

# Case Study: Modeling Liquid Mechanics

Conner Byrd, Eric Han, Ki Hyun, Sara Shao, Alex Shen, Mona Su, Dani Trejo, Steven Yuan

10/12/2021

## Introduction

Our key research objectives include understanding and predicting how turbulence affects the dynamics of water droplets and ice crystals (how they collide and mix) in clouds. With our machine learning model, we are trying to infer the volume distribution of clusters within clouds.

To do this, we began by doing some basic Exploratory Data Analysis on the three predictor variables: Reynolds number (**Re**), gravitational acceleration (**Fr**), and particle characteristic (**St**). (Our graphs are included in Appendix: Section 1)

## Methodology

```
head(train)
```

```
##      St  Re    Fr R_moment_1 R_moment_2 R_moment_3 R_moment_4
## 1 0.10 224 0.052 0.00215700 0.1303500 14.37400 1586.5000
## 2 3.00 224 0.052 0.00379030 0.4704200 69.94000 10404.0000
## 3 0.70 224 Inf 0.00290540 0.0434990 0.82200 15.5510
## 4 0.05 90 Inf 0.06352800 0.0906530 0.46746 3.2696
## 5 0.70 398 Inf 0.00036945 0.0062242 0.12649 2.5714
## 6 2.00 90 0.300 0.14780000 2.0068000 36.24900 671.6700
```

```
train_data <- train %>%
  mutate(Fr = as.ordered(Fr)) %>%
  mutate(Re = as.ordered(Re))
```

We decided at the beginning to treat Fr and Re as ordered factors. It makes sense to treat Fr as a categorical variable because we are given only three unique values, one of which is infinity, and in practice the three values are representative of different types of clouds. We decided to treat Re as a factored variable as well because we are also only given three unique values and because the differences between the three values are so large that it would be unwise to extrapolate our model to the ranges in between the values we are given. In terms of prediction and inference, we believe our models can still be generalized to Fr and Re values similar to the ones we are working with.

```
##
## Call:
## lm(formula = R_moment_1 ~ St + Re + Fr, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.038834 -0.008614  0.001702  0.009854  0.039423
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)  0.027329   0.002494  10.959 < 2e-16 ***
## St          0.012213   0.002078   5.877 8.42e-08 ***
## Re.L        -0.078880   0.003276 -24.081 < 2e-16 ***
## Re.Q         0.042715   0.002757  15.491 < 2e-16 ***
## Fr.L        -0.007219   0.002678  -2.696 0.00849 **
## Fr.Q         0.002056   0.003181   0.646 0.51987
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.01529 on 83 degrees of freedom
## Multiple R-squared:  0.9293, Adjusted R-squared:  0.9251
## F-statistic: 218.2 on 5 and 83 DF,  p-value: < 2.2e-16

##
## Call:
## lm(formula = R_moment_4 ~ St + Re + Fr, data = train_data)
##
## Residuals:
##          Min           1Q       Median           3Q          Max
## -2.495e+10 -1.019e+10 -4.413e+09   7.899e+09  4.555e+10
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.461e+09  2.328e+09   0.628 0.532039
## St           3.220e+09  1.940e+09   1.660 0.100777
## Re.L        -1.463e+10  3.058e+09  -4.786 7.32e-06 ***
## Re.Q         5.531e+09  2.574e+09   2.149 0.034582 *
## Fr.L        -1.014e+10  2.500e+09  -4.057 0.000112 ***
## Fr.Q         8.693e+09  2.970e+09   2.927 0.004411 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.427e+10 on 83 degrees of freedom
## Multiple R-squared:  0.4253, Adjusted R-squared:  0.3906
## F-statistic: 12.28 on 5 and 83 DF,  p-value: 6.371e-09
```

We decided to initially fit the most basic linear model with all three predictors to see what it would look like. The model for the first moment had a fairly high  $R^2$  value (0.929), but the model for the fourth moment had a much lower  $R^2$  value (0.425). Additionally, in our diagnostic plots (see Appendix: Section 2), we saw a pattern in the residuals vs fitted values plots where the models would consistently under predict in some areas and over predict in others, indicating non-linearity. In addition, looking at the Normal Q-Q plots, the normality assumption also seemed to be violated for higher moments. This is consistent with the fact that our histogram of St in our EDA was not normally distributed.

