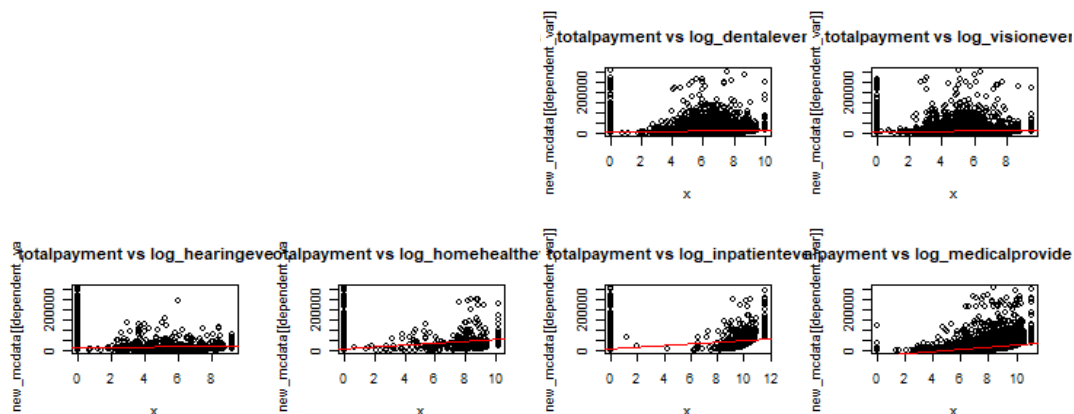
Since the data has some issues with linearity not being perfect, i.e medicareadvantagepayment vs outofpocketpayment, the data will have to be transformed. To see where the issue lies, histograms are drawn:

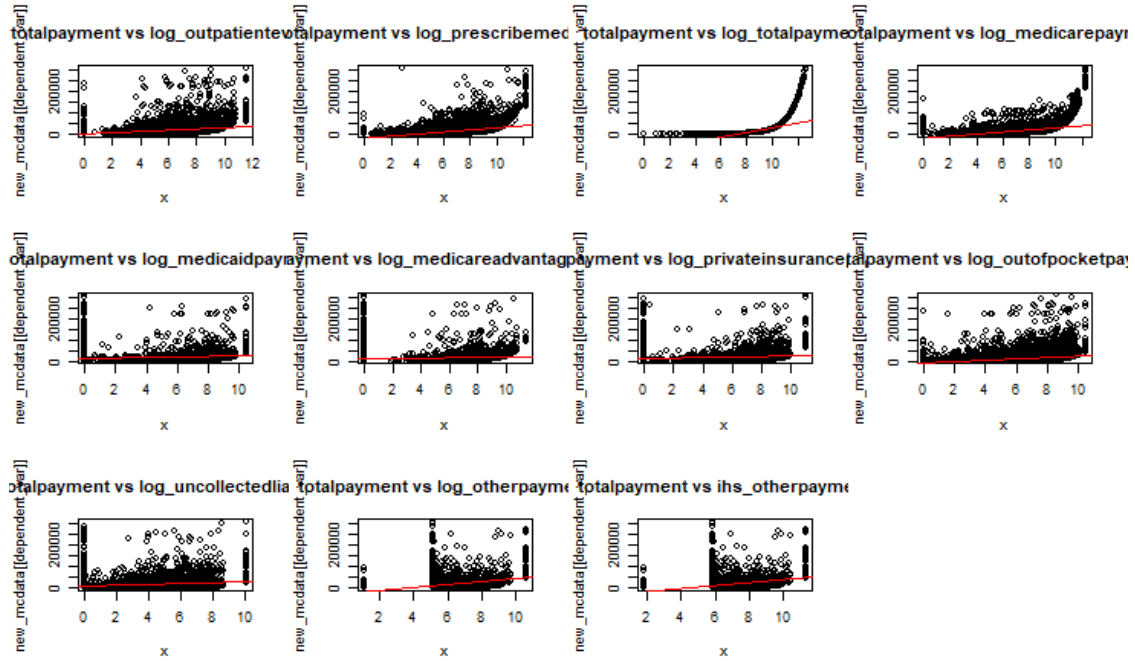


The issue happens to be an issue of the excess of zeros in the data which can be fixed with a `log()` function to help create a more normal distribution. This is how the data looks like after applying the `log()` function:



After applying the `log()` function, the linearity can be checked again. It indicates that the function did help reduce the non-linearity between the variables. Even though it isn't perfect linearity, it can still be used in a regression model





- First OLS Model

Totalpayment

$$\begin{aligned}
 = & B1 \text{ sex} + B2 \text{ numberchroniccond} + B3 \text{ log_dentalevent} \\
 & + B4 \text{ log_visionevent} + B5 \text{ log_hearingevent} \\
 & + B6 \text{ log_medicalproviderevent} + B7 \text{ log_outpatientevent} \\
 & + B8 \text{ log_prescribemedicine} + B9 \text{ log_medicarepayment} \\
 & + B10 \text{ log_medicaidpayment} + B11 \text{ log_medicareadvantagepayment} \\
 & + B12 \text{ log_privateinsurancepayment} + B13 \text{ log_outofpocketpayment} \\
 & + B14 \text{ log_uncollectedliability} + B15 \text{ log_otherpayment} + B16 \text{ age_g1} \\
 & + B17 \text{ age_g2} + B18 \text{ hisp} + B19 \text{ nonhispblack} + B20 \text{ nonhispwhite} \\
 & + B21 \text{ incomeless25}
 \end{aligned}$$

