

Predicting costs

Chris Eshleman

10/20/2019

Q: What's the cross-validation being used? Q: Can I relax the penalty I'm applying via lasso? YES - THERE'S SOMETHING TO THIS - Flesh it out!

The agency estimates project costs internally. It may be able to use statistics to improve those predictions.

Model variation in bids, withholding the two most recent quarters of observations. Then test those models on the withheld observations. How well do they help with predictions?

Methods and Data.

1. Basic (manually constructed linear model),
2. Lasso (penalized regression - machine learning).

The project-level data comes from internal agency cost estimation.

The economic data is just quarterly stuff from the usual suspects and is specific to national and regional economic and labor market conditions.

Load some programming tools that are commonly used for this analysis. Not all of these packages will be used, and at some point it'll be worth backing up and cleaning the list.

Overview.

This effort bridges exploration of the agency's data with potential predictors from exogenous sources. The key idea is that if economic indicators (numbers) can add measurable value to the agency's cost estimation methods, it can help set the agency up for better informed next steps. Those next steps are yet to be defined.

#![Here's what we have in mind.](Estimator_Data.png)

The measures of "accuracy" are, for now, bivariate correlation. We discuss below why it's not yet time for fancier metrics, but the process above represents a first cut at trying to add a little more statistical value to the process.

Data comes from a few sources and needs to be merged. Two usual suspects are the agency's internal project information and a set of economic indicators specific to Greater New York (18 counties on both sides of the Hudson River).

Add economic data.

Add permits and steel prices.

The City of New York's database on permitting covers commercial and residential activity.

I'd like to ask have robust data on the rest of the region, including (and namely) Jersey City, but it's weaker than the City's, which by itself is a decent barometer of construction activity in greater New York.

Prices of construction materials and labor also figure into the agency's internal cost estimation, and I'll use steel prices for now. Future models might try and rope in other pricing data points, but the estimators are generally already taking prices into account when setting their numbers so this is likely of second- and third-order importance but the modeling selection algorithms may suggest using them.

The Engineering News-Record tracks and aggregates construction cost data through an index. Use that to cover prices.

```
cci = read.csv("./cci.csv")
names(cci) = tolower(names(cci))
cci$avg. = NULL
cci = gather(cci, month, cci, jan:dec, factor_key=TRUE)
cci$month = gsub("(^[[:alpha:]])", "\\U\\1", cci$month, perl=TRUE)
cci$month = as.Date(paste(cci$month, "01", cci$year, sep="-"), format="%b-%d-%Y")
cci = cci[month(cci$month) == 2 | month(cci$month) == 5 | month(cci$month) == 8 | month(cci$month) == 1, ]

cci$Quarter = paste("Q", quarter(cci$month), sep="")
cci$Q = paste(year(cci$month), cci$Quarter, sep="-")
cci$Quarter = NULL
cci$month = NULL
cci$year=NULL

bids = merge(bids, cci, by = "Q", all.x=TRUE)

rm(cci)
```

Munge

Stats software treats different variables in different ways depending on individual formatting. It's worth taking a look at the data structure.

I'll need to tell the software to reformat some of the variables, namely the economic indicators.

One variable of interest is the bidding process. Institutional discussions and earlier modeling suggests the bidding process may influence the bids. Limits placed on the range of bidders, for example, could, on average and holding other things constant, increase the average (and lowest qualifying) bid - this is basic microeconomics. I'll simplify the bidding format variable by making it binary: "public" for projects without significant constraints and "other" for ones, such as projects closed to firms not deemed "small business enterprises," that aren't. First I'll clean it a bit to consolidate near-duplicate categories.

There's room to also eventually include the names (anonymized is fine) of each project estimator to help modeling. Past work has suggested there isn't major causal variation between estimators — they generally do a pretty equivalent job in estimating bids. But having their names included nonetheless may prove to offer some control value. We can leave that to future modeling.

Restating objective.

We want to understand whether and how we might help the agency estimate the actual cost of a project. That's invariably going to be represented by the low qualifying bid, and our starting point is the estimate coming from the Engineering Department. From here on we'll define "accuracy" as the ratio of dollars estimated over dollars bid. So a "1" would mean the engineering team nailed it, a "0.94" would mean they estimated 94 cents for every 1 dollar in the low bid, et cetera.

If there are outliers in there - projects that, for an unexplainable reason, was way off, consider removing it.

We may want to think of project size categorically.

```
#bids$decile = decile(vector = bids$Second.Bid)
#bids$decile = ordered(bids$decile, levels = 1:10)
```

```

bids = bids %>%
  mutate(quantile = ntile(as.numeric(Second.Bid), 10))

## Warning: `as_dictionary()` is soft-deprecated as of rlang 0.3.0.
## Please use `as_data_pronoun()` instead
## This warning is displayed once per session.

## Warning: `new_overscope()` is soft-deprecated as of rlang 0.2.0.
## Please use `new_data_mask()` instead
## This warning is displayed once per session.

## Warning: The `parent` argument of `new_data_mask()` is deprecated.
## The parent of the data mask is determined from either:
##
## * The `env` argument of `eval_tidy()`
## * Quosure environments when applicable
## This warning is displayed once per session.

## Warning: `overscope_clean()` is soft-deprecated as of rlang 0.2.0.
## This warning is displayed once per session.

Back up data.
bid = bids # backup my data frame

```

Analysis

Now the data is prepped. So split it into the training / test sets we talked about at the start. The bids start in 2015.