This information lead us to try using a GAM to model the relationship between the predictors and the 4 moments due to the increased flexibility GAMs provide. However, we knew that using GAMs made interpretability an issue, because interpreting a complex smooth function of a continuous predictor is very hard.

As a result, we decided to use variable transformations and interaction effects to make linear models with suitable model diagnostics for all 4 moments for the purpose of inference. We also tried doing forward selection with AIC to see if we could select a simpler model in case were overfitting, but our resulting models were the same as our input models (see Appendix: Section 2). We would also compare our final linear models with 4 GAM models (one for each moment) using 10-fold CV in order to find the best models for prediction.

## Results

### Final Linear Model

```
##
## Call:
## lm(formula = log(R_moment_1) ~ log(St) + Re + Fr + St * Fr +
##      Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.211809 -0.042926 -0.006391  0.038831  0.171243
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -5.349488   0.027649 -193.480 < 2e-16 ***
## log(St)      0.145668   0.014562  10.003 1.89e-15 ***
## Re.L        -4.028476   0.027949 -144.139 < 2e-16 ***
## Re.Q         0.644476   0.022457  28.698 < 2e-16 ***
## Fr.L        -0.102135   0.019793  -5.160 1.95e-06 ***
## Fr.Q         0.109493   0.026874   4.074 0.000113 ***
## St           0.095896   0.021136   4.537 2.13e-05 ***
## Fr.L:St      0.095376   0.017271   5.522 4.60e-07 ***
## Fr.Q:St     -0.064244   0.022042  -2.915 0.004692 **
## Re.L:Fr.L    0.242947   0.024770   9.808 4.39e-15 ***
## Re.Q:Fr.L   -0.076314   0.022588  -3.379 0.001159 **
## Re.L:Fr.Q   -0.077781   0.046306  -1.680 0.097169 .
## Re.Q:Fr.Q           NA          NA      NA      NA
## Re.L:St     -0.008897   0.021672  -0.411 0.682581
## Re.Q:St     -0.025325   0.018376  -1.378 0.172252
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07645 on 75 degrees of freedom
## Multiple R-squared:  0.999, Adjusted R-squared:  0.9988
## F-statistic: 5797 on 13 and 75 DF, p-value: < 2.2e-16
##
## Call:
## lm(formula = log(R_moment_2) ~ log(St) + Re + Fr + St * Fr +
##      Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.6442 -0.2697 -0.0561  0.3429  1.8016
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.300240   0.307283  -0.977  0.33167
## log(St)      1.500864   0.161841   9.274 4.50e-14 ***
## Re.L        -4.269814   0.310614 -13.746 < 2e-16 ***
## Re.Q         1.098368   0.249584   4.401 3.52e-05 ***
## Fr.L        -1.927717   0.219976  -8.763 4.20e-13 ***
## Fr.Q         0.946857   0.298671   3.170 0.00221 **
## St          -0.997463   0.234905  -4.246 6.16e-05 ***
```

```

## Fr.L:St      -0.007867    0.191944   -0.041   0.96742
## Fr.Q:St      -0.043393    0.244968   -0.177   0.85988
## Re.L:Fr.L     3.399373    0.275290   12.348   < 2e-16 ***
## Re.Q:Fr.L    -0.646249    0.251039   -2.574   0.01202 *
## Re.L:Fr.Q    -2.588037    0.514630   -5.029   3.27e-06 ***
## Re.Q:Fr.Q           NA           NA           NA           NA
## Re.L:St      -0.470624    0.240856   -1.954   0.05443 .
## Re.Q:St      -0.128587    0.204226   -0.630   0.53085
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8497 on 75 degrees of freedom
## Multiple R-squared:  0.9553, Adjusted R-squared:  0.9476
## F-statistic: 123.4 on 13 and 75 DF,  p-value: < 2.2e-16

##
## Call:
## lm(formula = log(R_moment_3) ~ log(St) + Re + Fr + St * Fr +
##     Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.7949 -0.4431 -0.1224  0.5575  2.9257
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.54701     0.51756   8.785 3.81e-13 ***
## log(St)       2.36189     0.27259   8.665 6.47e-13 ***
## Re.L         -4.89792     0.52317  -9.362 3.06e-14 ***
## Re.Q          1.52364     0.42038   3.624 0.000525 ***
## Fr.L         -3.75114     0.37051 -10.124 1.12e-15 ***
## Fr.Q          1.83513     0.50306   3.648 0.000486 ***
## St           -1.72871     0.39565  -4.369 3.95e-05 ***
## Fr.L:St      -0.09194     0.32329  -0.284 0.776889
## Fr.Q:St      -0.03580     0.41260  -0.087 0.931098
## Re.L:Fr.L     6.45803     0.46368  13.928 < 2e-16 ***
## Re.Q:Fr.L    -1.15825     0.42283  -2.739 0.007689 **
## Re.L:Fr.Q    -4.96308     0.86680  -5.726 2.01e-07 ***
## Re.Q:Fr.Q           NA           NA           NA           NA
## Re.L:St      -0.79750     0.40568  -1.966 0.053018 .
## Re.Q:St      -0.17827     0.34398  -0.518 0.605816
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.431 on 75 degrees of freedom
## Multiple R-squared:  0.9459, Adjusted R-squared:  0.9365
## F-statistic: 100.8 on 13 and 75 DF,  p-value: < 2.2e-16

##
## Call:
## lm(formula = log(R_moment_4) ~ log(St) + Re + Fr + St * Fr +
##     Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max

