

Additional Work

Model Selection

Because our model is focused on prediction rather than causal inference, we decided to undergo a rigorous variable selection process. After removing variables deemed too severely multicollinear, we're left with a “full” model consisting of 21 explanatory variables.

The reason we chose to screen our variables with VIF beforehand is that removing explanatory variables that are collinear not only helps with the assumptions of linear regression, but also helps computationally, as we have less variables to search through when searching for the best model.

Screening With LASSO

Before performing best subsets regression, we decided to run LASSO on our model in order to get a sense of variable importance in a predictive context.

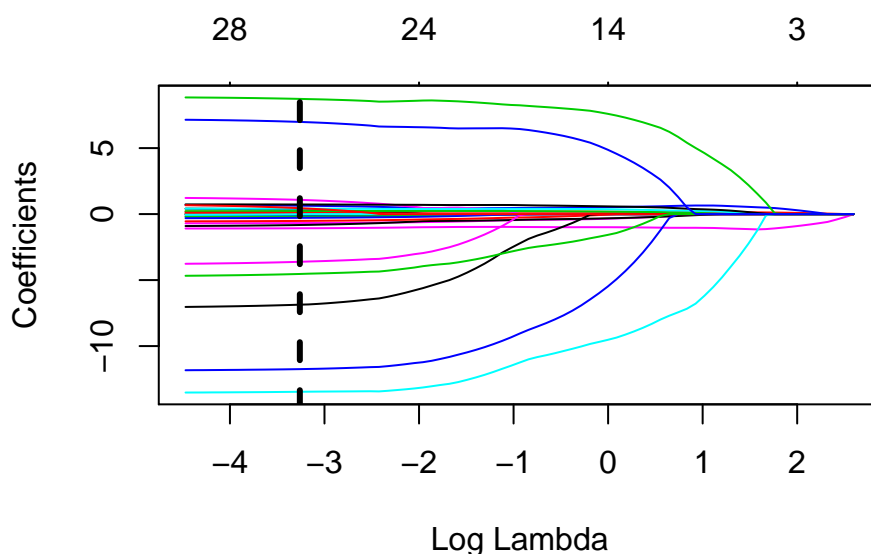


Figure 1: Lasso coefficient trails. The dotted line marks the optimal lambda.

Taking a look at the results of LASSO, we see that none of our coefficients have been zeroed out, meaning that we will need to take a look at other variable selection methods if we want to shrink our model. As a result, we explore a different method of model shrinkage: best subsets regression.

Best Subsets Regression

Best subsets regression exhaustively searches every combination of variables for every possible model size and selects the best models for each model size according to different criteria. The criteria we considered were Adjusted R^2 , Mallows' C_p , and BIC. We chose these three criteria since they're supported by the `R` function `regsubsets`, and we wanted to use the same library for the sake of consistency in model selection.

	Adjusted R^2	Mallows' C_p	BIC
Number of Variables	24	22	16

Table 1: The number of variables in the “best” model as chosen by various criteria. Note that the number of variables has increased due to dummy variables being added to the model.

As seen in **Table 1**, Adjusted R^2 as our criterion resulted in the largest model, while BIC as our criterion resulted in the smallest model. Taking a closer look at the actual models that were selected, we see that

some of our dummy variables for our only categorical variable, region, ended up being dropped by best subset regression. Because it's not possible to write a formula that drops some of these dummy variables as well as the fact that the majority of dummy variables were kept for all 3 models, we chose to keep Division in all 3 of our models even if some of the dummy variables ended being dropped. This isn't too consequential as in the Adjusted R^2 model and the Mallows' Cp Model, only the dummy variable associated with the South Atlantic division is dropped, while in the BIC model, only the dummy variables associated with the South Atlantic division and New England division are dropped.

Cross Validation

After creating our models (2 distinct ones in this case), it's clear the the criteria don't agree on which model is the best. In order to assess the performance of our models, we need to evaluate the predictive ability of our models on data they have never seen before. Rather than using a train-test split of our data, we decided to use cross validation since cross validation tends to smooth out noise or randomness, and also provides more precision while reducing bias as we have more data for fitting the models. Leave-one-out CV is too computationally expensive due to the large number of rows, so we went with k-fold CV instead with a fold size of 10. We also computed the MSE from CV for the full model as well as the LASSO model to serve as comparisons.

	Full Model	LASSO Model	Adjusted R^2 Model	Mallow's Cp Model	BIC Model
MSE	358	358	358	357	359

Table 2: The MSE from k-fold CV of our various models. Note that the MSE of our full model and LASSO model are the same since the two models are the same (albeit it's definitely possible for two different models to have the same MSE).

As seen in **Table 2**, the model with the lowest MSE ended up being our Mallows' Cp model, so we'll go ahead and choose that model as our "final" model for this step.