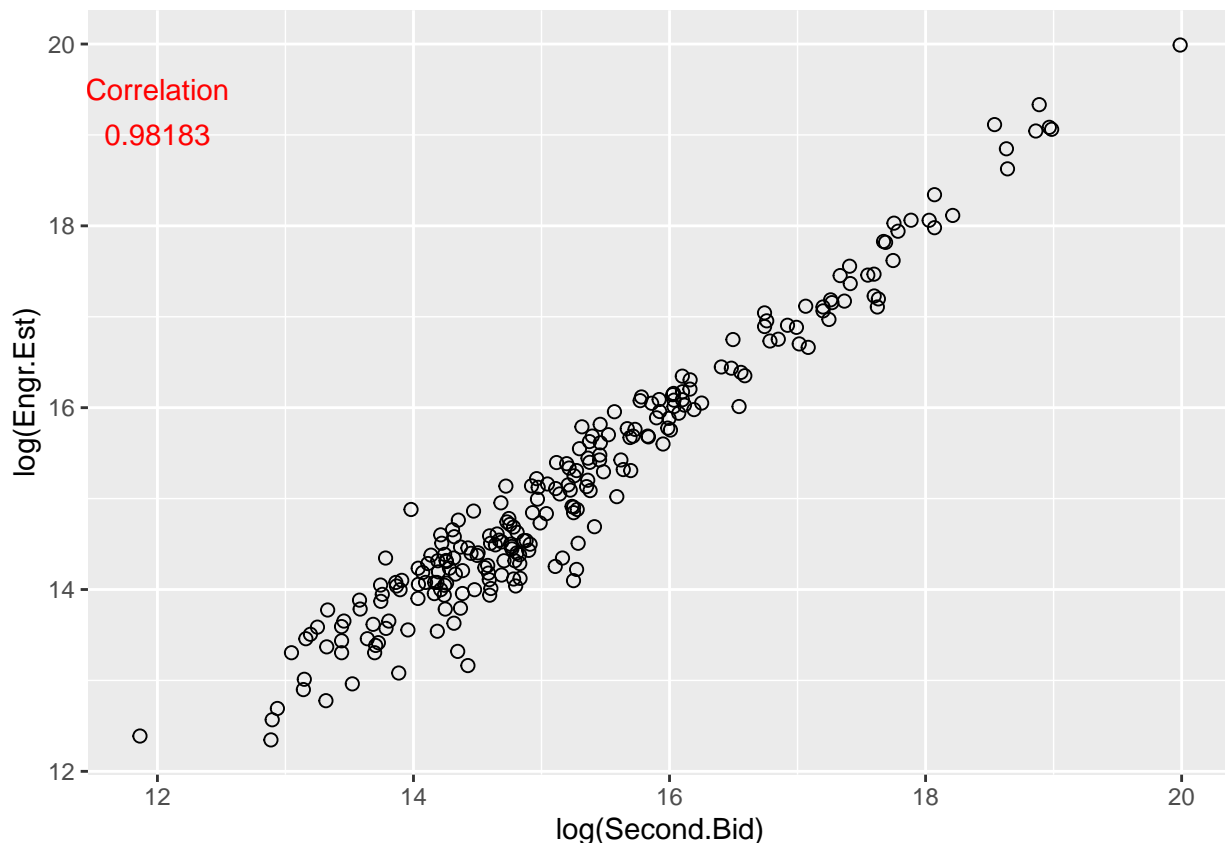
Motivation.

Why do we think developing a controlled multivariate (complicated) model will be worth it? Well, the average gap is \$2.5 million, or 20%, off of our estimates. What's the raw (uncontrolled) bivariate relationship between engineering estimates and low bids?

```

##
## Call:
## lm(formula = bids$Second.Bid ~ bids$Engr.Est)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -65502890  -996757   -514372   419628  53757694
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  9.428e+05  5.319e+05   1.772  0.0776 .
## bids$Engr.Est 8.855e-01  1.112e-02  79.663 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7762000 on 237 degrees of freedom
## Multiple R-squared:  0.964, Adjusted R-squared:  0.9638
## F-statistic: 6346 on 1 and 237 DF, p-value: < 2.2e-16

```



(The logarithm treatment is just to distribute it across the plot (one of the observations is an outlier).)

The in-house engineers' guesses predict more than 98% of the variation in low bids.

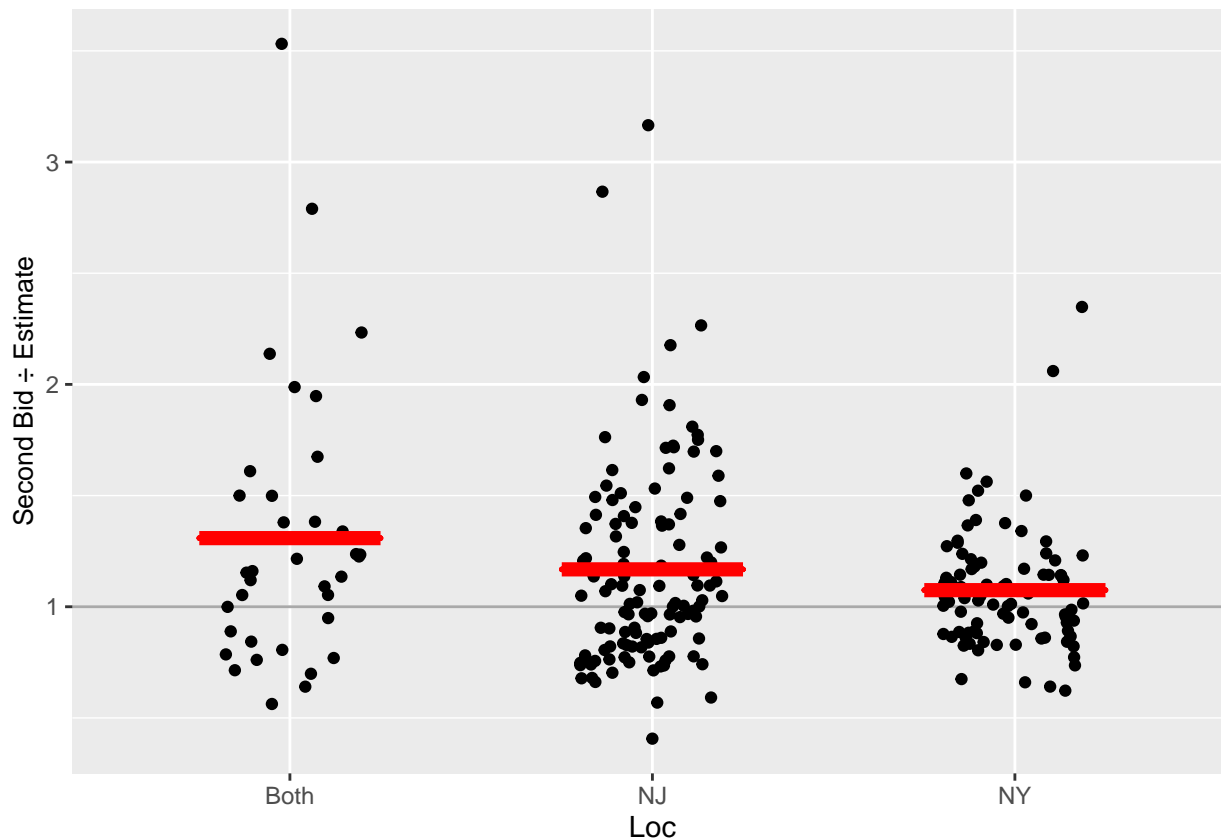
Some of the remaining variation can be explained with some guesswork. Calculate summary stats for key variables - this will help later:

```
loc.summary = aggregate(accuracy ~ Loc, mean, data=na.omit(bids))
type.summary = aggregate(accuracy ~ Typeology, mean, data=na.omit(bids))
format.summary = aggregate(accuracy ~ Format, mean, data=na.omit(bids))
decile.summary = aggregate(accuracy ~ quantile, mean, data=na.omit(bids))
ld.summary = aggregate(accuracy ~ LD, mean, data=na.omit(bids))
```