```

```
## -9.6675 -0.6183 -0.1392  0.7410  3.8875
##
## Coefficients: (1 not defined because of singularities)
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  9.28287    0.70476  13.172 < 2e-16 ***
## log(St)      3.10948    0.37119   8.377 2.28e-12 ***
## Re.L        -5.63270    0.71240  -7.907 1.80e-11 ***
## Re.Q         1.94835    0.57242   3.404 0.001070 **
## Fr.L        -5.55550    0.50452 -11.012 < 2e-16 ***
## Fr.Q         2.71540    0.68500   3.964 0.000167 ***
## St          -2.36662    0.53876  -4.393 3.62e-05 ***
## Fr.L:St     -0.17635    0.44023  -0.401 0.689864
## Fr.Q:St     -0.01987    0.56184  -0.035 0.971880
## Re.L:Fr.L    9.47956    0.63138  15.014 < 2e-16 ***
## Re.Q:Fr.L   -1.66019    0.57576  -2.883 0.005130 **
## Re.L:Fr.Q   -7.28081    1.18031  -6.169 3.21e-08 ***
## Re.Q:Fr.Q      NA         NA         NA      NA
## Re.L:St     -1.08269    0.55241  -1.960 0.053716 .
## Re.Q:St     -0.22245    0.46840  -0.475 0.636226
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.949 on 75 degrees of freedom
## Multiple R-squared:  0.9457, Adjusted R-squared:  0.9363
## F-statistic: 100.6 on 13 and 75 DF,  p-value: < 2.2e-16
```

A one percent increase in Stokes number is associated with 0.146% increase in R moment 1, holding all other predictors constant. When the Reynolds number is 224, the R moment 1 is expected to decrease by 403% from when the Reynolds number is 90, holding all other predictors constant. When the Reynolds number is 398, the R moment 1 is expected to increase by 64% compared to when the Reynolds number is 90. When the Reynolds number is 224 and the Froud number is 0.3, the R moment 1 is expected to be an additional 24% lower compared to when either of those conditions are not met.

