# Optimization Project 3:
# Feature Selection using MIQP vs Lasso

## Group Members:

Krish Engineer (ke7466)
Rudraksh Garg (rg44778)
Ryan Lee (rl32227)
Samarth Mishra (sm79247)

## Purpose

One of the most common problems in predictive analytics is variable selection for regression. Direct variable selection using optimization has long been dismissed by the statistics/analytics community because of computational difficulties. This computational issue was part of the motivation for the development of LASSO and ridge regression. However, in the recent past there have been tremendous advancements in optimization software, specifically the ability to solve mixed integer quadratic programs (MIQP) in order to find the best set of coefficients for regression.

## Objective

In this study, we tested the feasibility of direct variable selection using two approaches:
1. Mixed Integer Quadratic Programming (MIQP)
2. LASSO Regression

The MIQP variable selection problem for regression will be solved using Gurobi and its results will be compared to the results given by LASSO regression. The comparison of these results will then be used to advise management on which variable selection method should be preferred.

## Approach 1 : Direct Variable Selection - MIQP

### Objective Function

This problem can be structured as an optimization problem where our objective is to minimize the squared loss incurred during regression. So our objective becomes minimizing the sum of squared errors as specified below:

$$\min_{\beta, z} \sum_{i=1}^{n} (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_m x_{im} - y_i)^2$$

### Objective Function Manipulation

$\beta$ can be constructed as a vector of size $m+1$. This is $m+1$ because we have $m$ variables and 1 intercept ($\beta_0$). X can be constructed as an n by ($m+1$) matrix where the first column is all ones and the second

column onwards has the data from each of the n records. The target vector would be the output from each of the n records. We can then construct a vector of size n using the following:

$$(X\beta - y)$$

This looks similar to our objective function. Since we are using the squared errors, we would have to take the following instead:

$$(X\beta - y)^T * (X\beta - y)$$

We can simplify this for Gurobi and pose our objective function as a MIQP problem as given below:

$$\min_{\beta,z} \beta^T (X^T X)\beta + (-2\, y^T X)\beta$$

With this setup, while β may look as being of size m+1, it actually needs to be of dimension 2m+1 as there are m+1 additional binary variables $z_j$'s, that represent if those β variables should be included in the regression or not.

Given below are the Quadratic (Q) and Linear Term matrices (c) for this MIQP problem:

Quadratic Matrix (Q)

$Q = X^T X$ (Dimension: 2m+1 * 2m+1)
Here, only the upper left corner has non-zero values equal to $X^T X$ while all other values are zero (matrix manipulated to keep consistent in the Gurobi format )

Linear Term Matrix (c)

$c = -2y^T X$ (Dimension: 2m+1 * 1)

Here, only the first m components have non-zero values equal to $-2y^T X$ while while all other values are zero (matrix manipulated to keep consistent in the Gurobi format )

## Code Snippet for computing the Q and c matrices

Defining functions to compute Q and c matrices for the MIQP

```
1  ## Creating a function that outputs the Q matrix
2  def create_qmatrix(x_arr,m):
3      rows = 2*m+1
4      cols = 2*m+1
5      q_matrix = np.zeros((rows,cols))
6      q_matrix[:m+1,:m+1] = x_arr.T @ x_arr
7      return q_matrix
8
9  ## Creating a function that outputs the linear term
10 def create_linearTerm(x_arr,y_arr,m):
11     t2 = (-2)* y_arr.T @ x_arr
12     linear_term = np.pad(t2, (0, m), 'constant')
13     return linear_term
```

## Decision Variables:

$\beta_i$: Beta coefficients for each of the predictor variables in linear regression

$Z_j$: The z variables tell us whether a beta gets picked or not during variable selection

**Constraints**

1. Big-M Constraints:

   The z variables tell us whether a beta gets picked or not. Thus, a beta can exist only if z is non-zero.