Location likely has some predictive power that engineers may not be able to capture or fully predict. Basically, projects that span the Hudson River wind up costing more, on average, than ones plunked squarely in either New York or New Jersey. What's the raw (uncontrolled) relationship between bid accuracy and location?

```
##
## Call:
## lm(formula = bids$accuracy ~ bids$Loc)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.76135 -0.25572 -0.07396  0.18784  2.22307
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.30860    0.07132  18.349 < 2e-16 ***
## bids$LocNJ   -0.14044    0.08147  -1.724  0.08606 .
## bids$LocNY   -0.23401    0.08509  -2.750  0.00642 **
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4279 on 236 degrees of freedom
## Multiple R-squared:  0.0318, Adjusted R-squared:  0.02359
## F-statistic: 3.875 on 2 and 236 DF,  p-value: 0.02208
```



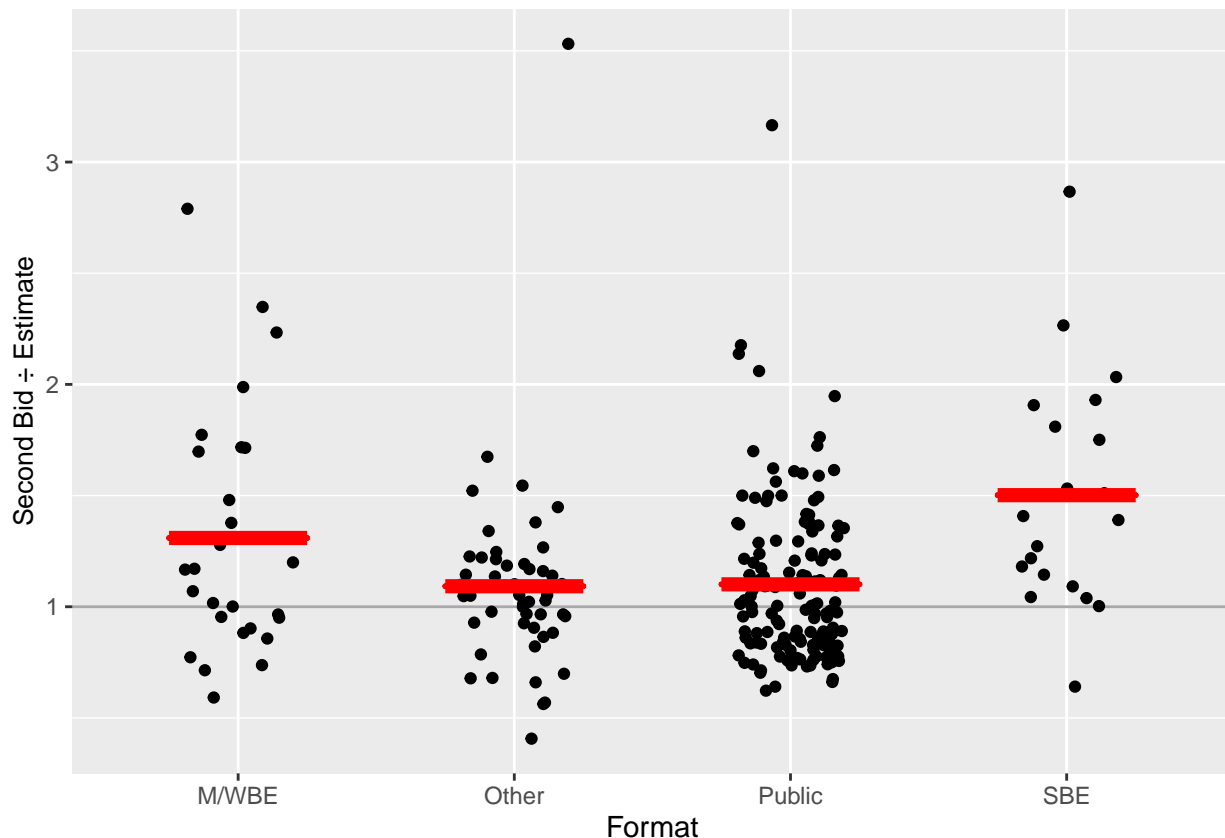
A project's bidding process can be constrained or open, with potential ramifications on the ability to estimate costs:

```
summary(lm(bids$accuracy ~ as.factor(bids$Format)))
```

```
##
## Call:
## lm(formula = bids$accuracy ~ as.factor(bids$Format))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.86097 -0.27198 -0.08562  0.16241  2.43943
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.30933    0.08035  16.295  <2e-16 ***
