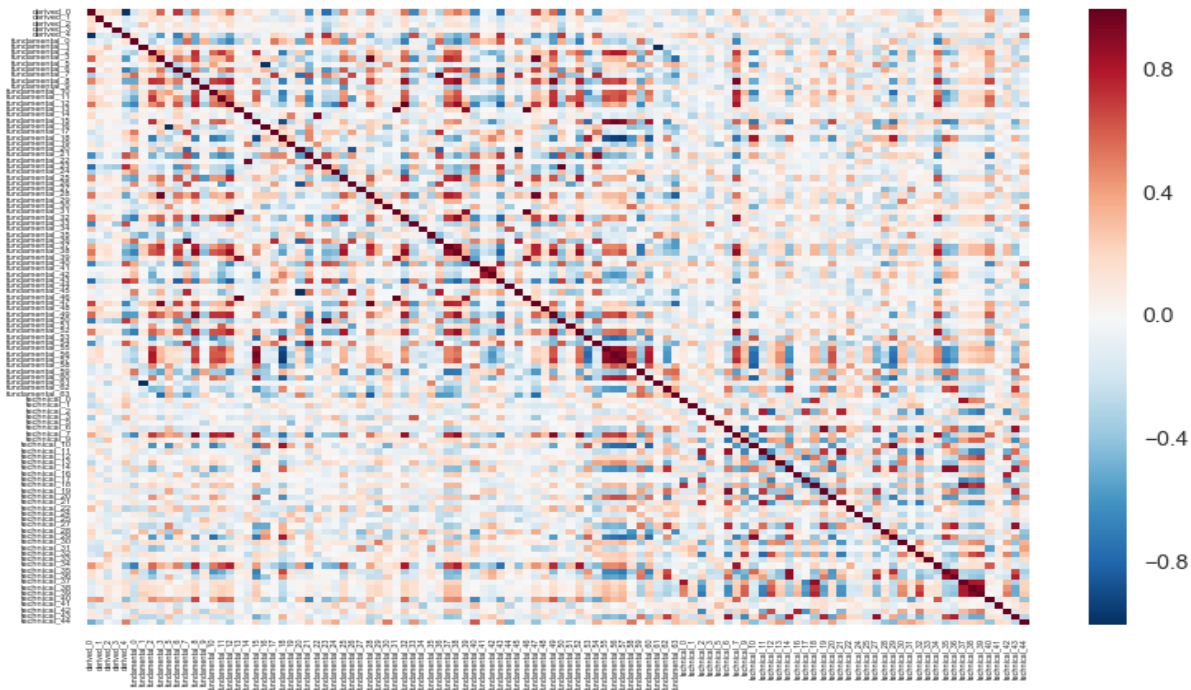


Results of Principal Component and Ridge Regression

Pair-wise correlation among features over the entire training set:



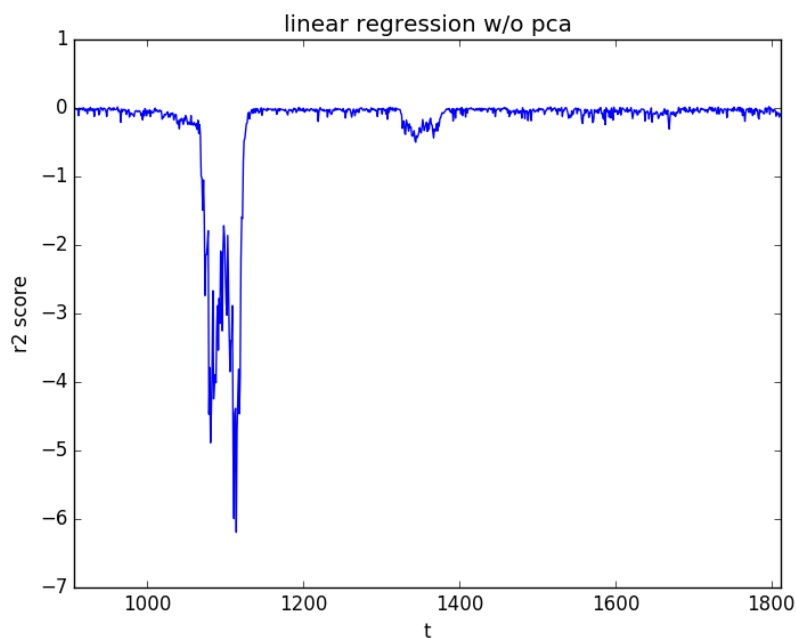
The abundance of collinearity among features implies pure linear regression may not work.