### Predicted Test Error For R Moments 1-4 With Linear Model

Predicted Mean-Squared Error For R Moment 1

```
## [1] 2.864351e-05
```

Predicted Mean-Squared Error For R Moment 2

```
## [1] 4621.783
```

Predicted Mean-Squared Error For R Moment 3

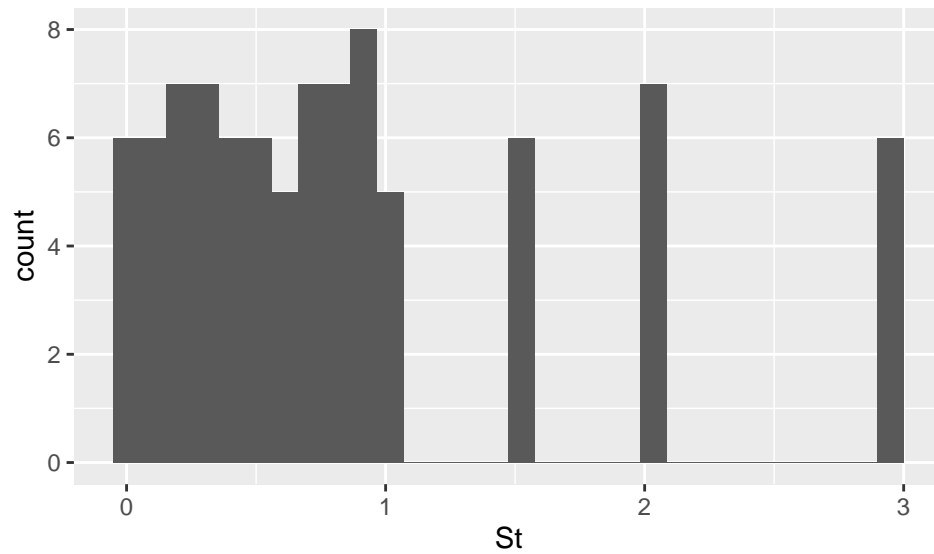
```
## [1] 578718160228
```

Predicted Mean-Squared Error For R Moment 4

```
## [1] 5.488601e+19
```

## Appendix

Figure 1.1: Distribution of St



We will try using a log transform on the St variable since the distribution for the St variable is not normally distributed.

Figure 1.2: R Moments and Re Colored by Fr

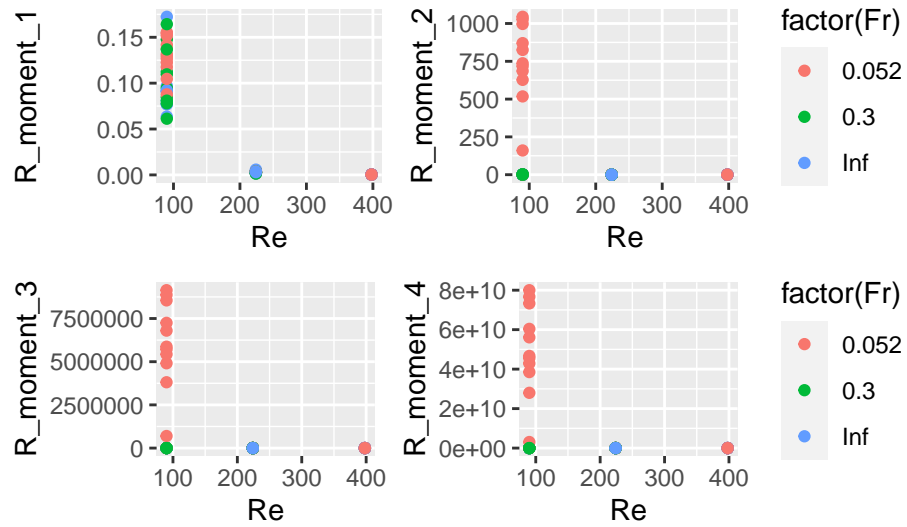
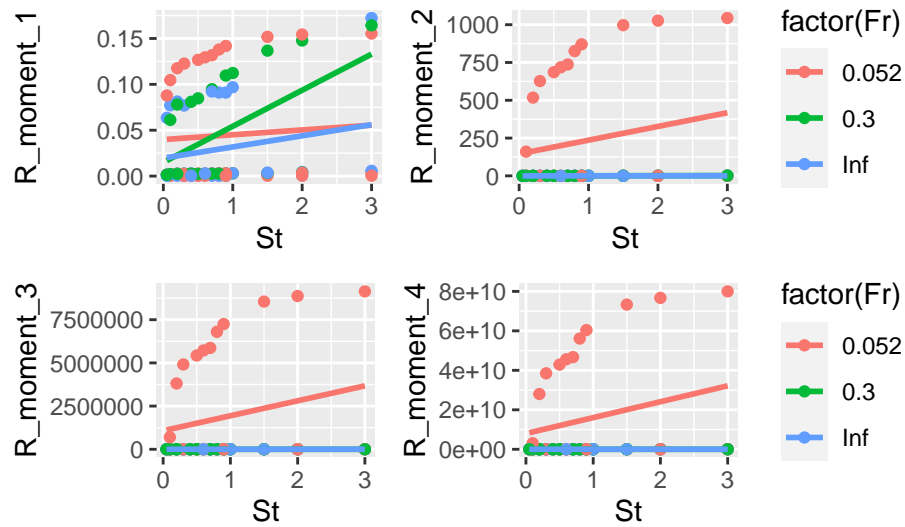


Figure 1.3: R Moments and St Colored by Fr



The graphs above show some evidence of interactions, so we will explore interaction terms in our model.