## as.factor(bids$Format)Other -0.21709    0.09937  -2.185   0.0299 *
## as.factor(bids$Format)Public -0.20862    0.08771  -2.379   0.0182 *
## as.factor(bids$Format)SBE    0.19246    0.12318   1.563   0.1195
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 0.4175 on 235 degrees of freedom
## Multiple R-squared:  0.08211,    Adjusted R-squared:  0.07039
## F-statistic: 7.007 on 3 and 235 DF,  p-value: 0.000156

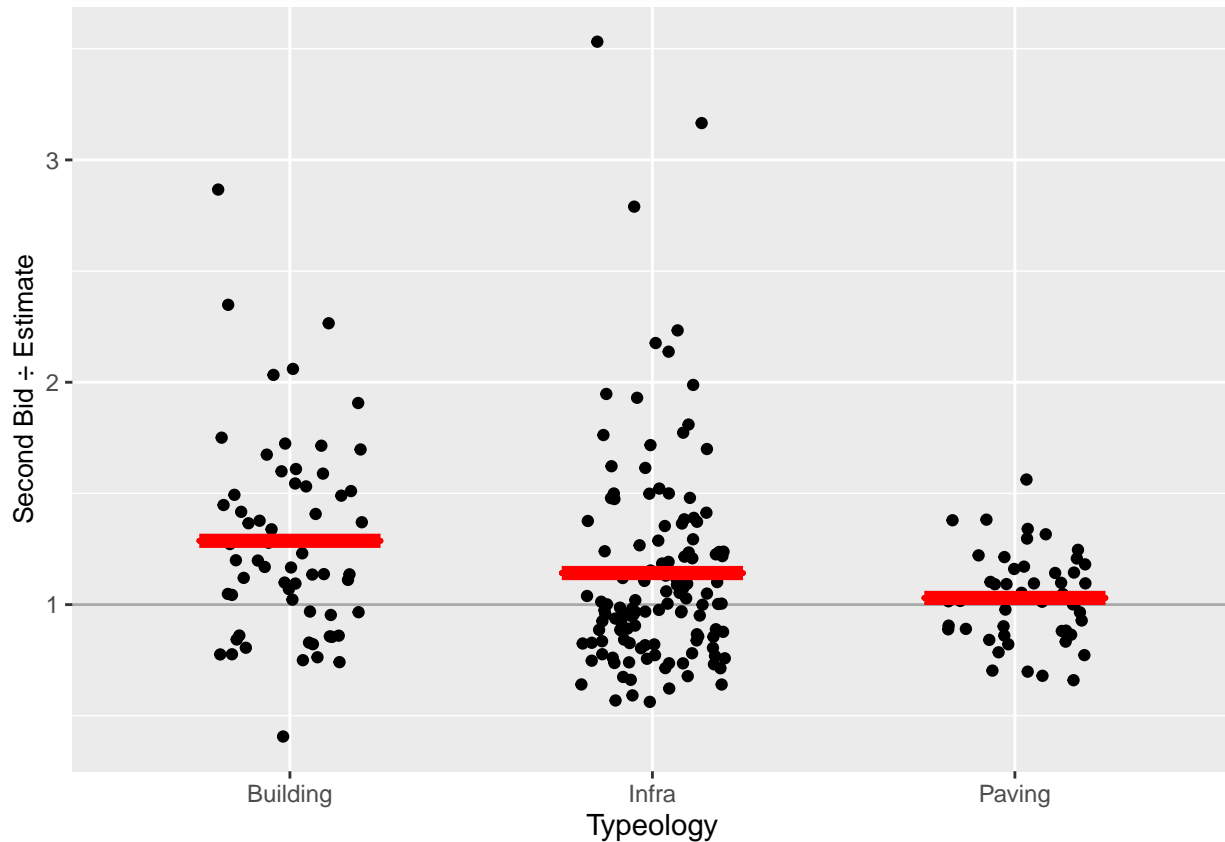
ggplot(bids, aes(x = Format, y = accuracy)) + #, shape = Loc)) +
  geom_hline(colour="dark gray", yintercept=1) +
  geom_jitter(width=0.2) +
  geom_crossbar(data=format.summary, aes(ymin = accuracy, ymax = accuracy),
               size=1,col="red", width = .5) +
  ylab("Second Bid ÷ Estimate") +
  theme(axis.title.y=element_text(size=10))
```



The signal is stronger regarding the type of project, which has an identifiable (if yet uncontrolled) relationship with estimating accuracy. What's the raw (uncontrolled) relationship between estimation accuracy and the type of project?

```
##
## Call:
## lm(formula = bids$accuracy ~ as.factor(bids$Typeology))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.88015 -0.26202 -0.08883  0.15642  2.38969
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.28697    0.05451   23.609 < 2e-16 ***
## as.factor(bids$Typeology)Infra -0.14499    0.06616   -2.192  0.02938 *
```

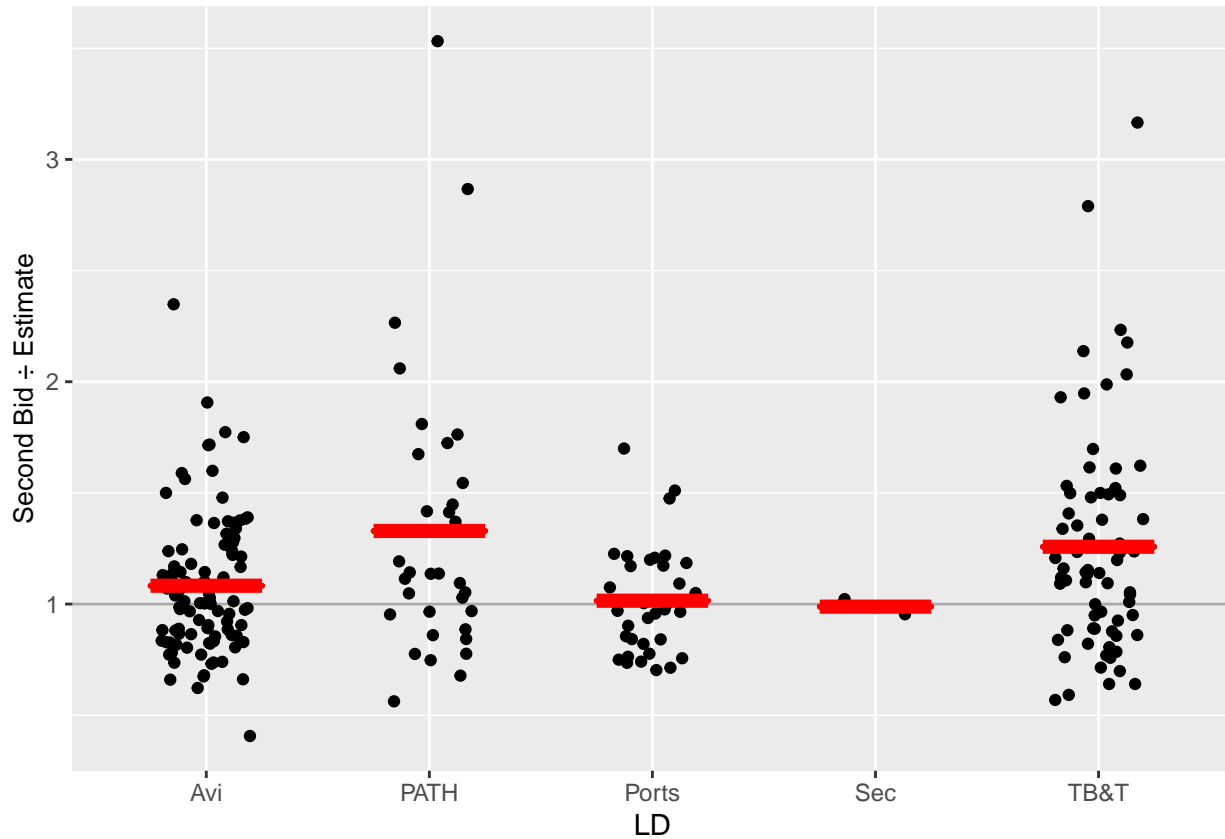
```
## as.factor(bids$Typeology)Paving -0.25693    0.08168  -3.146  0.00187 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4258 on 236 degrees of freedom
## Multiple R-squared:  0.04143,    Adjusted R-squared:  0.03331
## F-statistic:    5.1 on 2 and 236 DF,  p-value: 0.006784
```



By line department:

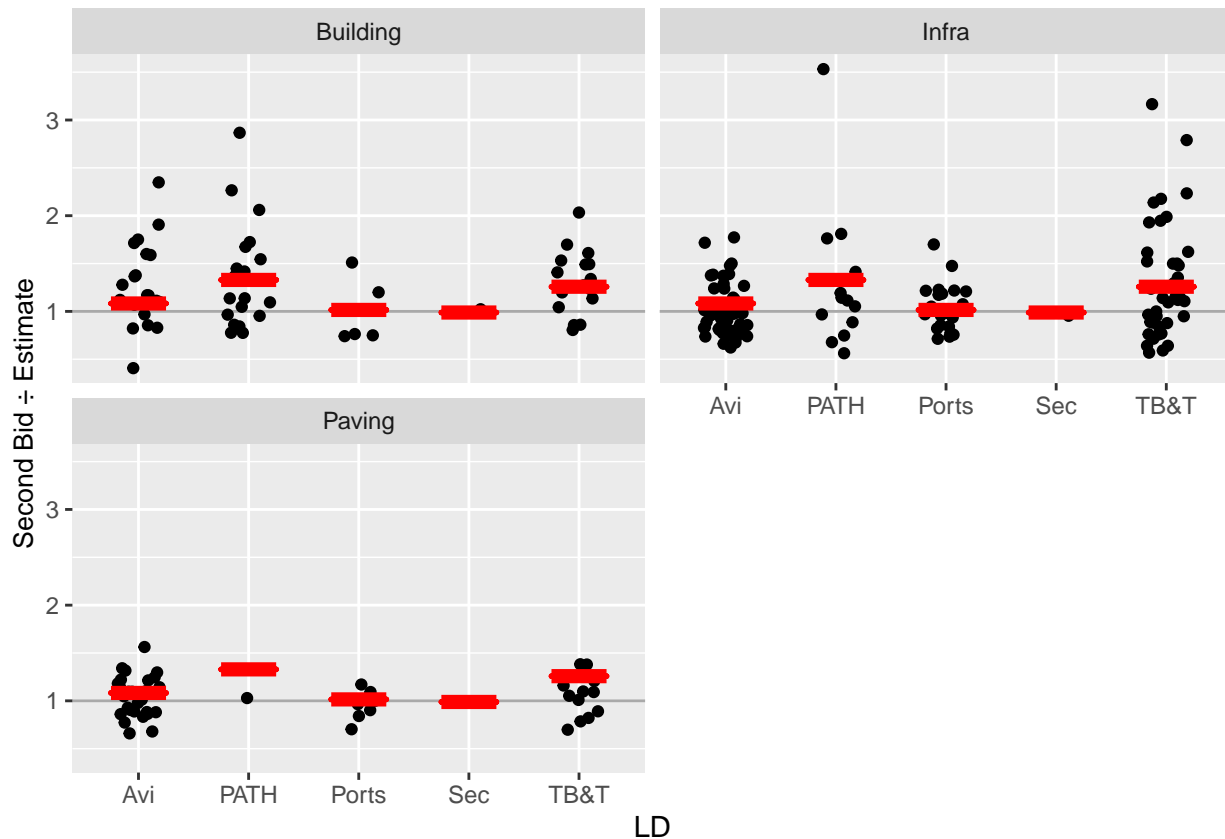
```
##
## Call:
## lm(formula = bids$accuracy ~ as.factor(bids$LD))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.76624 -0.25864 -0.06976  0.18420  2.20275
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.08270    0.04260   25.417 < 2e-16 ***
## as.factor(bids$LD)PATH  0.24622    0.08487    2.901  0.00407 **
## as.factor(bids$LD)Ports -0.06742    0.08137   -0.829  0.40819
## as.factor(bids$LD)Sec   -0.09435    0.30121   -0.313  0.75439
## as.factor(bids$LD)TB&T  0.17514    0.06627    2.643  0.00878 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 0.4217 on 234 degrees of freedom
## Multiple R-squared:  0.06762,    Adjusted R-squared:  0.05168
## F-statistic: 4.243 on 4 and 234 DF,  p-value: 0.002468
```