   Using "big-M" constraints, we can set a max that all coefficients could possibly be:

   $$s.t. -Mz_j \leq \beta_j \leq Mz_j \quad for\ j = 1, 2, 3, \dots, m$$

   Big M value taken in our MIQP model=100 (sufficiently large)

2. Number of variables to be chosen (k)

   We also need to specify the max number of variables that need to be selected. We would need binary variables that would act as switches that would state if a variable would be added or not. This can be framed using the following constraint where k is the maximum number of variables for our regression:

   $$\sum_{j=1}^{m} z_j \leq k$$
   $$z_j \ are\ binary$$

## Code Snippet for Computing the Beta Coefficients

```python
1  def get_betas(q_matrix, linear_term, k, M, m):
2
3  #   Setting up A Matrix
4
5      A = np.zeros((2*m+1,2*m+1)) # Inititalizing A matrix
6
7      A[0,-m:] = np.ones(m) # Constraint on z variables
8
9      # Big M Constraints
10     j=0
11     for i in range(1,m+1):
12         A[i,i] = 1
13         A[i,-m+j] = -M
14         j = j+1
15     d=0
16     for a,c in zip(range(i+1,2*m+1), range(1,m+1)):
17         A[a,c] = 1
18         print()
19         A[a,-m+d] = M
20         d = d + 1
21
22     # Setting up b matrix
23
24     b = [0]*(2*m+1)
25     b[0] = k
26     b_arr = np.array(b)
27
28     # Setting up sense matrix
29     sense = ['<']*(m+1)
30     sense.extend(['>']*(m))
31     sense_arr = np.array(sense)
32
33     # Running the optimzation model to compute beta coefficients
34     model = gp.Model()
35     mod_x = model.addMVar(2*m+1,vtype=np.array(['C']*(m+1) + ['B']*m), lb=np.array([np.NINF]+[-M]*(m)+ [0]*m))
36     mod_con = model.addMConstr(A, mod_x, sense_arr, b_arr)
37
38     #setting obj
39     model.setMObjective(q_matrix,linear_term,0,sense=gp.GRB.MINIMIZE)
40     model.Params.TimeLimit = 3600
41     #model.setParam('IntFeasTol', 1e-9) # binary var can be non zero
42     model.Params.OutputFlag = 0
43     model.optimize()
44     my_coeffs = mod_x.x
45
46     return (model.objVal, my_coeffs)
```

# Results

## Running the MIQP Model for different values of k with 10-fold CV

To solve the MIQP problem, we first started off by defining the cross-validation function to split the training data into 10 folds. The function was set to try values of k ranging from 5 to 50 since there were 50 variables to select from. We then set up the optimization problem through Gurobi in the format described in the Direct Variable Selection section of this report. Our objective was to find out the value of k that gives out the least Cross Validation error.

## Running the MIQP model for k=10 on entire training data

```
1  m = train_df.iloc[:,1:].shape[1]
2  Xtrain = train_df.iloc[:,1:]
3  idx = 0
4  n = Xtrain.shape[0]
5  new_col = [1]* n
6  Xtrain.insert(loc=idx, column='X0', value=new_col)
7  x_arr = Xtrain.to_numpy()
8
9  # Computing Q matrix using entire training data for k=10
10 q_matrix = create_qmatrix(x_arr,m)
11
12 y_arr = train_df['y'].to_numpy()
13
14 # Computing linear term using entire training data for k=10
15 linear_term = create_linearTerm(x_arr,y_arr,m)
16
17 # Computing Beta coefficients using entire training data for k=10
18 coeffs = get_betas(q_matrix, linear_term, 10,100,m)[1]
```

Given below are the Cross Validation error that we obtained for different values of k:

|   | Value of k | MSE |
|---|---|---|
| 0 | 5 | 3.959193 |
| 1 | 10 | 2.805040 |
| 2 | 15 | 3.134713 |
| 3 | 20 | 3.178051 |
| 4 | 25 | 3.451272 |
| 5 | 30 | 3.368715 |
| 6 | 35 | 3.320179 |
| 7 | 40 | 3.467398 |
| 8 | 45 | 3.465303 |
| 9 | 50 | 3.415552 |