Figure 2.1-2: Diagnostic Plots For Initial Linear Model of R Moment 1

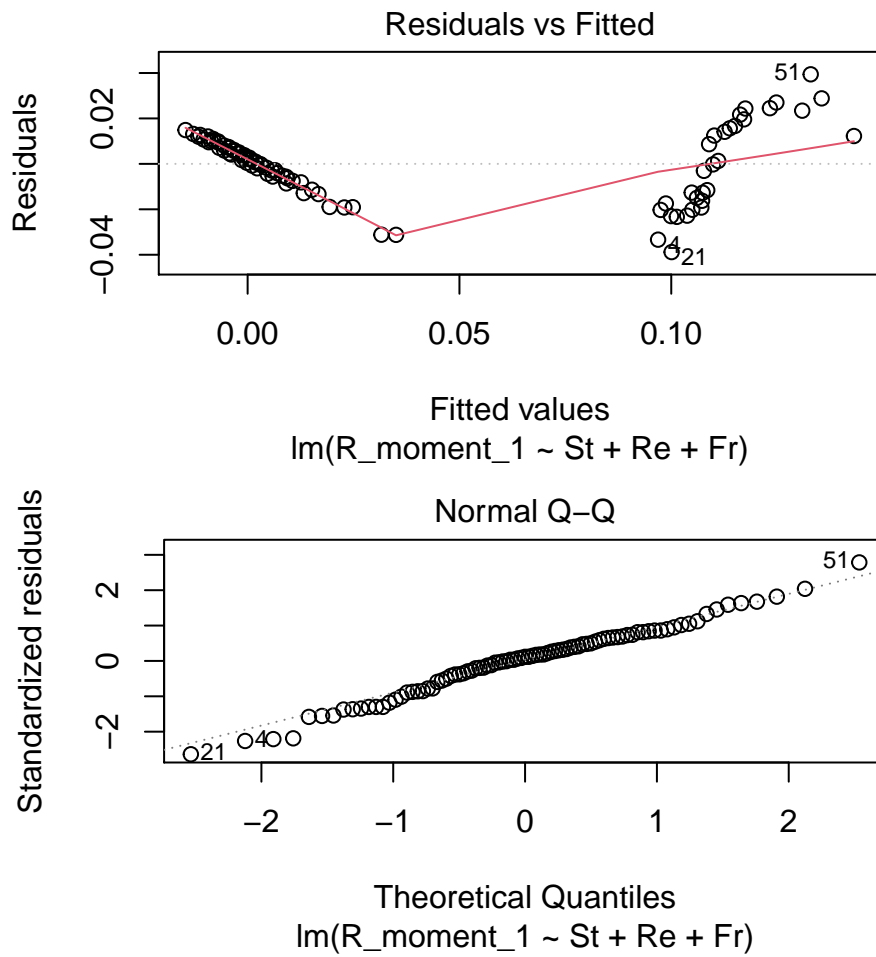
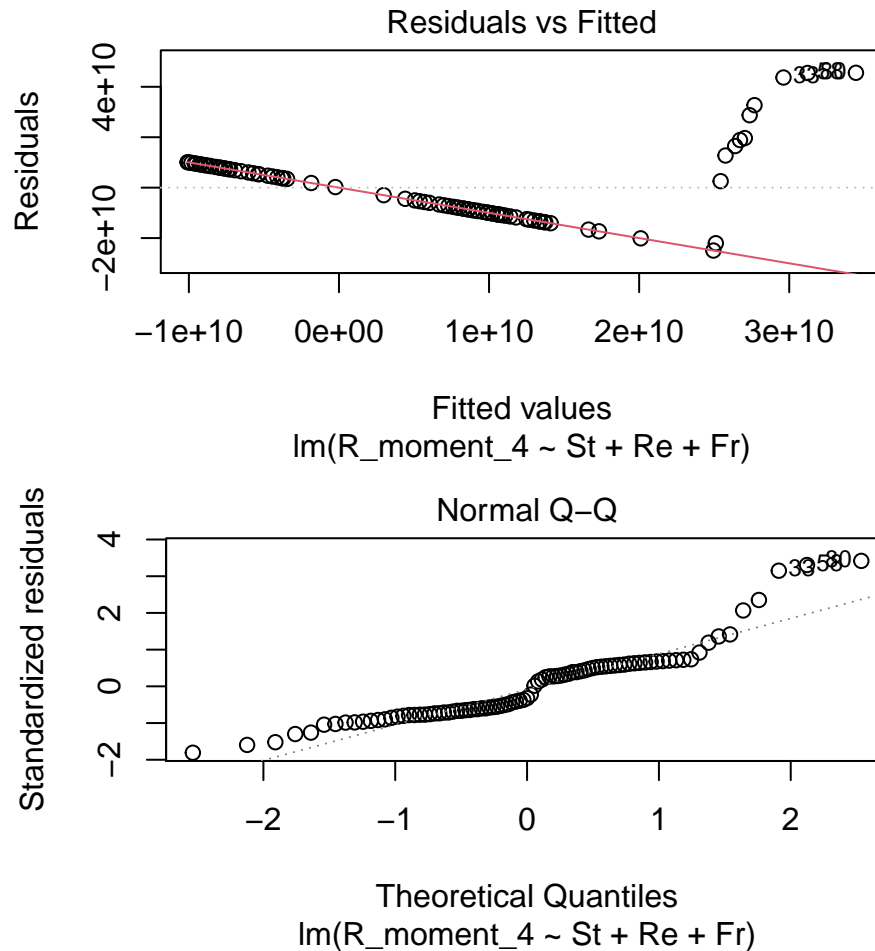