Infrastructure projects for TB&T seem difficult to target, but the results aren't systematically high or low, just spread across a wide range:

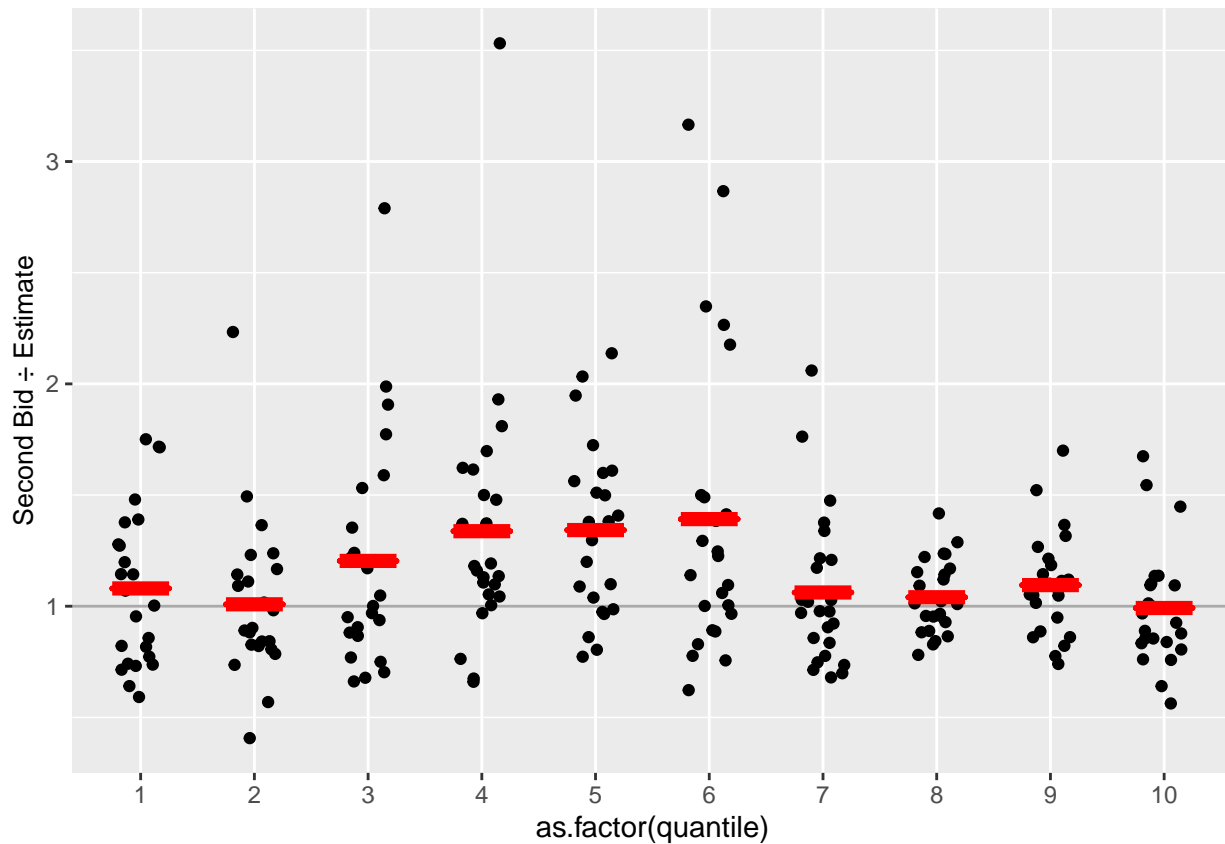
```
ggplot(bids, aes(x = LD, y = accuracy)) + #, shape = Typeology)) +
  geom_hline(colour="dark gray", yintercept=1) +
  geom_jitter(width=0.2) +
  geom_crossbar(data=ld.summary, aes(ymin = accuracy, ymax = accuracy),
               size=1,col="red", width = .5) +
  facet_wrap( ~ Typeology, ncol=2) +
  ylab("Second Bid ÷ Estimate") +
  theme(axis.title.y=element_text(size=10))
```

We see signals that projects of medium size (in dollar terms) may be less evasive than much larger or smaller projects. What's the relationship between low bids and accuracy, when we start considering the size of low bids?

```
##
## Call:
## lm(formula = bids$accuracy ~ as.factor(bids$quantile))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.76921 -0.25251 -0.07487  0.15180  2.19418
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.07999    0.08500   12.705  <2e-16 ***
## as.factor(bids$quantile)2 -0.07132    0.12021   -0.593    0.5536
## as.factor(bids$quantile)3  0.12379    0.12021    1.030    0.3042
## as.factor(bids$quantile)4  0.25750    0.12021    2.142    0.0332 *
## as.factor(bids$quantile)5  0.26254    0.12021    2.184    0.0300 *
## as.factor(bids$quantile)6  0.31188    0.12021    2.594    0.0101 *
## as.factor(bids$quantile)7 -0.01821    0.12021   -0.152    0.8797
## as.factor(bids$quantile)8 -0.03953    0.12021   -0.329    0.7426
## as.factor(bids$quantile)9  0.01498    0.12021    0.125    0.9009
## as.factor(bids$quantile)10 -0.08804    0.12151   -0.725    0.4695
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4164 on 229 degrees of freedom
```

```
## Multiple R-squared:  0.1102, Adjusted R-squared:  0.0752
## F-statistic:  3.15 on 9 and 229 DF,  p-value: 0.001329
```



Visualizing and testing the data iteratively like this has offered some initial insight into what might be accounting for variation in accuracy.