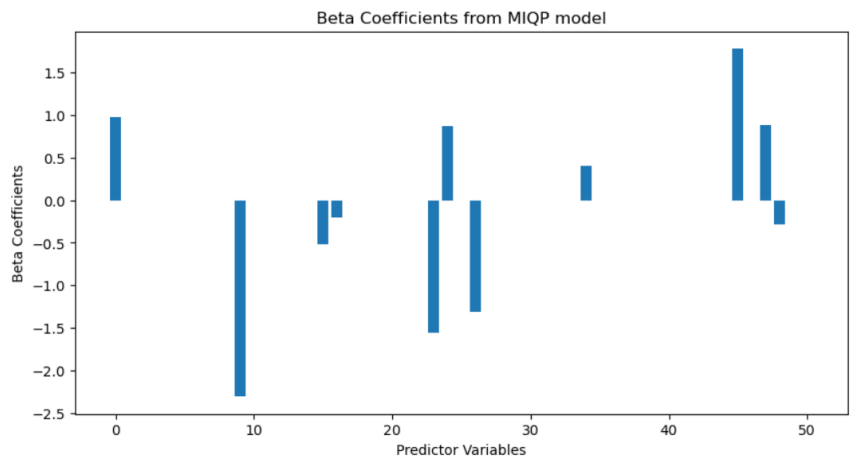**We see that the least Cross Validation error was obtained for k=10.**

Hence we chose k=10 for our final MIQP model and ran the MIQP model on the entire training dataset. We determined the $\beta$'s for each feature and then used it to predict the y values in the test set.

Given below are the Training and Test MSE's obtained:
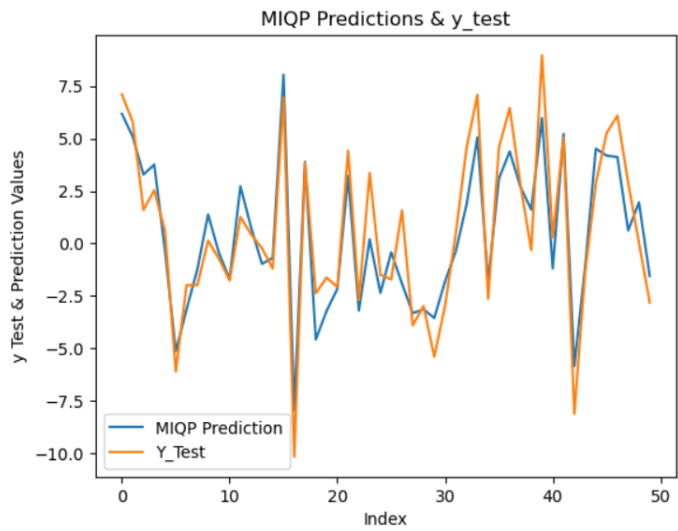
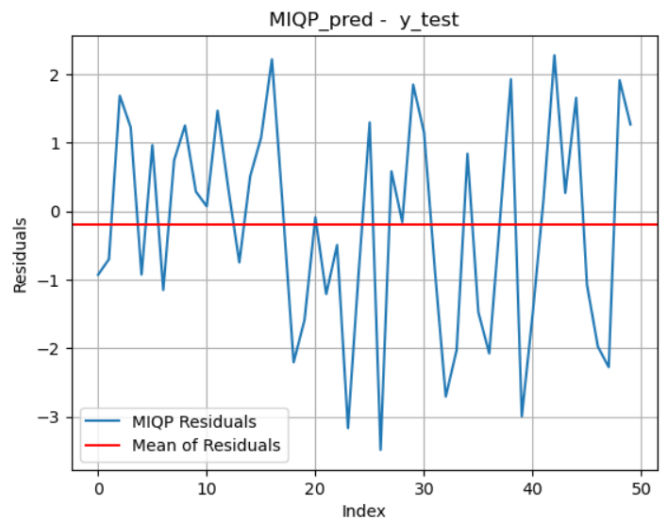**Training MSE obtained: 2.392**
**Test MSE obtained: 2.337**

**Beta Coefficients Obtained for MIQP**



**MIQP Predictions vs Actual y values (Test Set)**



**MIQP Prediction Residuals (Test Set)**

## Approach 2 : Indirect Variable Selection using LASSO

We applied the lasso regularization to a linear regression model trained on our training set to compare the outcomes of our MIQP with those of other indirect feature selection techniques. By including the absolute sum of coefficients penalty in the regression loss function, Lasso aids in the reduction of the regression coefficients. The indirect variable selection or LASSO Regularization method is posed as

$$\min_\beta \sum_{i=1}^{n} ( \beta_0 + \beta_1 x_{i1} + \ldots\ldots\ldots + \beta \; x_i \; - y_i)^2 + \lambda \sum_{j=1}^{m} |\beta|$$

Where $\lambda$ is a hyperparameter that can be chosen during cross validation. The LASSO regularization model has the advantage of 'shrinking' the $\beta s$ closer to zero, leading to variance reduction i.e. prevents overfitting. Furthermore, if $\lambda$ is large enough, several values of $\beta$ will be forced to be equal to 0.

There are more and more variables with zero coefficients as the penalty term grows. The penalty term is taken into consideration using the regularization parameter. This is how we can employ the Lasso to choose our features in an indirect manner.

# Results

## Running Lasso Regression doing 10-Fold CV to determine optimal $\lambda$

We used Lasso CV to do a 10-fold cross validation on the training set and determined the optimal value of $\lambda$.

**Code Snippet for Lasso regression using 10-fold CV**

```
1  from sklearn.linear_model import LassoCV
2
3  lasso_model = LassoCV(cv=10,random_state=0,normalize=True)
4  lasso_model.fit(X_train, y_train)
```
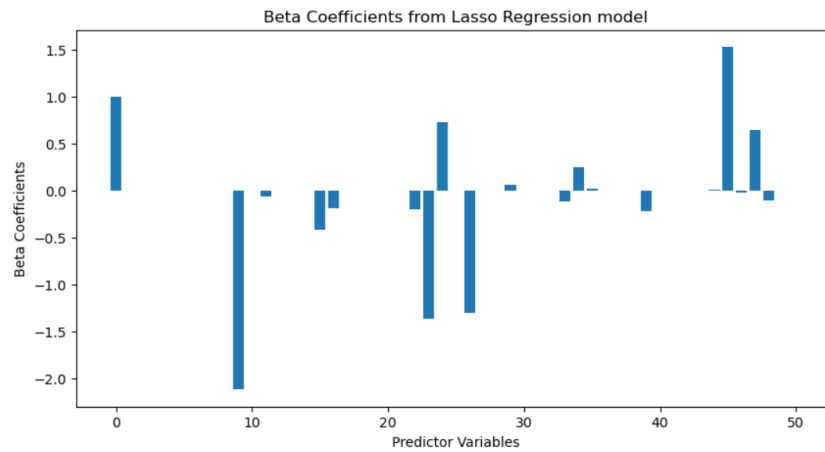
**Optimal lambda obtained: 0.0057**