This problem is addressed by two methods: principal component and ridge regression.

To assess the performance of each, I divide the data sets into 60%:40% and use the first portion as training and the second as testing.

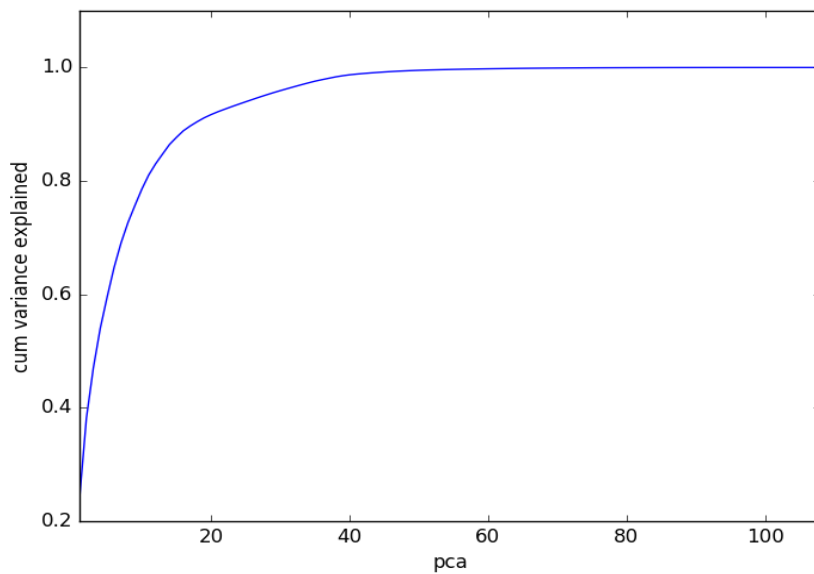
A time series of r^2 scores is generated from the testing set: the best score is 1 and the worst score can be arbitrarily negative.

Below is the time series of r^2 scores with pure linear regression as our baseline:

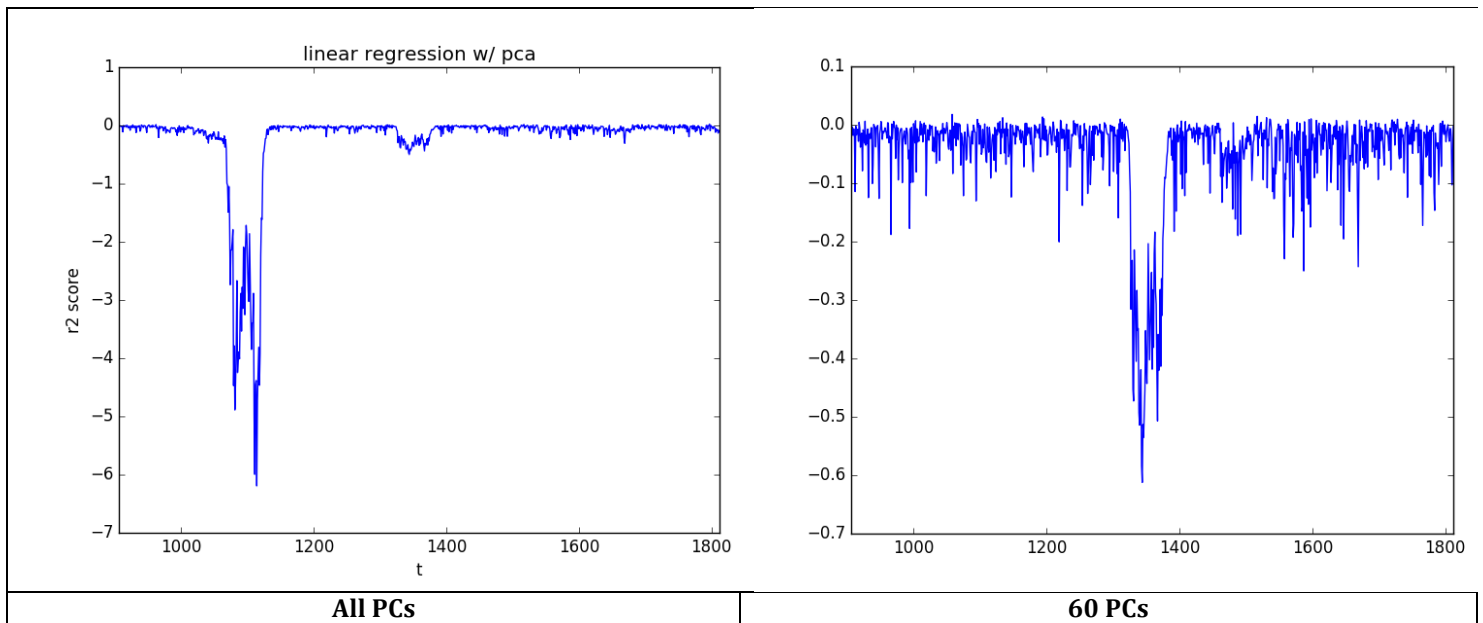


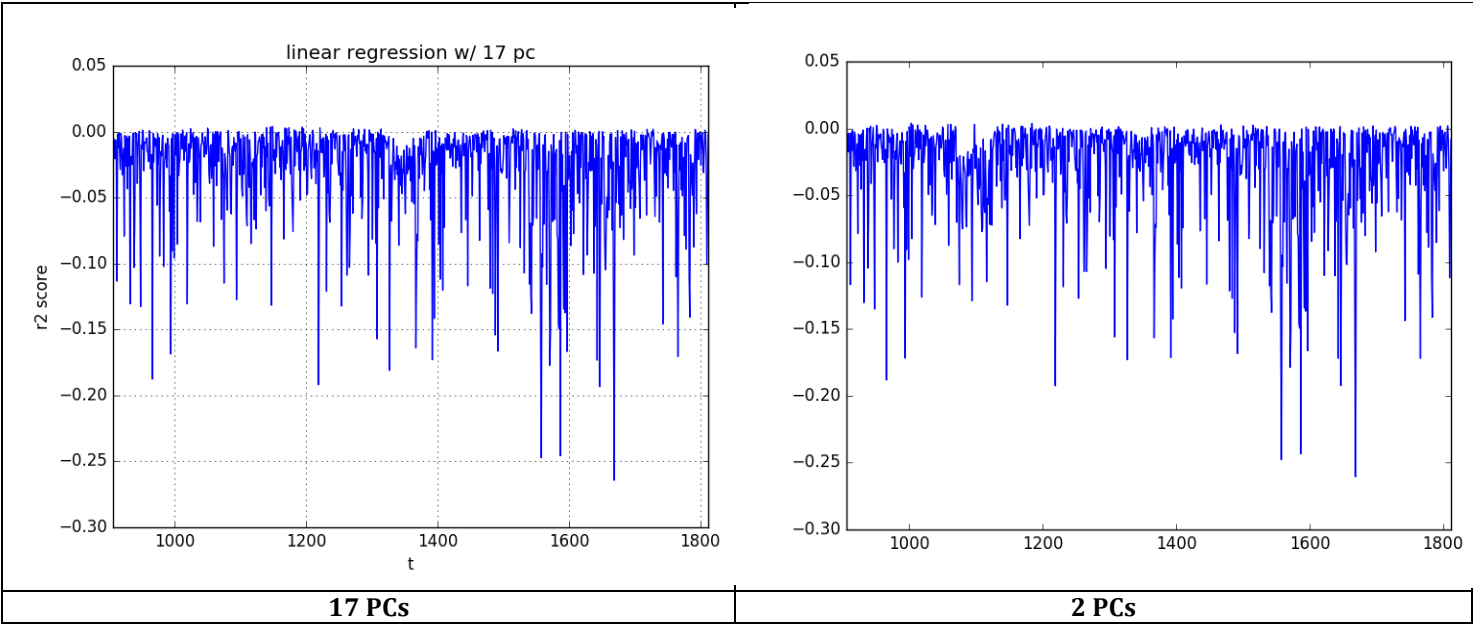
Principal Component Analysis

Variance explained with Principal Component: 17 PCs explains ~ 90% variance



Time series of r^2 scores of testing set using different # PCs: Overall scores have improved by reducing # PCs and stabilize at ~ 17 PCs.