Modeling and prediction

Try and use exogenous covariates to predict an alternative engineering estimate, without using the low bid information, that might be closer to the low bid. Call it “expected low bid” or something so we can remember what we’re trying to get.

A. Base model (manual selection)

Interpretation: specification was manual and intuitive.

Note: ensure the “accuracy” variable we calculated earlier is dropped before modeling or I’ll be introducing some dual (reverse) causality, which will confuse the models.

Note 2: when a number appears in the output without context, it is likely an information criterion (and AIC), which may or may not provide value post-modeling.

First remove accuracy.

```
train$accuracy=NULL
test$accuracy=NULL
```

Run model.

```
##
## Call:
## lm(formula = Second.Bid ~ Engr.Est + Employment.in.construction +
##     Format2 + Typeology + permits_1 + cci + as.factor(quantile),
##     data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -65506055  -1141743    81828   1117149   54487762
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.985e+06  2.045e+07   0.342   0.7330
## Engr.Est          8.813e-01  1.627e-02  54.155  <2e-16 ***
## Employment.in.construction -1.017e+05  1.821e+05  -0.558   0.5771
## Format2Public      -1.852e+06  1.126e+06  -1.645   0.1015
## TypeologyInfra     -1.820e+06  1.274e+06  -1.428   0.1547
## TypeologyPaving    -1.637e+06  1.575e+06  -1.040   0.2997
## permits_1          5.519e+01  1.014e+02   0.544   0.5867
## cci                2.768e+03  6.714e+03   0.412   0.6806
## as.factor(quantile)2  1.051e+06  2.339e+06   0.449   0.6537
## as.factor(quantile)3  1.064e+06  2.318e+06   0.459   0.6465
## as.factor(quantile)4  1.622e+06  2.326e+06   0.697   0.4864
## as.factor(quantile)5  1.738e+06  2.371e+06   0.733   0.4645
## as.factor(quantile)6  2.271e+06  2.340e+06   0.971   0.3329
## as.factor(quantile)7  1.971e+06  2.381e+06   0.828   0.4088
## as.factor(quantile)8  2.439e+06  2.390e+06   1.021   0.3084
## as.factor(quantile)9  5.295e+06  2.371e+06   2.233   0.0265 *
## as.factor(quantile)10 2.802e+06  2.963e+06   0.946   0.3454
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7817000 on 222 degrees of freedom
## Multiple R-squared:  0.9658, Adjusted R-squared:  0.9633
## F-statistic: 391.9 on 16 and 222 DF,  p-value: < 2.2e-16
```

The equation above throws extra information at the engineering estimate and tries to predict the low qualifying bid. If the result is noticeably closer to the low qualifying bid than the original estimate, you can use the delta as a post-estimation fudge factor to adjust your final estimate.

The numbers calculated above include a few metrics to use in comparing the three estimates against the objective data point at the low bid, which is what we're trying to predict. Ghose three estimates I'm talking about are: 1. The original, raw engineering cost estimate, 2. The first alternative, where we built a model by hand to try and use a few more data points to enhance the original estimate, and; 3. The second alternative, a kitchen sink model that throws even more data points at the question. This followed an effort to use a penalized regression to identify the best covariates, but that penalization algorithm actually suggested there isn't much we can do to enhance the original estimate. (Note: this will prove prescient.)

I'll build a table near the very end of this script that summarize the metrics I'm using to understand how well these modeling efforts work. The metrics will be: A. A basic t-test to understand whether there's even a statistically significant difference between the estimate I'm getting and the enhanced estimate I'm modeling with it, B. A correlation between the two numbers, to try and understand the magnitude of that difference (if we can trust it really exists), C. Two measures of the predictive modeling power of the models, an adjusted R-squared and the mean squared error (MSE). Both are common metrics of power. The first can be viewed discretely for each model but the second only provides a relative measure between models.

What is the summary of the predicted values? How does it compare to the summary of low bids?

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	-1569853	1720892	3478279	14925516	9814439	425157238

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	142000	1556557	3185420	14925516	9143818	479645000

Looking at the two summaries above, what's more accurate, the original engineering estimate or the first enhanced prediction? Neither, really. In fact the estimate and prediction aren't even (statistically) significantly different. Maybe something more robust can come with a little creativity.

(Note: the model should control to prevent negative values. To be done next time.)

It might have been worth trying with the log of prices, only because the statistical fit becomes multiplicative instead of linear, but that produced similar results.

B. Machine learning

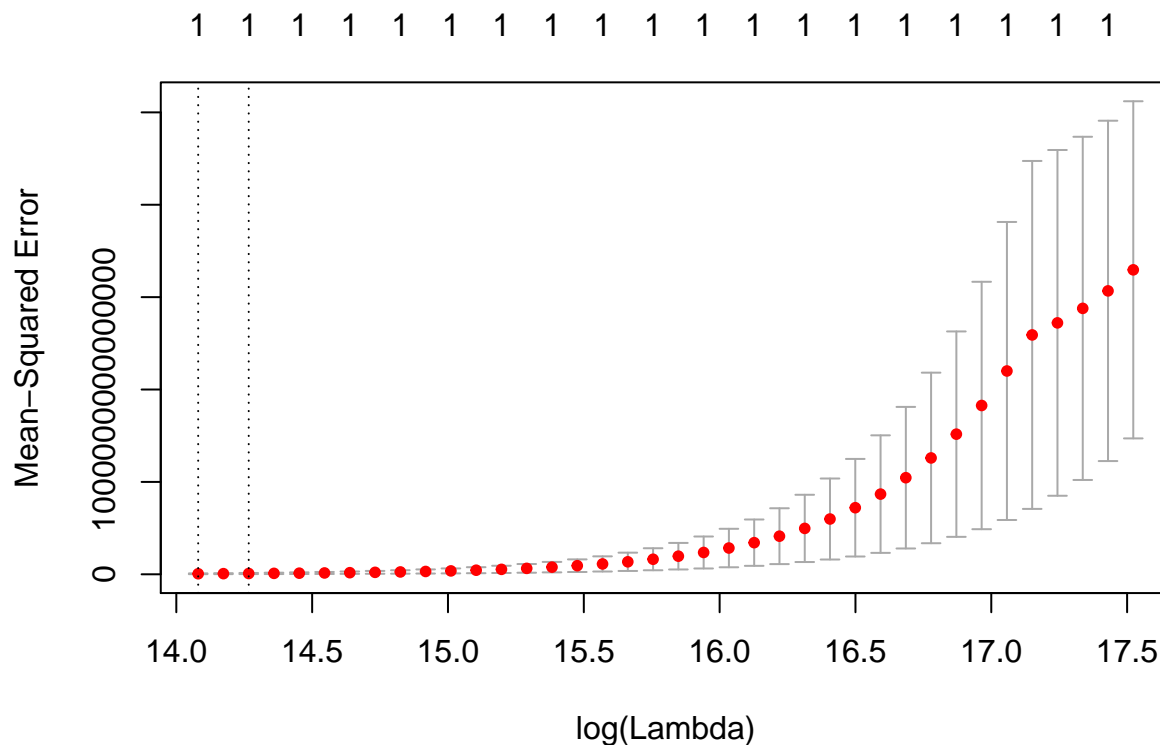
Strip the bids data of unusable stuff, then set controls and run. The package I'm using, glmnet, requires a little extra preparation.