Every variable was included in the model to see which variables are the most significant. Significance is based on the p-value of 0.05 or less. Not every categorical variable was included to prevent the dummy variable trap.

```

Call:
lm(formula = totalpayment ~ sex + numberchroniccond + log_dentalevent +
    log_visionevent + log_hearingevent + log_medicalproviderevent +
    log_outpatientevent + log_prescribemedicine + log_medicarepayment +
    log_medicaidpayment + log_medicareadvantagepayment + log_privateinsurancepayment +
    log_outofpocketpayment + log_uncollectedliability + log_otherpayment +
    age_g1 + age_g2 + hisp + nonhispblack + nonhispwhite + incomeless25,
    data = new_mcdata)

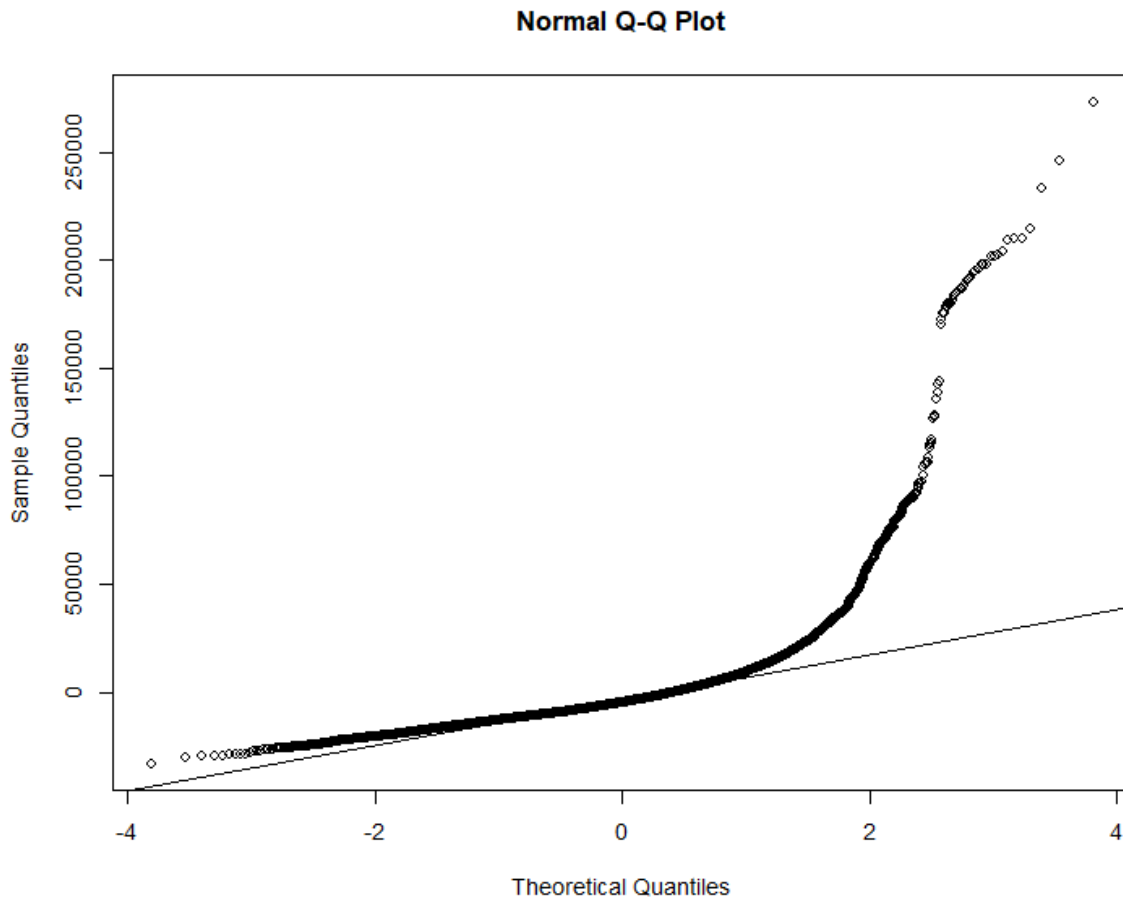
Residuals:
    Min       1Q   Median       3Q      Max
-32973 -10507  -4759   3640  273324

Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)    -56714.88   2405.32  -23.579  < 2e-16 ***
sex              282.85    526.84    0.537  0.591362
numberchroniccond -1528.16   406.12   -3.763  0.000169 ***
log_dentalevent   -16.37     89.26   -0.183  0.854463
log_visionevent  -400.23    101.36   -3.949  7.93e-05 ***
log_hearingevent -189.88    146.69   -1.294  0.195549
log_medicalproviderevent  976.79    218.71    4.466  8.08e-06 ***
log_outpatientevent  924.79     99.33    9.310  < 2e-16 ***
log_prescribemedicine 1696.46    168.37   10.076  < 2e-16 ***
log_medicarepayment  2152.06    151.95   14.163  < 2e-16 ***
log_medicaidpayment 1208.08    145.86    8.282  < 2e-16 ***
log_medicareadvantagepayment  681.58     95.17    7.162  8.72e-13 ***
log_privateinsurancepayment  649.42    114.90    5.652  1.65e-08 ***
log_outofpocketpayment  991.97    203.92    4.864  1.17e-06 ***
log_uncollectedliability  855.41    117.15    7.302  3.13e-13 ***
log_otherpayment  4094.30    301.70   13.571  < 2e-16 ***
age_g1          5769.79    835.63    6.905  5.45e-12 ***
age_g2          925.19    579.70    1.596  0.110534
hisp            356.99   1410.10    0.253  0.800145
nonhispblack    3702.02   1420.97    2.605  0.009198 **
nonhispwhite    757.54   1185.66    0.639  0.522895
incomeless25   -711.64    712.28   -0.999  0.317780
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

At the Alpha level of 0.05, the variables sex, log_dentalevent, log_hearingevent, age_g1, age_g2, hisp, nonhispwhite and incomeless25 are not significant and will be omitted from the model. The model also has an Adjusted R-Squared Value of 0.3168, which indicates how good a model can predict. Of course, this model can become better by omitting the insignificant variables and fixing further aspects of it.

Another aspect to check in the model is the residuals of the model. Based on this following graph, the model is an okay fit with some outliers at the tails. The upper tail has the most outliers, but this is probably due to the huge outliers in many of the different forms of payment in our data. An additional reason could be due to the existence of zero inflated values withing the model which we'll talk about later.



- Second OLS Model

$$\text{Totalpayment} = B1 \text{numberchroniccond} + B2 \log_visionevent + B3 \log_medicalproviderevent + B4 \log_outpatientevent + B5 \log_prescribemedicine + B6 \log_medicarepayment + B7 \log_medicaidpayment + B8 \log_medicareadvantagepayment + B9 \log_privateinsurancepayment + B10 \log_outofpocketpayment + B11 \log_uncollectedliability + B12 \log_otherpayment + B13 \text{nonhisblack}$$

```

Call:
lm(formula = totalpayment ~ numberchroniccond + log_visionevent +
    log_medicalproviderevent + log_outpatientevent + log_prescribemedicine +
    log_medicarepayment + log_medicaidpayment + log_medicareadvantagepayment +
    log_privateinsurancepayment + log_outofpocketpayment + log_uncollectedliability +
    log_otherpayment + nonhisblack, data = new_mcdata)

Residuals:
    Min       1Q   Median       3Q      Max
-30409 -10597  -4972   3612  273432

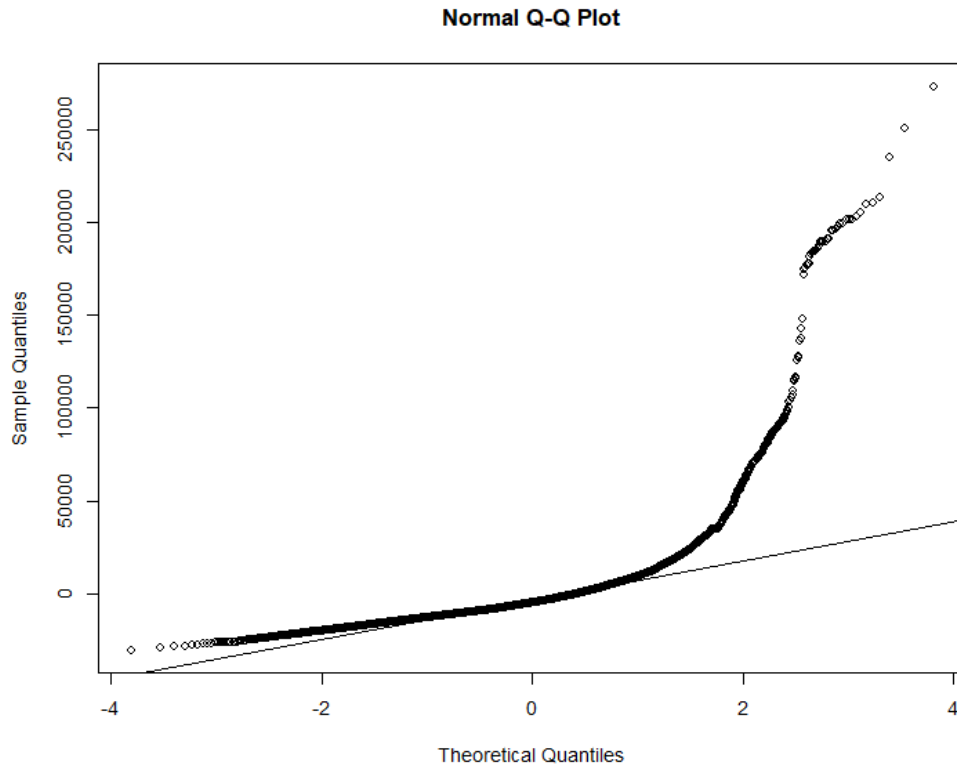
Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)    -53682.23   1987.97  -27.003  < 2e-16 ***
numberchroniccond  -1973.69    399.73   -4.938  8.08e-07 ***
log_visionevent   -485.69    100.57   -4.829  1.40e-06 ***
log_medicalproviderevent    983.90    218.75    4.498  6.97e-06 ***
log_outpatientevent    972.11     99.28    9.792  < 2e-16 ***
log_prescribemedicine   1815.44    166.66   10.893  < 2e-16 ***
log_medicarepayment   2071.99    151.28   13.696  < 2e-16 ***
log_medicaidpayment   1396.16    134.64   10.370  < 2e-16 ***
log_medicareadvantagepayment    630.73     93.15    6.771  1.38e-11 ***
log_privateinsurancepayment    593.82    110.71    5.364  8.41e-08 ***
log_outofpocketpayment    925.35    192.32    4.811  1.53e-06 ***
log_uncollectedliability    882.07    117.10    7.533  5.56e-14 ***
log_otherpayment    4073.47    300.91   13.537  < 2e-16 ***
nonhisblack       3587.75    901.21    3.981  6.93e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 22030 on 7309 degrees of freedom
Multiple R-squared:  0.3138,    Adjusted R-squared:  0.3125
F-statistic: 257.1 on 13 and 7309 DF,  p-value: < 2.2e-16

```

After adjusting the variables in the model, all the variables seem to be significant to totalpayment. They are significant at the Alpha level of 0.05, which means that they have a 95% chance of predicting the variable totalpayment. Even though the model was created to become better, the Adjusted R-Squared decreased to 0.3125.

This is the residual plot of the model after the adjustment. There's not much difference in the outliers. Since the outliers are causing a problem, another solution called Box-Cox Transformation will be used.



- Third OLS Model

$(Totalpayment^{best.lam})$

$$\begin{aligned}
 = & B0 + B1 numberchroniccond + B2 \log(visionevent) \\
 & + B3 \log(medicalproviderevent) + B4 \log(outpatientevent) \\
 & + B5 \log(prescribemedicine) + B6 \log(medicarepayment) \\
 & + B7 \log(medicaidpayment) + B8 \log(medicareadvantagepayment) \\
 & + B9 * \log(privateinsurancepayment) \\
 & + B10 \log(outofpocketpayment) + B11 \log(uncollectedliability) \\
 & + B12 \log(otherpayment) + B13 nonhisblack + \varepsilon
 \end{aligned}$$

If you want to understand why $totalpayment^{best.lam}$, you can watch Math et al's video on YouTube about the Box-Cox Transformation linked in the references of this paper. This is one of many possible solutions to the problem of zero-inflated values within our dataset. These are the results:

```

Call:
lm(formula = (totalpayment^best.lam) ~ numberchroniccond + log_visionevent +
  log_medicalproviderevent + log_outpatientevent + log_prescribemedicine +
  log_medicarepayment + log_medicaidpayment + log_medicareadvantagepayment +
  log_privateinsurancepayment + log_outofpocketpayment + log_uncollectedliability +
  log_otherpayment + nonhispblack, data = new_mcddata)

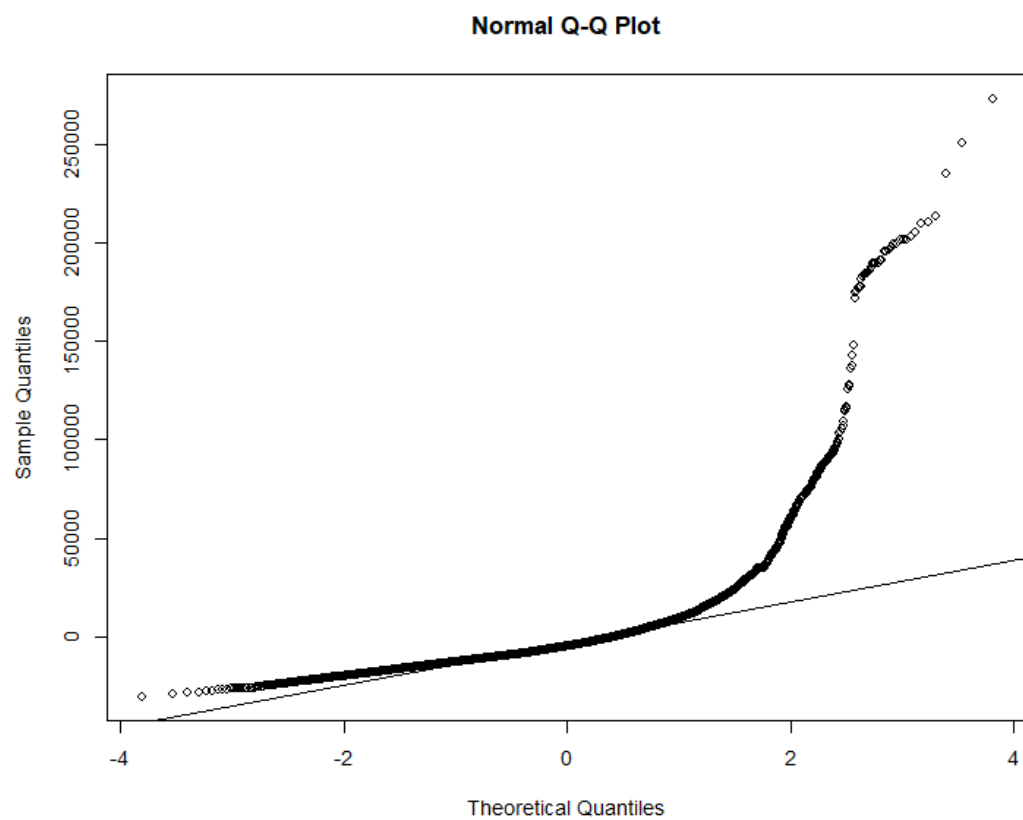
Residuals:
    Min       1Q   Median       3Q      Max
-2.1956 -0.6904 -0.2566  0.4528  6.3602