Figure 2.3-4: Diagnostic Plots For Inital Linear Model of R Moment 4



Because the linearity condition is not fulfilled in the above Residuals vs. Fitted plots, we will consider performing a log transformation on our response variables (R moments 1-4).

Figure 2.5: Forward Selection With AIC On Linear Models

```
## Start: AIC=-444.89
## log(R_moment_1) ~ log(St) + Re + Fr + St * Fr + Fr * Re + St *
## Re
##
## Call:
## lm(formula = log(R_moment_1) ~ log(St) + Re + Fr + St * Fr +
## Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.211809 -0.042926 -0.006391  0.038831  0.171243
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -5.349488   0.027649 -193.480 < 2e-16 ***
## log(St)      0.145668   0.014562  10.003 1.89e-15 ***
## Re.L        -4.028476   0.027949 -144.139 < 2e-16 ***
```



```

## Re.Q      0.644476  0.022457  28.698 < 2e-16 ***
## Fr.L      -0.102135  0.019793  -5.160 1.95e-06 ***
## Fr.Q      0.109493  0.026874   4.074 0.000113 ***
## St        0.095896  0.021136   4.537 2.13e-05 ***
## Fr.L:St   0.095376  0.017271   5.522 4.60e-07 ***
## Fr.Q:St   -0.064244  0.022042  -2.915 0.004692 **
## Re.L:Fr.L  0.242947  0.024770   9.808 4.39e-15 ***
## Re.Q:Fr.L -0.076314  0.022588  -3.379 0.001159 **
## Re.L:Fr.Q -0.077781  0.046306  -1.680 0.097169 .
## Re.Q:Fr.Q      NA      NA      NA      NA
## Re.L:St     -0.008897  0.021672  -0.411 0.682581
## Re.Q:St     -0.025325  0.018376  -1.378 0.172252
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07645 on 75 degrees of freedom
## Multiple R-squared:  0.999, Adjusted R-squared:  0.9988
## F-statistic: 5797 on 13 and 75 DF, p-value: < 2.2e-16

## Start:  AIC=-16.23
## log(R_moment_2) ~ log(St) + Re + Fr + St * Fr + Fr * Re + St *
##      Re

##
## Call:
## lm(formula = log(R_moment_2) ~ log(St) + Re + Fr + St * Fr +
##      Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.6442 -0.2697 -0.0561  0.3429  1.8016
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.300240   0.307283  -0.977  0.33167
## log(St)      1.500864   0.161841   9.274 4.50e-14 ***
## Re.L        -4.269814   0.310614 -13.746 < 2e-16 ***
## Re.Q         1.098368   0.249584   4.401 3.52e-05 ***
## Fr.L        -1.927717   0.219976  -8.763 4.20e-13 ***
## Fr.Q         0.946857   0.298671   3.170 0.00221 **
## St          -0.997463   0.234905  -4.246 6.16e-05 ***
## Fr.L:St     -0.007867   0.191944  -0.041 0.96742
## Fr.Q:St     -0.043393   0.244968  -0.177 0.85988
## Re.L:Fr.L    3.399373   0.275290  12.348 < 2e-16 ***
## Re.Q:Fr.L   -0.646249   0.251039  -2.574 0.01202 *
## Re.L:Fr.Q   -2.588037   0.514630  -5.029 3.27e-06 ***
## Re.Q:Fr.Q      NA      NA      NA      NA
## Re.L:St     -0.470624   0.240856  -1.954 0.05443 .
## Re.Q:St     -0.128587   0.204226  -0.630 0.53085
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8497 on 75 degrees of freedom
## Multiple R-squared:  0.9553, Adjusted R-squared:  0.9476
## F-statistic: 123.4 on 13 and 75 DF, p-value: < 2.2e-16

```

```

## Start: AIC=76.57
## log(R_moment_3) ~ log(St) + Re + Fr + St * Fr + Fr * Re + St *
## Re

##
## Call:
## lm(formula = log(R_moment_3) ~ log(St) + Re + Fr + St * Fr +
## Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.7949 -0.4431 -0.1224  0.5575  2.9257
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.54701     0.51756   8.785 3.81e-13 ***
## log(St)       2.36189     0.27259   8.665 6.47e-13 ***
## Re.L         -4.89792     0.52317  -9.362 3.06e-14 ***
## Re.Q          1.52364     0.42038   3.624 0.000525 ***
## Fr.L         -3.75114     0.37051 -10.124 1.12e-15 ***
## Fr.Q          1.83513     0.50306   3.648 0.000486 ***
## St           -1.72871     0.39565  -4.369 3.95e-05 ***
## Fr.L:St       -0.09194     0.32329  -0.284 0.776889
## Fr.Q:St       -0.03580     0.41260  -0.087 0.931098
## Re.L:Fr.L      6.45803     0.46368  13.928 < 2e-16 ***
## Re.Q:Fr.L     -1.15825     0.42283  -2.739 0.007689 **
## Re.L:Fr.Q     -4.96308     0.86680  -5.726 2.01e-07 ***
## Re.Q:Fr.Q      NA          NA      NA      NA
## Re.L:St       -0.79750     0.40568  -1.966 0.053018 .
## Re.Q:St       -0.17827     0.34398  -0.518 0.605816
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.431 on 75 degrees of freedom
## Multiple R-squared:  0.9459, Adjusted R-squared:  0.9365
## F-statistic: 100.8 on 13 and 75 DF, p-value: < 2.2e-16

## Start: AIC=131.52
## log(R_moment_4) ~ log(St) + Re + Fr + St * Fr + Fr * Re + St *
## Re

##
## Call:
## lm(formula = log(R_moment_4) ~ log(St) + Re + Fr + St * Fr +
## Fr * Re + St * Re, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.6675 -0.6183 -0.1392  0.7410  3.8875
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  9.28287     0.70476  13.172 < 2e-16 ***
## log(St)       3.10948     0.37119   8.377 2.28e-12 ***
## Re.L         -5.63270     0.71240  -7.907 1.80e-11 ***

```

```

## Re.Q      1.94835    0.57242    3.404 0.001070 **
## Fr.L      -5.55550    0.50452 -11.012 < 2e-16 ***
## Fr.Q      2.71540    0.68500    3.964 0.000167 ***
## St       -2.36662    0.53876   -4.393 3.62e-05 ***
## Fr.L:St   -0.17635    0.44023   -0.401 0.689864
## Fr.Q:St   -0.01987    0.56184   -0.035 0.971880
## Re.L:Fr.L  9.47956    0.63138   15.014 < 2e-16 ***
## Re.Q:Fr.L -1.66019    0.57576   -2.883 0.005130 **
## Re.L:Fr.Q -7.28081    1.18031   -6.169 3.21e-08 ***
## Re.Q:Fr.Q      NA          NA          NA          NA
## Re.L:St     -1.08269    0.55241   -1.960 0.053716 .
## Re.Q:St     -0.22245    0.46840   -0.475 0.636226
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.949 on 75 degrees of freedom
## Multiple R-squared:  0.9457, Adjusted R-squared:  0.9363
## F-statistic: 100.6 on 13 and 75 DF,  p-value: < 2.2e-16

```