The output below represents the algorithm's effort to look for variables that might be dependable in adding predictive power to the original engineering estimate.

Running a regularized regression below adopts base assumptions in the modeling software package (10 folds, standardized coefficients, gaussian distribution, MSE evaluation metric, et cetera).

How many variables survive the penalty process as the penalty grows?

```
plot(lasso)
```



At a lower penalty (higher tuning parameter) there are a dozen non-zero predictors remaining, once the parameter cranks up a little all but a couple fall away. A little fine-tuning identifies the two parameters that survive as the model gets further from ordinary:

```
lasso2 = cv.glmnet(x=XP, y=YP, alpha=.25)
coef(lasso2, s=lasso2$lambda.1se)
```

```
## 46 x 1 sparse Matrix of class "dgCMatrix"
##                                     1
## (Intercept)                    -2817346.4089514
## Engr.Est                        0.1919425
## Loc                             .
## LD                             .
## Typeology                       .
## Second.Bid                     0.7742366
## Consumer.price.index           .
## Employment.in.communications    .
## Employment.in.construction      .
## Employment.in.education.and.health .
## Employment.in.financial.and.business.services .
## Employment.in.financial.services .
## Employment.in.government        .
## Employment.in.other.services    .
## Employment.in.production.industries .
## Employment.in.professional.services .
## Employment.in.real.estate       .
## Employment.in.retail            .
## Employment.in.transport.services .
## Employment.in.wholesale         .
## Output.in.communications         .
## Output.in.construction           .
## Output.in.financial.services     .
## Output.in.government             97.8327334
## Output.in.retail                 .
## Output.in.education.and.health   .
## Output.in.financial.and.business.services .
## Output.in.other.services         .
## Output.in.production.industries .
## Output.in.professional.services .
## Output.in.real.estate           .
## Output.in.transport.services    .
## Output.in.wholesale              .
## Personal.disposable.income..nominal .
## Personal.disposable.income..real .
## Personal.income..nominal        .
## Retail.sales..nominal           .
## Retail.sales..real              .
## Total.employment                .
## Total.office.based.employment   .
## Total.output                    .
## permits                         .
## permits_1                       .
## cci                             .
## Format2                         -57244.4206265
## quantile                        51947.3854675
```

This suggests using anything beyond the engineer's estimate itself to better predict the lowest qualifying good adds more uncertainty (in the form of noise that's tough to explain) than it adds value. (The "penalty"

associated with adding variables is greater than the extra predictive power they bring.) The most obvious exception is project size, which isn't too surprising given the exploratory work done earlier.

```

bids = fastDummies::dummy_cols(bids, select_columns = "quantile")
summary(lm(bids$Second.Bid ~ bids$Engr.Est + bids$quantile_2 + bids$quantile_3 + bids$quantile_4 + bids$quantile_5 + bids$quantile_6 + bids$quantile_7 + bids$quantile_8 + bids$quantile_9 + bids$quantile_10))

##
## Call:
## lm(formula = bids$Second.Bid ~ bids$Engr.Est + bids$quantile_2 +
##     bids$quantile_3 + bids$quantile_4 + bids$quantile_5 + bids$quantile_6 +
##     bids$quantile_7 + bids$quantile_8 + bids$quantile_9 + bids$quantile_10)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -64927474  -490838    31601   573534  56729942
##
## Coefficients:
##              Estimate      Std. Error t value      Pr(>|t|)
## (Intercept)   76048.34661  1598483.07814    0.048    0.9621
## bids$Engr.Est    0.87692     0.01604  54.685 <0.0000000000000002
## bids$quantile_2 -34059.25935  2260576.56828  -0.015    0.9880
## bids$quantile_3  191658.26895  2260597.38380   0.085    0.9325
## bids$quantile_4  467164.92654  2260635.79910   0.207    0.8365
## bids$quantile_5  724044.57116  2260706.51997   0.320    0.7491
## bids$quantile_6  921683.48531  2261008.27650   0.408    0.6839
## bids$quantile_7  491236.92917  2262088.51151   0.217    0.8283
## bids$quantile_8  1195745.54423  2264742.28422   0.528    0.5980
## bids$quantile_9  4099789.39909  2289440.33695   1.791    0.0747
## bids$quantile_10 2003942.85274  2907737.97455   0.689    0.4914
##
## (Intercept)
## bids$Engr.Est ***
## bids$quantile_2
## bids$quantile_3
## bids$quantile_4
## bids$quantile_5
## bids$quantile_6
## bids$quantile_7
## bids$quantile_8
## bids$quantile_9 .
## bids$quantile_10
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7831000 on 228 degrees of freedom
## Multiple R-squared:  0.9648, Adjusted R-squared:  0.9632
## F-statistic: 624.1 on 10 and 228 DF,  p-value: < 0.00000000000000022

```