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   -1.141337   0.091697  -12.447 < 2e-16 ***
numberchroniccond -0.062365   0.018438   -3.382 0.000722 ***
log_visionevent -0.016764   0.004639   -3.614 0.000304 ***
log_medicalproviderevent 0.225721   0.010090   22.371 < 2e-16 ***
log_outpatientevent 0.088344   0.004579   19.293 < 2e-16 ***
log_prescribemedicine 0.195800   0.007687   25.470 < 2e-16 ***
log_medicarepayment 0.181348   0.006978   25.989 < 2e-16 ***
log_medicaidpayment 0.150238   0.006210   24.192 < 2e-16 ***
log_medicareadvantagepayment 0.089258   0.004297   20.774 < 2e-16 ***
log_privateinsurancepayment 0.080264   0.005107   15.717 < 2e-16 ***
log_outofpocketpayment 0.198580   0.008871   22.385 < 2e-16 ***
log_uncollectedliability 0.056850   0.005401   10.526 < 2e-16 ***
log_otherpayment 0.214636   0.013880   15.464 < 2e-16 ***
nonhispblack 0.130993   0.041569    3.151 0.001633 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.016 on 7309 degrees of freedom
Multiple R-squared:  0.736,    Adjusted R-squared:  0.7355
F-statistic: 1567 on 13 and 7309 DF,  p-value: < 2.2e-16

```

All our model variables are still significant to the variable total payment. The only difference is that now our Adjusted R-Squared has the value of 0.7355 which indicates a much better fit this model is. It isn't perfect but it is better than what we calculated in our first model. Our residuals though haven't changed a lot either from our first model. Other solutions like the Poisson Regression could help fix the Adjust R-Squared of the model and give it an even better fit, but there's implications that come with it. One of the implications is that the interpretation of the results is different due to the different type of regression being used. Due to that reason, this model will be chosen as an explanation of the variable total payment.



Conclusion

In conclusion, our final model is:

$$\begin{aligned}
 & \text{Transformed Total Payment} = \\
 & -1.11596 + (-0.06229 * \text{Number of Chronic Conditions}) + (-0.01670 \\
 & \quad * \text{Log of Vision Events}) + (0.22373 * \text{Log of Medical Provider Events}) \\
 & \quad + (0.08858 * \text{Log of Outpatient Events}) + (0.19537 \\
 & \quad * \text{Log of Prescribed Medicine}) + (0.18112 \\
 & \quad * \text{Log of Medicare Payment}) + (0.14992 * \text{Log of Medicaid Payment}) \\
 & \quad + (0.08908 * \text{Log of Medicare Advantage Payment}) + (0.08028 \\
 & \quad * \text{Log of Private Insurance Payment}) + (0.19716 * \text{Log of Out-of-Pocket Payment}) \\
 & \quad + (0.05725 * \text{Log of Uncollected Liability}) \\
 & \quad + (0.21521 * \text{Log of Other Payment}) + (0.13121 * \text{Non-Hispanic Black})
 \end{aligned}$$

Where each unit increase in each variable causes the Transformed Total Payment to increase by the coefficient number next to it. Even though a variable like Number of Chronic Conditions has a negative effect on Total Payment, may seem counterintuitive, but the other variables that explain the payment explain why this number would be less. Another counterintuitive variable is Vision Events and the lowering of the Total Payment, this could be possibly due to the relatively lower cost associated with it. Meanwhile the rest of the variables, except Non-Hispanic Black, have a positive effect on the Total Payment. This suggests that the higher their values are, the higher the Total payment is. The last variable in our model, Non-Hispanic Black suggests that beneficiaries who identify as Black tend to have higher healthcare expenditure in comparison to all other races. This could be in relation to socio-economic and higher numbers of chronic conditions (Gaskin et al., 2021, p.144)

My model highlights the importance of factors like insurance coverage, chronic conditions and socioeconomic factors may play a role into beneficiaries' healthcare expenditures. By using the model to calculate their impact on expenditure, policies can be improved and preventative healthcare strategies can be implemented for those associated with high costs.

Resources

- Kamb, L. (2022, March 30). WA expanding health care options for undocumented immigrants. The Seattle Times. <https://www.seattletimes.com/seattle-news/politics/wa-expanding-health-care-options-for-undocumented-immigrants/>
- U.S. Department of Health and Human Services. (n.d.). Who is eligible for Medicaid? HHS.gov. Retrieved June 14, 2024, from <https://www.hhs.gov/answers/medicare-and-medicaid/who-is-eligible-for-medicaid/index.html>
- Centers for Medicare & Medicaid Services. (2021, November 17). Medicare Current Beneficiary Survey Cost & Supplement [Data set]. Data.cms.gov. <https://data.cms.gov/medicare-current-beneficiary-survey-mcbs/medicare-current-beneficiary-survey-cost-supplement>
- California Health Advocates. (2023, March 3). Have you been contacted for the Medicare Current Beneficiary Survey? Verify your participation & prevent fraud. California Health Advocates. <https://cahealthadvocates.org/have-you-been-contacted-for-the-medicare-current-beneficiary-survey-verify-your-participation-prevent-fraud/>
- Johnston, K. J., Pinska, L., Munro, H., Parker, D., Nguyen, H., & Quiter, E. (2019). Determining out-of-pocket expenditure towards care in the Medicare Current Beneficiary Survey. PLOS ONE. <https://doi.org/10.1371/journal.pone.0222539>
- Xu, W. Y., Fonseca, V., Finch, M. D., Kim, M. J., & Frett, B. (2021). Quality, health, and spending in Medicare Advantage and traditional Medicare. The American Journal of Managed Care (AMJC). <https://doi.org/10.37765/ajmc.2021.88581>
- Math et al. (2022, January 24). The derivative of a logarithm: why the log derivative is one over x [Video]. YouTube. <https://www.youtube.com/watch?v=vGOEpjz2Ks>
- Gaskin, D. J., Dinwiddie, G. Y., & Chan, K. S. (2021). Racial and ethnic disparities in health care access and utilization under the Affordable Care Act. Medical Care Research and Review. <https://doi.org/10.1177/1077558720965111>

- Comparison of OLS models

	Dependent variable:		
	(1)	totalpayment (2)	(3)
sex	282.853 (526.839)		
numberchroniccond	-1,528.159*** (406.121)	-1,973.691*** (399.728)	-1,973.691*** (399.728)
log_dentalevent	-16.373 (89.261)		
log_visionevent	-400.230*** (101.355)	-485.691*** (100.572)	-485.691*** (100.572)
log_hearingevent	-189.881 (146.688)		
log_medicalproviderevent	976.790*** (218.710)	983.899*** (218.750)	983.899*** (218.750)
log_outpatientevent	924.787*** (99.330)	972.108*** (99.276)	972.108*** (99.276)
log_prescribemedicine	1,696.461*** (168.373)	1,815.444*** (166.661)	1,815.444*** (166.661)
log_medicarepayment	2,152.065*** (151.947)	2,071.990*** (151.279)	2,071.990*** (151.279)
log_medicaidpayment	1,208.082*** (145.862)	1,396.163*** (134.637)	1,396.163*** (134.637)
log_medicareadvantagepayment	681.581*** (95.166)	630.726*** (93.149)	630.726*** (93.149)
log_privateinsurancepayment	649.423*** (114.903)	593.822*** (110.713)	593.822*** (110.713)
log_outofpocketpayment	991.968*** (203.923)	925.353*** (192.321)	925.353*** (192.321)
log_otherpayment	4,094.299*** (301.697)	4,073.465*** (300.908)	4,073.465*** (300.908)
age_g1	5,769.789*** (835.635)		
age_g2	925.195 (579.699)		
hisp	356.994 (1,410.102)		
nonhispblack	3,702.017*** (1,420.965)	3,587.755*** (901.215)	3,587.755*** (901.215)
nonhispwhite	757.543 (1,185.663)		
incomeless25	-711.642 (712.283)		
Constant	-56,713.880*** (2,405.316)	-53,681.230*** (1,987.973)	-53,682.230*** (1,987.973)
Observations	7,323	7,323	7,323
R2	0.319	0.314	0.314
Adjusted R2	0.317	0.313	0.313
Residual Std. Error	21,966.410 (df = 7301)	22,034.100 (df = 7309)	22,034.100 (df = 7309)
F Statistic	162.642*** (df = 21; 7301)	257.057*** (df = 13; 7309)	257.057*** (df = 13; 7309)

- Code:

```
# All the needed libraries
```

```
library("stargazer")
```

```
library("readxl")
```

```
library("zoo")
```

```
library("dplyr")
```

```
library("lubridate")
```

```
library("purrr")
```

```
library("ggplot2")
```

```
library("moments")
```

```
library("car")
```

```
library("lmtest")
```

```
library("corrplot")
```

```
library("fastDummies") #used for changing categorical data into dummy variables
```

```
library("gt") #used for formatting tables nicely when knitting
```

```
library("GGally")
```

```
library("psych")
```

```
library("moments") #for testing skewness
```

```
library("nortest")
```

```
#installing packages
```

```
#install.packages("gt")
```

```
#install.packages("fastDummies")
```

```
#install.packages("GGally")
```

```
#install.packages("corrplot")
```

```
#install.packages("psych")
```

```
#install.packages("nortest")
```

```
library(caret)
```

```
#Setting working directory for the data
```

```
setwd("D://Bellevue College//24 Spring//ECON 400//Week 8 - Term Paper")
```

```
mcdata <- read_excel("cspuf2021.xls")
```

#In this section, I'll be formatting the dataset for further usage

#renaming variables for easy use <https://www.geeksforgeeks.org/how-to-rename-multiple-columns-in-r/>

```
new_names <- c("age", "sex", "race",
               "income", "numberchroniccond", "dentalevent", "visionevent",
               "hearingevent", "homehealthevent", "inpatientevent",
               "medicalproviderevent", "outpatientevent", "prescribemedicine",
               "totalpayment", "medicarepayment", "medicaidpayment",
               "medicareadvantagepayment", "privateinsurancepayment", "outofpocketpayment",
               "uncollectedliability", "otherpayment")

names(mcdata) <- new_names
```

#Starting with the age variable

Changing categorical variables to separate dummy variables that represent each group

Create new columns that represent each new group

```
mcdata$age_g1 <- c(0)
mcdata$age_g2 <- c(0)
mcdata$age_g3 <- c(0)
```

#to delete the mistake columns

```
#mydata2 = select(mcdata, -28, -29)
#mcdata <- mydata2
#rm(mydata2)
```

#Dividing the age variable amongst the 3 new columns

```
mcdata$age_g1 <- ifelse(mcdata$age == 1, 1, 0)
```

```
mcddata$age_g2 <- ifelse(mcddata$age == 2, 1, 0)
```

```
mcddata$age_g3 <- ifelse(mcddata$age == 3, 1, 0)
```

```
#sex variable
```

```
# Changing categorical variables to 0 & 1 dummy variables within the same column
```

```
# value of 1 = male, value of 0 = female
```

```
mcddata$sex <- ifelse(mcddata$sex == 1, 1, 0)
```

```
print(head(mcddata$sex, n=5))
```

```
#race variable:
```

```
#1:Non-Hispanic white
```

```
#2:Non-Hispanic black
```

```
#3:Hispanic
```

```
#4:Other
```

```
# Create new columns that represent each new group
```

```
mcddata$nonhispwhite<-c(0)
```

```
mcddata$nonhispblack<-c(0)
```

```
mcddata$hisp<-c(0)
```

```
mcddata$otherrace<-c(0)
```

```
#Assigning dummy variables to each column
```

```
mcddata$nonhispwhite <- ifelse(mcddata$race == 1, 1, 0)
```

```
mcddata$nonhispblack <- ifelse(mcddata$race == 2, 1, 0)
```

```
mcddata$hisp <- ifelse(mcddata$race == 3, 1, 0)
```

```
mcddata$otherrace <- ifelse(mcddata$race == 4, 1, 0)
```

```
#income variable
```

```
# 1:<$25,000
```

```
# 2:>=$25,000
```

```
# Create new columns that represent each new group
```

```
mcddata$incomeless25<-c(0)
```

```

mcdata$incomemoreequal25<-c(0)
#Using if else statement to assign values
mcdata$incomeless25 <- ifelse(mcdata$income == 1, 1, 0)
mcdata$incomemoreequal25 <- ifelse(mcdata$income == 2, 1, 0)

# Removing any possible N/A values
mcdata <- na.omit(mcdata)

# Dropping the categorical columns after switching them to dummy columns
#https://www.listendata.com/2015/06/r-keep-drop-columns-from-data-frame.html
drops <- c("age", "race", "income")
# Creating a new data set after the changes
new_mcdata= mcdata[ , !(names(mcdata)%in% drops)]

#Visualizing/viewing the data

# Checking for correlation then saving it as a csv file due to the 22*22 design
cor_matrix <- cor(select_if(new_mcdata, is.numeric))
cor_matrix
write.csv(cor_matrix, "cor_matrix.csv", row.names = TRUE)

#Testing for linearity of total payment vs each variable
par(mfrow = c(4, 4))
dependent_var <- "totalpayment"
independent_vars <- setdiff(names(new_mcdata), dependent_var)

# Create a function to plot linearity
par(mfrow = c(4,4))

```

```

plot_linearity <- function(x, y, plot_title) {
  plot(y ~ x, data = new_mcddata, main = plot_title, xlab = deparse(substitute(x)), ylab =
deparse(substitute(y)))
  abline(lm(y ~ x, data = new_mcddata), col = "red")
}
# Loop through each independent variable and plot linearity
for (var in independent_vars) {
  plot_title <- paste(dependent_var, "vs", var)
  plot_linearity(new_mcddata[[var]], new_mcddata[[dependent_var]], plot_title)
}

#Viewing issue with linearity via histograms
#Histograms
#https://www.geeksforgeeks.org/histograms-in-r-language/
#https://bookdown.org/dli/rguide/histogram.html
# Dental event
hist(new_mcddata$dentalevent, main = "Dental Event", xlab = "Dental Event ($)")
# Vision event
hist(new_mcddata$visionevent, main = "Vision Event", xlab = "Vision Event ($)")
# Hearing event
hist(new_mcddata$hearingevent, main = "Hearing Event", xlab = "Hearing Event ($)")
# Home health event
hist(new_mcddata$homehealthevent, main = "Home Health Event", xlab = "Home Health
Event ($)")
# Inpatient event
hist(new_mcddata$inpatientevent, main = "Inpatient Event", xlab = "Inpatient Event ($)")
# Medical provider event
hist(new_mcddata$medicalproviderevent, main = "Medical Provider Event", xlab = "Medical
Provider Event ($)")
# Outpatient event

```

```

hist(new_mdata$outpatientevent, main = "Outpatient Event", xlab = "Outpatient Event
($)")

# Prescribe medicine

hist(new_mdata$prescribemedicine, main = "Prescribe Medicine", xlab = "Prescribe
Medicine ($)")

# Total payment

hist(new_mdata$totalpayment, main = "Total Payment", xlab = "Total Payment ($)")

# Medicare payment

hist(new_mdata$medicarepayment, main = "Medicare Payment", xlab = "Medicare
Payment ($)")

# Medicaid payment

hist(new_mdata$medicaidpayment, main = "Medicaid Payment", xlab = "Medicaid
Payment ($)")

# Medicare Advantage payment

hist(new_mdata$medicareadvantagepayment, main = "Medicare Advantage Payment",
xlab = "Medicare Advantage Payment ($)")

# Private insurance payment

hist(new_mdata$privateinsurancepayment, main = "Private Insurance Payment", xlab =
"Private Insurance Payment ($)")

# Out-of-pocket payment

hist(new_mdata$outofpocketpayment, main = "Out-of-Pocket Payment", xlab = "Out-of-
Pocket Payment ($)")

# Uncollected liability

hist(new_mdata$uncollectedliability, main = "Uncollected Liability", xlab = "Uncollected
Liability ($)")

# Other payment

hist(new_mdata$otherpayment, main = "Other Payment", xlab = "Other Payment ($)")

#Numerical data is skewed to the right due to most of them having one heavy outlier

#Trying to eliminate skewness by using either log or sqrt forms

# +1 is added to handle the 0 and any possible negative values

```


#Sqrt was tested, but didn't handle the skewness issue so I stuck with taking the log of all the payments that were in \$

Dental event

```
new_mdata$log_dentalevent <- log(new_mdata$dentalevent + 1)
```

```
hist(new_mdata$log_dentalevent, main = "Log-Transformed Dental Event", xlab = "Log(Dental Event ($) + 1)")
```

Vision event

```
new_mdata$log_visionevent <- log(new_mdata$visionevent + 1)
```

```
hist(new_mdata$log_visionevent, main = "Log-Transformed Vision Event", xlab = "Log(Vision Event ($) + 1)")
```

Hearing event

```
new_mdata$log_hearingevent <- log(new_mdata$hearingevent + 1)
```

```
hist(new_mdata$log_hearingevent, main = "Log-Transformed Hearing Event", xlab = "Log(Hearing Event ($) + 1)")
```

Home health event

```
new_mdata$log_homehealthevent <- log(new_mdata$homehealthevent + 1)
```

```
hist(new_mdata$log_homehealthevent, main = "Log-Transformed Home Health Event", xlab = "Log(Home Health Event ($) + 1)")
```

Inpatient event

```
new_mdata$log_inpatientevent <- log(new_mdata$inpatientevent + 1)
```

```
hist(new_mdata$log_inpatientevent, main = "Log-Transformed Inpatient Event", xlab = "Log(Inpatient Event ($) + 1)")
```

Medical provider event

```
new_mdata$log_medicalproviderevent <- log(new_mdata$medicalproviderevent + 1)
```

```
hist(new_mdata$log_medicalproviderevent, main = "Log-Transformed Medical Provider Event", xlab = "Log(Medical Provider Event ($) + 1)")
```

Outpatient event

```
new_mdata$log_outpatientevent <- log(new_mdata$outpatientevent + 1)
```

```
hist(new_mdata$log_outpatientevent, main = "Log-Transformed Outpatient Event", xlab = "Log(Outpatient Event ($) + 1)")
```

Prescribe medicine

```
new_mdata$log_prescribemedicine <- log(new_mdata$prescribemedicine + 1)
```

```

hist(new_mdata$log_prescribemedicine, main = "Log-Transformed Prescribe Medicine",
xlab = "Log(Prescribe Medicine ($) + 1)")

# Total payment
new_mdata$log_totalpayment <- log(new_mdata$totalpayment + 1)

hist(new_mdata$log_totalpayment, main = "Log-Transformed Total Payment", xlab =
"Log(Total Payment ($) + 1)")

# Medicare payment
new_mdata$log_medicarepayment <- log(new_mdata$medicarepayment + 1)

hist(new_mdata$log_medicarepayment, main = "Log-Transformed Medicare Payment",
xlab = "Log(Medicare Payment ($) + 1)")

# Medicaid payment
new_mdata$log_medicaidpayment <- log(new_mdata$medicaidpayment + 1)

hist(new_mdata$log_medicaidpayment, main = "Log-Transformed Medicaid Payment",
xlab = "Log(Medicaid Payment ($) + 1)")

# Medicare Advantage payment
new_mdata$log_medicareadvantagepayment <-
log(new_mdata$medicareadvantagepayment + 1)

hist(new_mdata$log_medicareadvantagepayment, main = "Log-Transformed Medicare
Advantage Payment", xlab = "Log(Medicare Advantage Payment ($) + 1)")

# Private insurance payment
new_mdata$log_privateinsurancepayment <- log(new_mdata$privateinsurancepayment +
1)

hist(new_mdata$log_privateinsurancepayment, main = "Log-Transformed Private
Insurance Payment", xlab = "Log(Private Insurance Payment ($) + 1)")

# Out-of-pocket payment
new_mdata$log_outofpocketpayment <- log(new_mdata$outofpocketpayment + 1)

hist(new_mdata$log_outofpocketpayment, main = "Log-Transformed Out-of-Pocket
Payment", xlab = "Log(Out-of-Pocket Payment ($) + 1)")

# Uncollected liability
new_mdata$log_uncollectedliability <- log(new_mdata$uncollectedliability + 1)

hist(new_mdata$log_uncollectedliability, main = "Log-Transformed Uncollected
Liability", xlab = "Log(Uncollected Liability ($) + 1)")

# Other payment

```

```

new_mdata$otherpayment <- new_mdata$otherpayment +
abs(min(new_mdata$otherpayment)) + 1

new_mdata$log_otherpayment <- log(new_mdata$otherpayment)

new_mdata$ihs_otherpayment <- asinh(new_mdata$otherpayment)

hist(new_mdata$log_otherpayment, main = "Log-Transformed Other Payment", xlab =
"Log(Other Payment ($) + 1)")

#Extracting the column and rows separately from the dim() function
dims<-dim(new_mdata)
num_rows<-dims[1]
num_columns<-dims[2]
#Printing each one separately
print_rows<- print(paste0("Number of rows: ", num_rows))
print_columns<- print(paste0("Number of columns: ", num_columns))

#Trying to eliminate skewness by using either log or sqrt forms
# +1 is added to handle the 0 and any possible negative values
#Sqrt was tested, but didn't handle the skewness issue so I stuck with taking the log of all the
payments that were in $

# Dental event
new_mdata$log_dentalevent <- log(new_mdata$dentalevent + 1)
hist(new_mdata$log_dentalevent, main = "Log-Transformed Dental Event", xlab =
"Log(Dental Event ($) + 1)")

# Vision event
new_mdata$log_visionevent <- log(new_mdata$visionevent + 1)
hist(new_mdata$log_visionevent, main = "Log-Transformed Vision Event", xlab =
"Log(Vision Event ($) + 1)")

# Hearing event

```

```

new_mdata$log_hearingevent <- log(new_mdata$hearingevent + 1)

hist(new_mdata$log_hearingevent, main = "Log-Transformed Hearing Event", xlab =
"Log(Hearing Event ($) + 1)")

# Home health event

new_mdata$log_homehealthevent <- log(new_mdata$homehealthevent + 1)

hist(new_mdata$log_homehealthevent, main = "Log-Transformed Home Health
Event", xlab = "Log(Home Health Event ($) + 1)")

# Inpatient event

new_mdata$log_inpatientevent <- log(new_mdata$inpatientevent + 1)

hist(new_mdata$log_inpatientevent, main = "Log-Transformed Inpatient Event", xlab =
"Log(Inpatient Event ($) + 1)")