Model Diagnostics

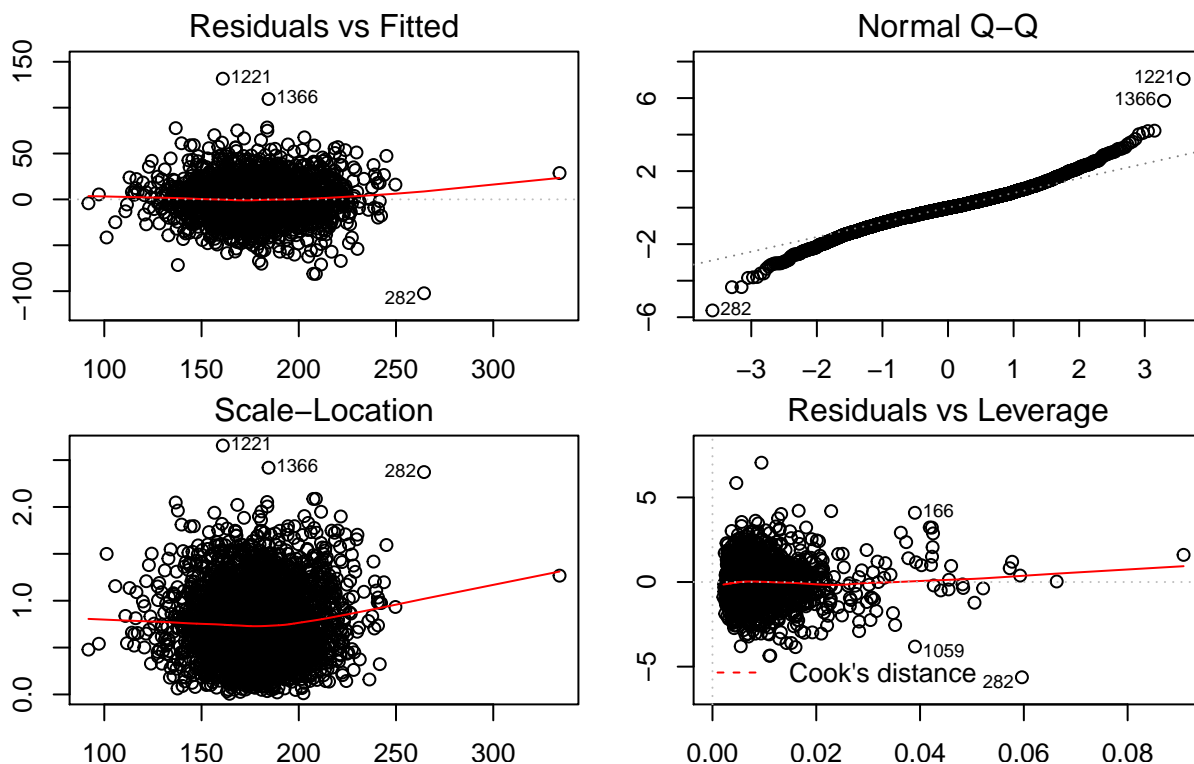


Figure 2: Diagnostic plots of the chosen model.

We notice an outlier in our residual plots in **Figure 2** that reveal a point with a somewhat high leverage. After investigating the possibility of an encoding error, we discovered that this point belonged to Union County, Florida which is known to have a disproportionately high cancer death rate compared to the rest of the United States, so we left that data point in. Something that was concerning during EDA was that a few of our explanatory variables didn't have normal distributions. When we tried applying a Box-Cox transformation to them, our variable selection process ended up not dropping any variables and had a high R^2 , but our model diagnostics revealed that assumptions of linear regression were being broken due to residuals not being randomly scattered about the fitted line. Applying a Box-Tidwell transformation at this stage of model selection was considered, but due to powers being pushed to infinity during iterations and being unable to diagnose this issue as Box-Tidwell wasn't covered in class, we decided to not continue pursuing this particular transformation. The diagnostics of the model reveal that the assumptions of linear regression are mostly followed anyways, so a transformation wouldn't necessarily create a huge improvement.

Adding States to the Model

While we have region as one of the variables in our model, it's possible that certain states may go against the trend of the region. As a result, we will consider adding states as variables to our model. Doing so will allow the coefficient of a state to "counteract" the coefficient of its region in the event that a state is significantly different than its region. In order to decide which states to add to our model, we will use forward selection using AIC and BIC as our criteria. We chose to do forward selection rather than best subsets regression here due to the large number of additional columns we have added via one hot encoding the state variable. We chose AIC and BIC as our criteria since they're supported by the step function and we want to use the same library for the sake of consistency during model selection.

The results of our forward selection reveal that adding states does in fact add precision to our model - AIC adds 20 states to our model and BIC adds 7 states to our model. AIC adds significantly more variables than BIC, though that's not surprising considering that BIC penalizes model complexity more heavily. The states chosen by both BIC and AIC tend to be in the Southern and Midwest regions of the United States, perhaps revealing that these regions contain many outlier states. In order to determine which model fits the data better, we again ran k-fold cross validation (with a fold size of 10) on these two models and because the AIC model had a lower MSE, we chose the AIC model as our "final" model for this stage of the model selection process. Our model now has 43 total variables (counting the dummy variables for Division as separate variables) with the addition of the 20 state variables.

Exploring Interactions