We used this value of $\lambda$ to fit the LASSO model on the entire training set and determined Beta coefficients. Using the Beta coefficients, we determined the Lasso prediction on the test set. Given below are the Training and Test MSE's obtained:
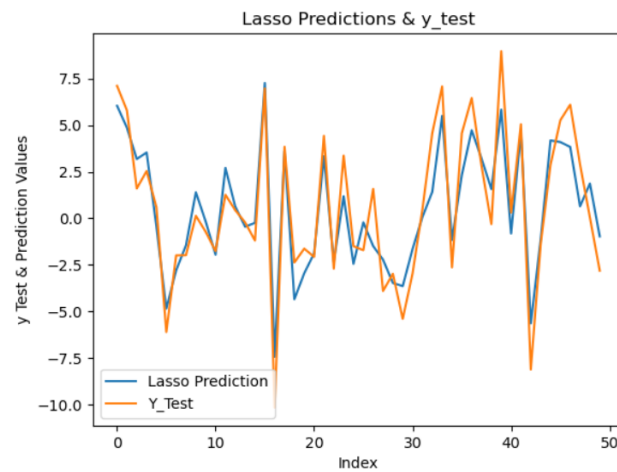
**Training MSE obtained: 2.401**
**Test MSE obtained: 2.360**
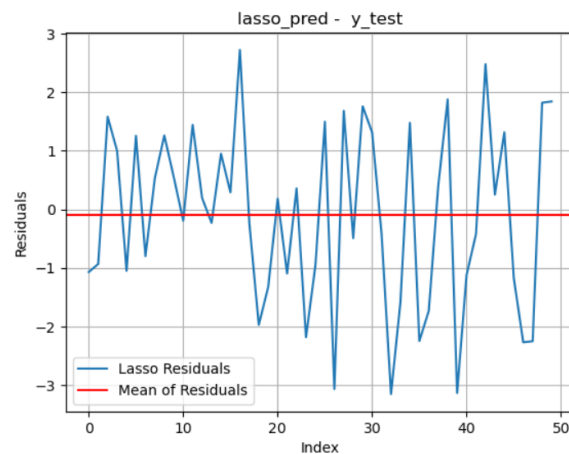
**Beta Coefficients Obtained for Lasso Model**



We see that Lasso model selects 18 out of the 50 predictor variables, more than what we chose through our most optimal MIQP model

**Lasso Predictions vs Actual y values (Test Set)**



**Lasso Prediction Residuals (Test Set)**

## Recommendation

We are in a better position to use specialized optimization techniques as opposed to only depending on easily accessible ML libraries as computation resources become more freely available. Due to the assumptions they make for faster convergence in their internal optimizations, the existing modules and packages may have some limitations. In our use-case, our unique optimization can produce a model that is less complex and has a higher holdout sample MSE. A slightly more reliable model is created by our MIQP approach not taking into account LASSO features that only had coefficient values close to zero in the lasso regression.

However, there are advantages and disadvantages to each approach depending on the use case. This must be taken into consideration when deciding which one to pick. With this in mind, the company may choose any approach based on the needs of the business case.

Finding the best value of k computationally takes much longer for the direct variable selection technique (MIQP model) than it does for the LASSO model. The direct variable selection model might be challenging in instances where there are many features to choose from. However, this strategy performs only slightly better. In the workplace, this model can be run overnight or on a cloud server and provide results without affecting productive hours, however that still may be costly in terms of time and money. It is also worth noting that MIQP chose less variables than LASSO which might be considered an advantage as there is less data dependency and lower complexity in the model.

Because there is only a very small difference in MSE and mean residuals between our MIQP and LASSO, in this situation, with the given data, it might not be worth the extra training time for the potentially insignificant improvement in accuracy. However, as there is still a tradeoff between accuracy and computation cost, the decision between the two should be determined by the application and use cases of these variable selection methods.

In contrast to the MIQP model, LASSO regression is a built-in Python feature that can be quickly accessible and utilized with few lines of code. For clients and stakeholders who might not be familiar with conventional optimization techniques, the LASSO model's outcomes are also simpler to explain and performs almost as well as MIQP in terms of MSE. Out of regularization and feature selection models, LASSO is more popular as it balances accuracy, time, and cost much better than direct variable selection through MIQP.