# Medical provider event

new_mdata$log_medicalproviderevent <- log(new_mdata$medicalproviderevent + 1)

hist(new_mdata$log_medicalproviderevent, main = "Log-Transformed Medical
Provider Event", xlab = "Log(Medical Provider Event ($) + 1)")

# Outpatient event

new_mdata$log_outpatientevent <- log(new_mdata$outpatientevent + 1)

hist(new_mdata$log_outpatientevent, main = "Log-Transformed Outpatient Event",
xlab = "Log(Outpatient Event ($) + 1)")

# Prescribe medicine

new_mdata$log_prescribemedicine <- log(new_mdata$prescribemedicine + 1)

hist(new_mdata$log_prescribemedicine, main = "Log-Transformed Prescribe
Medicine", xlab = "Log(Prescribe Medicine ($) + 1)")

# Total payment

new_mdata$log_totalpayment <- log(new_mdata$totalpayment + 1)

hist(new_mdata$log_totalpayment, main = "Log-Transformed Total Payment", xlab =
"Log(Total Payment ($) + 1)")

# Medicare payment

new_mdata$log_medicarepayment <- log(new_mdata$medicarepayment + 1)

hist(new_mdata$log_medicarepayment, main = "Log-Transformed Medicare
Payment", xlab = "Log(Medicare Payment ($) + 1)")

# Medicaid payment

new_mdata$log_medicaidpayment <- log(new_mdata$medicaidpayment + 1)

```

```

hist(new_mdata$log_medicaidpayment, main = "Log-Transformed Medicaid
Payment", xlab = "Log(Medicaid Payment ($) + 1)")

# Medicare Advantage payment

new_mdata$log_medicareadvantagepayment <-
log(new_mdata$medicareadvantagepayment + 1)

hist(new_mdata$log_medicareadvantagepayment, main = "Log-Transformed Medicare
Advantage Payment", xlab = "Log(Medicare Advantage Payment ($) + 1)")

# Private insurance payment

new_mdata$log_privateinsurancepayment <-
log(new_mdata$privateinsurancepayment + 1)

hist(new_mdata$log_privateinsurancepayment, main = "Log-Transformed Private
Insurance Payment", xlab = "Log(Private Insurance Payment ($) + 1)")

# Out-of-pocket payment

new_mdata$log_outofpocketpayment <- log(new_mdata$outofpocketpayment + 1)

hist(new_mdata$log_outofpocketpayment, main = "Log-Transformed Out-of-Pocket
Payment", xlab = "Log(Out-of-Pocket Payment ($) + 1)")

# Uncollected liability

new_mdata$log_uncollectedliability <- log(new_mdata$uncollectedliability + 1)

hist(new_mdata$log_uncollectedliability, main = "Log-Transformed Uncollected
Liability", xlab = "Log(Uncollected Liability ($) + 1)")

# Other payment

new_mdata$otherpayment <- new_mdata$otherpayment +
abs(min(new_mdata$otherpayment)) + 1

new_mdata$log_otherpayment <- log(new_mdata$otherpayment)

new_mdata$lhs_otherpayment <- asinh(new_mdata$otherpayment)

hist(new_mdata$log_otherpayment, main = "Log-Transformed Other Payment", xlab =
"Log(Other Payment ($) + 1)")

```

```

#Testing for linearity of total payment vs each variable again after changes
par(mfrow = c(4, 4))
dependent_var <- "totalpayment"
independent_vars <- setdiff(names(new_mcddata), dependent_var)

# Create a function to plot linearity
par(mfrow = c(4,4))
plot_linearity <- function(x, y, plot_title) {
  plot(y ~ x, data = new_mcddata, main = plot_title, xlab = deparse(substitute(x)), ylab =
deparse(substitute(y)))
  abline(lm(y ~ x, data = new_mcddata), col = "red")
}
# Loop through each independent variable and plot linearity
for (var in independent_vars) {
  plot_title <- paste(dependent_var, "vs", var)
  plot_linearity(new_mcddata[[var]], new_mcddata[[dependent_var]], plot_title)
}

```

#Regression Model

```
# ols model
```

```
# Variables age_g3, otherrace and incomemoreequal25 have been removed due to perfect
collinearity with their opposite variables
```

```

ols1 <- lm(formula = totalpayment ~ sex + numberchroniccond + log_dentalevent+
log_visionevent+ log_hearingevent+
          log_medicalproviderevent +log_outpatientevent + log_prescribemedicine+
log_medicarepayment+

```

```
log_medicaidpayment+ log_medicareadvantagepayment+
log_privateinsurancepayment + log_outofpocketpayment+
```

```
log_uncollectedliability+log_otherpayment+age_g1+age_g2+hisp+nonhispblack+nonhispwhite+
incomeless25
```

```
, data = new_mcddata)
```

```
summary(ols1)
```

```
r1<-residuals(ols1)
```

```
qqnorm(r1)
```

```
qqline(r1)
```

```
ols2 <- lm(formula = totalpayment ~ numberchroniccond + log_visionevent +
log_medicalproviderevent +
```

```
log_outpatientevent + log_prescribemedicine + log_medicarepayment +
log_medicaidpayment +
```

```
log_medicareadvantagepayment + log_privateinsurancepayment +
log_outofpocketpayment +
```

```
log_uncollectedliability + log_otherpayment + nonhispblack, data = new_mcddata)
```

```
summary(ols2)
```

```
r2<-residuals(ols2)
```

```
qqnorm(r2)
```

```
qqline(r2)
```

```
library("MASS")
```

```
new_mcddata$totalpayment <- new_mcddata$totalpayment + 1
```

```
bc <- boxcox(ols2, lambda = seq(-3, 3))
```

```
best.lam = bc$x[which(bc$y==max(bc$y))]
```

```

# Update the model formula with the best lambda value

fullmodel.inv <- lm((totalpayment^best.lam) ~ numberchroniccond + log_visionevent +
log_medicalproviderevent +
log_outpatientevent + log_prescribemedicine + log_medicarepayment +
log_medicaidpayment +
log_medicareadvantagepayment + log_privateinsurancepayment +
log_outofpocketpayment +
log_uncollectedliability + log_otherpayment + nonhispblack, data =
new_mcddata)

summary(fullmodel.inv)

r3<-residuals(ols3)

qqnorm(r3)

qqline(r3)

stargazer(ols1,ols2, ols3, type = "text")

```