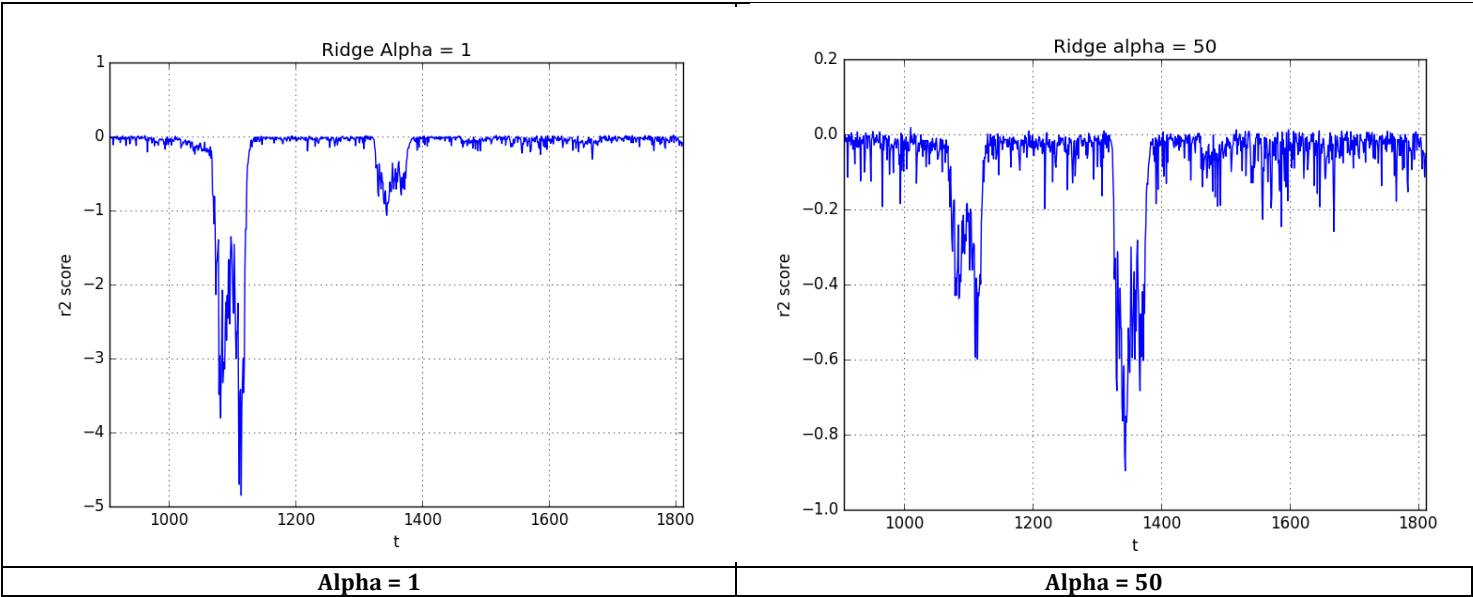


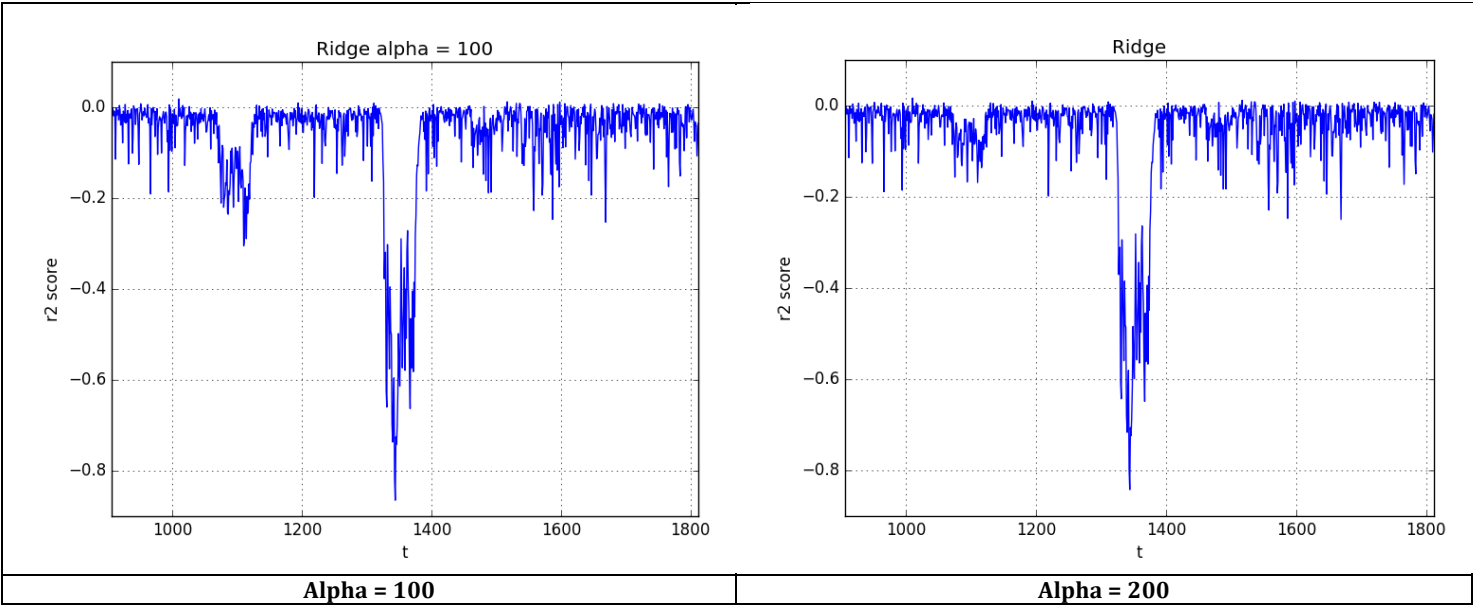
Statistics of r2 scores using 17 PCs:

count	907.000000
mean	-0.029043
std	0.038069
min	-0.264366
25%	-0.036402
50%	-0.015020
75%	-0.004523
max	0.003937

Ridge Regression

Time series of r2 scores of Ridge Regression using different alphas (penalty coefficient): scores have improved by increasing alpha. But alpha cannot be too big or else the model will become trivial.





Statistics of r2 scores with alpha = 200:

count	907.000000
mean	-0.061561
std	0.117495
min	-0.842012
25%	-0.059893
50%	-0.023498
75%	-0.009553
max	0.016547

Coefficients of features with alpha = 200:

